

# An automated assessment system for data structures and algorithms with diverse question types and knowledge-driven evaluation

Hoang Ngoc Long<sup>\*</sup>, Do Van Nhon  
HongBang International University

## ABSTRACT

*This paper proposes an automated assessment system for Data Structures and Algorithms (DSA) that supports eight diverse question types-Multiple-Choice Questions (MCQ), Fill-in-the-Blank Number (FN), Fill-in-the-Blank Number Sequence (FNS), Fill-in-the-Blank Programming Syntax (FPS), Fill-in-the-Blank Short String (FSS), Fill-in-the-Blank Expression/Formula (FE), Matching Pairs (MP), and Programming Exercises (PE)-across seven difficulty levels (Easy to Expert). Leveraging Bloom's Taxonomy and psychometric parameters (difficulty:  $p$ , discrimination:  $d$ ), the system uses knowledge matrices to integrate Course Learning Outcomes (CLOs), six content chapters, and cognitive levels, enabling adaptive test generation for quizzes, midterms, and exams. Automated scoring employs exact matching for MCQ/MP/FN/FNS, pattern matching for FPS/FSS, symbolic evaluation for FE, and test-case evaluation for PE. The system tracks knowledge progression across tests, per topic, and per CLO, providing comprehensive feedback. Initial implementations have shown the system's ability to improve assessment quality, decrease grading efforts, and facilitate data-driven teaching approaches in technical fields such as DSA.*

**Keywords:** *automated assessment, data structures and algorithms, e-learning, intelligent tutoring system, knowledge representation*

## 1. INTRODUCTION

Intelligent Tutoring Systems (ITS) leverage AI for personalized instruction and assessment, integrating domain knowledge with student activities. In DSA, ITS aid mastery of complex concepts like algorithmic design [1, 2], but comprehensive DSA assessment systems are limited. Essential to ITS, assessment modules evaluate proficiency and track progress, helping students identify strengths and weaknesses. While ITS excel in STEM [1, 2, 3], DSA-specific frameworks lack support for diverse cognitive and practical skills. Bloom's Taxonomy organizes assessments across cognitive levels, with studies [4] showing its value in balancing theory and practice. However, existing systems often use basic formats like Multiple-Choice Questions (MCQs) [4, 5] and manual grading, limiting scalability. Research in [5] proposes an intelligent system for course content access and MCQ testing but lacks testing theories for question design. Commercial tools like the International English Test [6], VioEdu [7] support general education but not IT disciplines like DSA. Learning Management Systems (LMS) like Elearning [8] enable testing but lack deep performance analysis. Recent advancements have introduced automated test generation and scoring, yet many systems

remain constrained by rigid question formats and insufficient alignment with CLOs. The study [9] has offered new insights on how an automated assessment can contribute to both curricula and pedagogical practice. The curriculum in this study corresponds to the use of formative exercises in a C++ programming course and pedagogic representation shows meaningful integration of an e-assessment tool to support learners and develop their knowledge competencies. Automated platforms like Moodle [8], CodeRunner [10], and LeetCode [11] have advanced assessment practices. Moodle [8] supports MCQs and short-answer questions, CodeRunner [10] excels in programming exercises (PE), and LeetCode [11] focuses on coding challenges. PrairieLearn, an open-source platform for STEM education, supports custom question variants with randomization and mastery-based adaptivity, allowing student retries until proficiency, and auto-grades code exercises via test cases as well as text-based explanations [12]. Yet, it relies on custom implementations for formats like fill-in-the-blank or matching pairs, with manual Bloom alignment, no integrated CLO tracking, and limited flexible difficulty for DSA assessments. However, these systems lack support for diverse

---

Corresponding author: Hoang Ngoc Long  
Email: [longhn@hiu.vn](mailto:longhn@hiu.vn)

question types (e.g., Fill-in-the-Blank Number Sequences (FNS) or Programming Syntax (FPS)) and flexible difficulty levels, hindering comprehensive evaluation. Moreover, alignment with CLOs and tracking knowledge progression by topic or across

tests remain underdeveloped, resulting in no complete solution for assessing Information Technology subjects within ITS frameworks. To highlight the novelty, Table 1 compares our system with existing platforms on key aspects.

**Table 1.** Comparison with Existing Assessment Platforms

System	Question Types	Adaptivity	Auto PE	Bloom Alignment
Ours	8 types (MCQ-PE)	Matrix (7 levels)	Yes	Explicit
Moodle	MCQ, short answer	Limited	Partial	Manual
CodeRunner	Code exercises	Limited	Yes	Manual
PrairieLearn	Custom, MCQ, code	Mastery-based	Yes	Manual
LeetCode	Coding challenges	Difficulty rec.	Yes	Manual

Assessment systems should integrate diverse question types, support multiple cognitive levels, and offer automated scoring to reduce instructor workload. This study addresses these gaps by proposing an assessment framework with seven difficulty levels, eight question types, and knowledge matrices integrating CLOs, six content chapters, Bloom's Taxonomy, and psychometric parameters (difficulty:  $p$ , discrimination:  $d$ ). The framework supports automated test generation, scoring, and evaluation per test, across tests, by topic, and per CLO, offering a scalable solution for DSA education. The remainder of this paper is structured as follows: Section 2 presents the knowledge representation model for assessment; Section 3 outlines the core assessment problems and details the algorithms for matrix generation, test creation, scoring, and progression tracking; Section 4 covers implementation, architecture, and evaluation; and Section 5 concludes with future directions.

**2. KNOWLEDGE REPRESENTATION MODEL FOR ASSESSMENT**

**2.1. Knowledge base model for data structures and algorithms**

**A. Modeling the knowledge base for the content of DSA course**

**Definition 2.1:** The knowledge organization model for representing the course content of DSA is defined as:  $(T, Q, R, CLO)$ . In which:

- $T$  (*Topics*) is the set of topics in the course, organized in a hierarchical tree structure based on content relationships. It includes six chapters: (C1) Overview, (C2) Searching & Sorting, (C3) Linked Lists, (C4) Stacks & Queues, (C5) Binary Search Trees, and (C6) Hash Tables).
- $Q$  (*Questions*) is the set of diverse questions in the course, expanded from objective multiple-choice questions to eight main types: Multiple-Choice Questions (MCQ), Fill-in-the-Blank Number (FN),

Fill-in-the-Blank Number Sequence (FNS), Fill-in-the-Blank Programming Syntax (FPS), Fill-in-the-Blank Short String (FSS), Fill-in-the-Blank Expression/Formula (FE), Matching Pairs (MP), and Programming Exercises (PE). For each question  $q \in Q$ , the structure is defined in detail as:  $(question\_type, question\_content, answers, correct\_answers, difficult\_level, discrimination\_level, bloom\_level, tolerance\_range, partial\_credit, synonyms, syntax\_variants, format\_variants, test\_cases)$ , where,  $question\_type$  is the type of question, determining its format and scoring method, including MCQ, FN, FNS, FPS, FSS, FE, MP, and PE;  $question\_content$  is the content of the question;  $answers$  is the set of possible answers;  $correct\_answers$  is the set of correct answers;  $difficult\_level$  is the difficulty level of the question;  $discrimination\_level$  is the discrimination level, reflecting the question's ability to differentiate students by ability; initial values for  $difficult\_level$  ( $p$ ) and  $discrimination\_level$  ( $d$ ) are expert-estimated by IT lecturers; these parameters are refined post-test using student data:  $p = (\text{number of correct responses} / \text{total attempts})$ ,  $d = (\text{proportion correct in upper 27\%} - \text{proportion correct in lower 27\%})$  per classical test theory [4]; questions with low  $d$  ( $< 0.2$ ) or extreme  $p$  ( $\leq 0.1$  or  $\geq 0.8$ ) are flagged for refinement or removal;  $bloom\_level$  represents the cognitive level of the question according to the Bloom's Taxonomy scale of the course's learning outcomes;  $tolerance\_range$  represents acceptable error range for numerical (FN) or formula-based (FE) questions to handle minor calculation errors;  $partial\_credit$  is the mechanism for awarding partial points if the answer is partially correct;  $syntax\_variants$  is the set of valid syntax variations for FPS questions;  $synonyms$  is the set of synonymous terms or phrases for FSS

questions; *format\_variants* is the set of valid formats for FE questions; *test\_cases* is the set of test cases for PE questions. This extended structure supports flexible automated scoring, with mechanisms tailored to each question type: exact matching for MCQ, FN, FNS, and MP; pattern matching for FPS, FSS; symbolic evaluation for FE; and test-case evaluation for PE, enhanced by attributes like *tolerance\_range*, *partial\_credit*.

- *R (Relationships)* is the set of relationships between the course topics (*T*) and questions (*Q*), represented through knowledge matrices. These relationships map each question to one or more chapters, CLOs. Formally,  $R = \{(question\_id, topic, clo)\}$ . In which, *question\_id* is the identifier of a question from *Q*; *topic* is the topic from *T* associated with the question; *clo* is the course learning outcome from associated with the question. A *question\_id* can have multiple tuples in *R*, each with a unique (*topic, clo*).
- *CLOs* includes five objectives: understand algorithm concepts, complexity, and representation (*CLO1, Chapter 1*); analyze search problems, constraints, and solutions (*CLO2, Chapter 2*); analyze sorting problems, constraints, and solutions (*CLO3, Chapter 2*); master basic data structures for programming (*CLO4, Chapters 3-6*); implement DSA for simple problems (*CLO5, Chapters 3-6*).

This knowledge base model has been applied to organize a bank of questions for the DSA course. This question bank is used to generate a test suitable for the requirements and abilities of students, and the results of this test serve as the basis for diagnosing students' knowledge levels.

**B. Modeling user of the system**

**Definition 2.2:** The user model in the system is designed to represent and manage individual student profiles and their interactions with the assessment system, enabling personalized test generation, scoring, and progress tracking. The model is defined as: (*PROFILE, TEST\_LIST*); where:

- *PROFILE* contains the personal and academic information of each student.
- *TEST\_LIST* is the set of tests that a student has attempted, with each *test*  $\in$  *TEST\_LIST* structured

as: (*test\_id, question\_scores, total\_score*). In which, *test\_id*: a unique identifier linking the test to the test definition (*L, question\_sets, maximum\_scores, time, score*) from the assessment framework, where *L* represents the difficulty level, *question\_sets* categorizes questions by Bloom's levels (*Remembering, Understanding, Applying, Analyzing*), and *maximum\_scores* defines the maximum points for each question in the specific test; *question\_scores*: a set of pairs (*question\_id, achieved\_score*), where *question\_id* is the identifier of a question from the question bank *Q*, associated with the extended structure; *achieved\_score* is the score earned by the student for the corresponding question, calculated based on the *maximum\_score* defined in the test and the student's response; *total\_score* represents cumulative score of the test, computed as the sum of all *achieved\_scores* in *question\_scores*.

**2.2. Question bank model and knowledge linkages**

The question bank model (*Q*), a core component of the DSA assessment system, serves as a dynamic repository for tailored test creation and automated evaluation within the knowledge organization model (*T, Q, R, CLO*). It includes eight question types-MCQ, FN, FNS, FPS, FSS, FE, MP, PE-spanning six chapters (C1-C6). Each question is structured with attributes as defined in Definition 2.1 to support scoring mechanisms (e.g., test-case evaluation for PE). The design of the question bank should support operations like removing or revising low-quality questions and incorporating high-quality ones to ensure continuous improvement in both the number and quality of questions [4]. Knowledge linkages via the relationship set *R* connect *Q* to topics *T* and *CLOs* through knowledge matrices, incorporating *bloom\_level*, *difficult\_level*, and *question\_type*. These matrices, generated per test difficulty, guide question selection. Table 2 shows a Level 3 (Medium) mid-term matrix, distributing 30 questions across Bloom's levels and types. This two-dimensional matrix has question types as columns and Bloom's Taxonomy levels as rows, with each cell containing an array of tuples. Each tuple pairs a chapter (topic) with the number of questions allocated to that chapter.

**Table 2.** Knowledge Matrix for a Mid-Term Exam at Level 3 - Medium

Bloom Levels	Question Types								Total
	MCQ	FN	FNS	FPS	FSS	FE	MP	PE	
R	-	(C2, 3), (C3, 3), (C4, 3)	-	-	-	-	-	-	9
U	-	(C2, 4)	(C3, 4)	(C4, 4)	-	-	-	-	12
AP	-	(C2, 2)	(C3, 2)	(C4, 2)	-	-	-	-	6

Bloom Levels	Question Types							Total	
	MCQ	FN	FNS	FPS	FSS	FE	MP		PE
AN	-	-	-	(C3, 1)	-	-	-	(C2, 1), (C4, 1)	3
Total	-	15	6	7	-	-	-	2	30

A test is typically created by choosing questions from a question bank, which is a comprehensive collection of numerous questions. It has structure and defined in Definition 2.3.

**Definition 2.3:** The test structure in the system defines the organization and composition of assessments used to evaluate student performance in the DSA course. It is designed to ensure a balanced evaluation across topics, learning outcomes, and cognitive levels, while supporting automated scoring and personalized feedback. The test structure is formally defined as:

- $(L, question\_sets, maximum\_scores, time, score)$ ; where:
- $L$  (Test Difficulty Level) represents the overall difficulty of the test, categorized into seven levels ranging from Easy to Expert. The difficulty level  $L$  guides the selection of questions and influences the assignment of  $maximum\_scores$ , ensuring alignment with the target student population's ability.
  - $question\_sets$  is the collection of questions including  $(Q_R, Q_U, Q_{AP}, Q_{AN}, ordered\_questions)$ . In which,  $Q_R, Q_U, Q_{AP}, Q_{AN}$  are the set of questions at the Remembering (R), Understanding (U), Applying (A), Analyzing (A) Bloom's levels correspondingly;  $ordered\_questions$  is

the collection of questions in particular order of  $[Non\_PE\_Questions, PE\_Questions]$ .

- $maximum\_scores$  is a set of pairs  $(question\_id, maximum\_score)$ , where  $question\_id$  identifies a question in  $question\_sets$ , and  $maximum\_score$  is the maximum points assigned to that question in the specific test. The  $maximum\_score$  can be assigned by the user or is determined based on the question's complexity, Bloom's level, and the test's difficulty level ( $L$ ), allowing for flexibility across different tests.
- $time$  is the duration of the test, measured in minutes;  $score$  is the desired total score for the test in points. The parameters  $time$  and  $score$  will be given by user.

The test structure integrates with the knowledge organization model  $(T, Q, R, CLO)$  by using the knowledge matrix to select questions from  $Q$  that align with the topics ( $T$ ) and CLOs targeted by the test. The distribution of questions across  $Q_R, Q_U, Q_{AP}$  and  $Q_{AN}$  is determined by the difficulty test level ( $L$ ) as in Table 3. For instance, a Level 3 (Medium) test might have 30% R, 40% U, 20% AP, and 10% AN, as seen in prior examples.

**Table 3.** Difficult level of the test

Difficulty of Test Level (L)	Bloom Proportions (R/U/AP/AN)	Difficult level (p)	Discrimination level (d)	Question Types	Course Content
Level 1 (Easy)	50 / 30 / 10 / 10	$p \geq 0.7$	$0.2 \leq d \leq 1$	MCQ, FSS, MP	C1, C2
Level 2 (Moderate)	40 / 30 / 20 / 10	$0.6 \leq p < 0.7$	$0.2 \leq d \leq 1$	MCQ, FN, MP	C1, C2
Level 3 (Medium)	30 / 40 / 20 / 10	$0.4 \leq p < 0.6$	$0.2 \leq d \leq 1$	FN, FPS, FNS, PE	C2–C4
Level 4 (Difficult)	20 / 30 / 30 / 20	$0.3 \leq p < 0.4$	$0.2 \leq d \leq 1$	FN, FPS, PE, FE	C2–C5
Level 5 (Very Diff.)	10 / 20 / 30 / 40	$p < 0.3$	$0.2 \leq d \leq 1$	FN, PE, MP, FE	C2–C5
Level 6 (Advance)	5 / 15 / 30 / 50	$p < 0.25$	$0.3 \leq d \leq 1$	FN, PE, FPS, MP	C2–C5
Level 7 (Expert)	0 / 10 / 30 / 60	$p < 0.2$	$0.4 \leq d \leq 1$	FN, PE, FPS, FE	C1–C6

When a student completes a test, their performance is recorded in the  $TEST\_LIST$  of the user model as  $(test\_id, question\_scores, total\_score)$ , where  $question\_scores$  contains pairs  $(question\_id, achieved\_score)$ . The  $achieved\_score$  is calculated using the question's scoring mechanism and adjusted by attributes like  $tolerance\_range$  or  $partial\_credit$ . This structure supports the system's goals of generating adaptive tests, providing detailed feedback, and tracking student progress across CLOs, with the flexibility to adjust difficulty and scoring based on educational needs.

**3. ASSESSMENT PROBLEMS AND ALGORITHM DESIGN**

This section outlines five core challenges in

developing an automated assessment system for DSA, using the knowledge model  $(T, Q, R, CLO)$ , question bank  $Q$ , and test structure  $(L, question\_sets, maximum\_scores, time, score)$  to enable adaptive testing and evaluation:

- 1) Knowledge Matrix Generation: building a matrix balancing questions across topics ( $T$ ), CLOs, Bloom's levels (R, U, AP, AN), and types (MCQ, FN, FNS, FPS, FSS, FE, MP, PE) per test's difficulty level ( $L$ ).
- 2) Test Generation: creating tests meeting difficulty ( $L$ ), chapter coverage, and Bloom's distribution, aligned with CLOs and psychometric parameters.
- 3) Evaluation and Scoring: designing scoring for eight

question types, considering attributes (e.g., tolerance\_range for FN, test\_cases for PE), for accurate assessment.

4) Knowledge Assessment per Test: analyzing performance on single tests to assess topic and CLO mastery, providing immediate feedback.

5) Knowledge Progression Across Tests: tracking progress over multiple tests to identify trends and personalize learning.

**3.1. Knowledge matrix generation problem and algorithm**

**Definition 3.1:** The knowledge matrix generation problem creates a balanced question distribution across Bloom's levels and question types, modeled

as:  $(L, Q, T, N) \rightarrow KM$ , where  $L$  is the test difficulty level (Table 3),  $Q$  is the question bank,  $T$  is the set of topics,  $N$  is the number of questions, and  $KM$  is a 2D knowledge matrix.

Algorithm 1 generates the knowledge matrix based on the proportions defined in Table 3. When the number of available questions is insufficient or cell counts do not match, the system prompts the instructor via the User Interface to adjust the matrix. Instructors can modify question counts per (*bloom\_level, question\_type, chapter*) while ensuring the total matches  $N$ . This semi-automated approach preserves the required Bloom proportions while allowing pedagogical adjustments.

```

Algorithm 1: GENERATE_KNOWLEDGE_MATRIX
Input:  $L, Q, T, N$ 
Output:  $KM$ 
Procedure:
1. EXTRACT Bloom_Proportions, Covered_Question_Types, Covered_Chapters FROM  $L$ 
2. INITIALIZE  $KM$  AS empty matrix [Bloom_Levels][Covered_Question_Types]
   WHERE Bloom_Levels = {R, U, AP, AN}
3. COMPUTE question counts per Bloom level:
    $N_R \leftarrow \text{FLOOR}(N * \text{Bloom\_Proportions}[R])$ 
    $N_U \leftarrow \text{FLOOR}(N * \text{Bloom\_Proportions}[U])$ 
    $N_{AP} \leftarrow \text{FLOOR}(N * \text{Bloom\_Proportions}[AP])$ 
    $N_{AN} \leftarrow \text{FLOOR}(N * \text{Bloom\_Proportions}[AN])$ 
   Total_Allocated  $\leftarrow N_R + N_U + N_{AP} + N_{AN}$ 
   IF Total_Allocated  $\neq N$  THEN
     PROMPT user to adjust  $N$ 
     RECALCULATE  $N_R, N_U, N_{AP}, N_{AN}$ 
   ENDIF
4. FOR EACH Bloom_Level IN Bloom_Levels DO
    $N_{\text{Bloom}} \leftarrow \text{SELECT} \{N_R, N_U, N_{AP}, N_{AN}\}$  FOR Bloom_Level
   FOR EACH Question_Type IN Covered_Question_Types DO
     Cell_Total  $\leftarrow 0$ 
     WHILE Cell_Total  $< N_{\text{Bloom}}$  DO
       Random_Chapter  $\leftarrow \text{RANDOM}(\text{Covered\_Chapters})$ 
       Default_Count  $\leftarrow \text{RANDOM\_INTEGER}(0, N_{\text{Bloom}} - \text{Cell\_Total})$ 
       APPEND (Random_Chapter, Default_Count) TO
          $KM[\text{Bloom\_Level}][\text{Question\_Type}]$ 
       Cell_Total  $\leftarrow \text{Cell\_Total} + \text{Default\_Count}$ 
     ENDWHILE
     IF Cell_Total  $> 0$  THEN
       VALIDATE  $KM[\text{Bloom\_Level}][\text{Question\_Type}]$ 
     ENDIF
     PROMPT user to adjust  $KM[\text{Bloom\_Level}][\text{Question\_Type}]$ 
   ENDFOR
   Bloom_Total  $\leftarrow \text{SUM}(\text{count FOR } (\text{chapter}, \text{count}) \text{ IN } KM[\text{Bloom\_Level}])$ 
   IF Bloom_Total  $\neq N_{\text{Bloom}}$  THEN
     PROMPT user to adjust  $KM[\text{Bloom\_Level}]$ 
   ENDIF
ENDFOR
5. Total_Questions  $\leftarrow \text{SUM}(\text{count FOR row IN } KM \text{ FOR cell IN row FOR } (\text{chapter}, \text{count}) \text{ IN cell})$ 
   IF Total_Questions  $\neq N$  THEN
     PROMPT user to adjust  $KM$ 
   ENDIF
6. RETURN  $KM$ 
    
```

**Figure 1.** The algorithm of knowledge matrix generation for DSA assessment

**3.2. Test generation problem and algorithm**

**Definition 3.2:** The test generation problem is modeled as:  $(\text{knowledge\_matrix}, Q, \text{time}, \text{score}) \rightarrow (L, \text{question\_sets}, \text{maximum\_scores}, \text{time}, \text{score})$ , where the knowledge matrix is obtained from Algorithm 1, and the output follows the test structure defined in Definition 2.3.

Algorithm 2 generates a test by selecting questions

from the question bank according to the knowledge matrix. Questions are filtered based on Bloom's level, question type, chapter, difficulty  $p$ , and discrimination  $d$  to ensure balance across cognitive levels, topics, and psychometric parameters. In this algorithm, the *base\_score* values are justified by teaching experience in DSA courses, scaled by *difficulty\_factor* and *bloom\_factor*. These weighting

parameters are currently fixed but can be made configurable in future versions.

Algorithm 2: GENERATE_TEST
<b>Input:</b> <i>knowledge_matrix (KM), Q, time, score</i>
<b>Output:</b> <i>L, question_sets, maximum_scores, time, score</i>
<b>Procedure:</b>
<pre> 1. <b>EXTRACT FROM</b> KM: L ← KM.Difficulty_Level, Covered_Question_Types ← KM.Covered_Question_Types,    Covered_Chapters ← KM.Covered_Chapters, Bloom_Proportions ← KM.Bloom_Proportions 2. <b>INITIALIZE</b> question_sets <b>AS</b> {QR ← [], QU ← [], QAP ← [], QAN ← [],    ordered_questions ← []} 3. <b>FOR EACH</b> Bloom_Level IN {R, U, AP, AN} <b>DO</b>    Current_Set ← <b>SELECT</b> {QR, QU, QAP, QAN} <b>FOR</b> Bloom_Level    <b>FOR EACH</b> Question_Type IN Covered_Question_Types <b>DO</b>      <b>FOR EACH</b> (Chapter, Count) IN KM[Bloom_Level][Question_Type] <b>DO</b>        Q_Filtered ← {q ∈ Q   q.bloom_level = Bloom_Level <b>AND</b>          q.question_type = Question_Type <b>AND</b>          q.chapter = Chapter <b>AND</b>          q.difficulty_level <b>MATCHES</b> L <b>AND</b>          q.discrimination_level <b>MATCHES</b> L}        <b>IF</b> SIZE(Q_Filtered) &lt; Count <b>THEN</b>          <b>PROMPT</b> user to adjust Count          <b>UPDATE</b> Count        <b>ENDIF</b>        Selected_Questions ← <b>RANDOM_SUBSET</b>(Q_Filtered, Count)        <b>APPEND</b> Selected_Questions <b>TO</b> Current_Set      <b>ENDFOR</b>    <b>ENDFOR</b> 4. All_Questions ← <b>CONCATENATE</b> (QR, QU, QAP, QAN)    PE_Questions ← {q ∈ All_Questions   q.question_type = PE}    Non_PE_Questions ← <b>SHUFFLE</b>(All_Questions \ PE_Questions)    question_sets.ordered_questions ← [Non_PE_Questions, PE_Questions]    <b>REDISTRIBUTE</b> question_sets.ordered_questions <b>TO</b> {QR, QU, QAP, QAN} <b>BY</b> bloom_level 5. <b>INITIALIZE</b> maximum_scores <b>AS</b> empty set    Total_Default_Score ← 0    <b>FOR EACH</b> Question IN question_sets.ordered_questions <b>DO</b>      <b>IF</b> Question.question_type = MCQ <b>THEN</b>        base_score ← 1      <b>ELSE IF</b> Question.question_type IN {FN, FNS, FPS, FSS, FE, MP} <b>THEN</b>        base_score ← 3      <b>ELSE IF</b> Question.question_type = PE <b>THEN</b>        base_score ← 5      <b>ENDIF</b>      Difficulty_Factor ← <b>MAP</b>(L, {Level 1: 0.5, Level 2: 0.7, Level 3: 1.0,        Level 4: 1.2, Level 5: 1.5, Level 6: 1.8, Level 7: 2.0})      Bloom_Factor ← <b>MAP</b>(Question.bloom_level, {R: 1.0, U: 1.2, AP: 1.5, AN: 2.0})      Question.default_score ← base_score × Difficulty_Factor × Bloom_Factor      Total_Default_Score ← Total_Default_Score + Question.default_score    <b>ENDFOR</b>    <b>FOR EACH</b> Question IN question_sets.ordered_questions <b>DO</b>      Question.maximum_score ← <b>ROUND</b>((Question.default_score /        Total_Default_Score) × score)      <b>IF</b> Question.maximum_score &lt; 1 <b>THEN</b>        Question.maximum_score ← 1      <b>ENDIF</b>      <b>APPEND</b> (Question.question_id, Question.maximum_score) <b>TO</b> maximum_scores    <b>ENDFOR</b>    Final_Score ← <b>SUM</b>(q.maximum_score <b>FOR</b> q <b>IN</b> question_sets.ordered_questions)    <b>IF</b> Final_Score ≠ score <b>THEN</b>      <b>PROMPT</b> user to adjust maximum_scores    <b>ENDIF</b> 6. <b>RETURN</b> (L, question_sets, maximum_scores, time, score) </pre>

Figure 2. The algorithm of test generation for DSA assessment

### 3.3. Evaluation and scoring of question types

Four scoring mechanisms—exact matching, pattern recognition, symbolic evaluation, and test-case-based evaluation—are used to assess responses for different question types, aligned with Bloom's Taxonomy. These are detailed in Definition 3.3 and Algorithm 3.

- Exact matching is applied to MCQ, MP, FN, and FNS. This mechanism compares responses directly to

correct answers. MCQ and MP award partial credit based on correct selections or pairs (e.g., 2/3 correct yields 2/3 score). FNS checks sequence elements with partial credit within a tolerance (e.g.,  $\pm 0.001$ ). FN scores full points if within tolerance, else zero.

- Pattern recognition is applied to FPS and FSS. Student responses are normalized (e.g., removes spaces for FPS, lowercase for FSS) and matched

against predefined correct answers or valid variants. Full marks are awarded only if a match is found.

- Symbolic evaluation is used for FE questions. Responses are normalized and parsed into symbolic expressions. The system checks for mathematical equivalence with the correct answer using expression simplification. Equivalent expressions receive full marks.
- Test-case-based evaluation is used exclusively for PE. It combines static analysis (40%) via *cppcheck* and dynamic unit testing (60%) via *GoogleTest*. Compilation failure results in zero score. Successful

compilations are evaluated for code style issues and functional correctness against test cases, with partial credit awarded for passed tests.

**Definition 3.3:** The evaluation and scoring problem computes the achieved score of a student's response based on the question type and its corresponding scoring mechanism. It is modeled as:  $(question, student\_response, maximum\_score) \rightarrow (achieved\_score)$ , where the input *question* comes from the question bank (Definition 2.1) and the *maximum\_score* is defined in the test structure (Definition 2.3). The detailed process is presented in Algorithm 3.

```

Algorithm 3: EVALUATE_AND_SCORE
Input: question, student_response, maximum_score
Output: achieved_score
Procedure:
1. INITIALIZE achieved_score  $\leftarrow$  0
2. CASE question.question_type OF
  MCQ, MP:
    correct_count  $\leftarrow$  COUNT(student_response MATCHES question.correct_answers)
    achieved_score  $\leftarrow$  (correct_count / LENGTH(question.correct_answers))  $\times$ 
      maximum_score  $\times$  question.partial_credit
  FN:
    is_valid  $\leftarrow$  PARSE(student_response) WITHIN (question.correct_answers  $\pm$ 
      question.tolerance_range)
    achieved_score  $\leftarrow$  is_valid ? maximum_score : 0
  FNS:
    correct_items  $\leftarrow$  COUNT(student_response[i] MATCHES question.correct_answers[i]
      WITHIN question.tolerance_range FOR i IN
      RANGE(LENGTH(question.correct_answers)))
    achieved_score  $\leftarrow$  (correct_items / LENGTH(question.correct_answers))  $\times$ 
      maximum_score  $\times$  question.partial_credit
  FPS, FSS:
    normalized_response  $\leftarrow$  NORMALIZE(student_response, question.question_type)
    variants  $\leftarrow$  question.question_type = FPS ?
      question.syntax_variants : question.synonyms
    achieved_score  $\leftarrow$  (normalized_response MATCHES question.correct_answers
      OR ANY v IN variants) ? maximum_score : 0
  FE:
    normalized_response  $\leftarrow$  STRIP(student_response)
    IF IS_VALID_EXPRESSION(normalized_response) THEN
      response_expr  $\leftarrow$  PARSE_TO_EXPRESSION(normalized_response)
      correct_expr  $\leftarrow$  PARSE_TO_EXPRESSION(question.correct_answers)
      achieved_score  $\leftarrow$  (SIMPLIFY(response_expr - correct_expr) = 0 OR
        normalized_response IN question.syntax_variants) ? maximum_score : 0
    ELSE
      achieved_score  $\leftarrow$  0
    ENDIF
  PE:
    issues  $\leftarrow$  COUNT_ISSUES(cppcheck, student_response)
    static_score  $\leftarrow$  MAX(0, 1 - (issues / 10))  $\times$  0.4  $\times$  maximum_score
    passed_tests  $\leftarrow$  COUNT_PASSED_TESTS(GoogleTest, question.test_cases,
      student_response)
    unit_score  $\leftarrow$  (passed_tests / LENGTH(question.test_cases))  $\times$  0.6  $\times$ 
      maximum_score  $\times$  question.partial_credit
    achieved_score  $\leftarrow$  static_score + unit_score
3. RETURN achieved_score
    
```

Figure 3. The algorithm of evaluation and scoring for diverse question types

### 3.4. Knowledge assessment per test

**Definition 3.4:** The per-test knowledge assessment problem evaluates a student's performance on a single test and generates feedback. It is modeled as:  $(test, student\_profile, question\_scores, student\_responses, R) \rightarrow (performance\_metrics, feedback)$ , where the inputs follow the test

structure (Definition 2.3) and user model (Definition 2.2), and *R* maps questions to topics and CLOs (Definition 2.1).

Algorithm 4 aggregates scores, computes performance metrics across topics, CLOs, and Bloom's levels, and generates detailed feedback, including PE-specific insights.

Algorithm 4: ASSESS_PER_TEST
<b>Input:</b> <i>test, student_profile, question_scores, student_responses, R</i> <b>Output:</b> <i>performance_metrics, feedback</i>
<b>Procedure:</b> 1. <b>INITIALIZE</b> performance_metrics $\leftarrow$ {topic_scores $\leftarrow$ {}, clo_scores $\leftarrow$ {}, bloom_scores $\leftarrow$ {}, bloom_max_scores $\leftarrow$ {}, total_score $\leftarrow$ 0}, feedback $\leftarrow$ "" 2. <b>FOR EACH</b> (question_id, achieved_score) <b>IN</b> question_scores <b>DO</b> question $\leftarrow$ test.question_sets[question_id] max_score $\leftarrow$ test.maximum_scores[question_id] total_score $\leftarrow$ total_score + achieved_score response $\leftarrow$ student_responses[question_id] status $\leftarrow$ achieved_score = max_score ? "Correct" : "Incorrect" <b>APPEND TO</b> feedback: "Question {question_id}: {status}" ({achieved_score}/{max_score})\nCorrect Answer: {question.correct_answers}\n" <b>IF</b> status = "Incorrect" <b>THEN</b> <b>IF</b> question.question_type = PE <b>THEN</b> issues $\leftarrow$ <b>COUNT_ISSUES</b> (cppcheck, response) failed $\leftarrow$ <b>COUNT_FAILED_TESTS</b> (GoogleTest, question.test_cases, response) <b>APPEND TO</b> feedback: "Static Issues: {issues}\nFailed Tests: {failed}/{LENGTH(question.test_cases)}\n" <b>ELSE</b> <b>APPEND TO</b> feedback: "Your Answer: {response}\n" <b>ENDIF</b> <b>ENDIF</b> relations $\leftarrow$ {(topic, clo)   (q_id, topic, clo) <b>IN</b> R <b>WHERE</b> q_id = question_id} <b>FOR EACH</b> (topic, clo) <b>IN</b> relations <b>DO</b> normalized_score $\leftarrow$ (achieved_score / max_score) / LENGTH(relations) topic_scores[topic] $\leftarrow$ (topic_scores[topic, 0] + normalized_score) / (COUNT(question_scores[topic, []]) + 1) clo_scores[clo] $\leftarrow$ (clo_scores[clo, 0] + normalized_score) / (COUNT(question_scores[clo, []]) + 1) <b>ENDFOR</b> bloom_level $\leftarrow$ question.bloom_level bloom_scores[bloom_level] $\leftarrow$ bloom_scores[bloom_level, 0] + achieved_score bloom_max_scores[bloom_level] $\leftarrow$ bloom_max_scores[bloom_level, 0] + max_score <b>ENDFOR</b> <b>FOR EACH</b> bloom_level <b>IN</b> bloom_scores <b>DO</b> <b>IF</b> bloom_max_scores[bloom_level] > 0 <b>THEN</b> bloom_scores[bloom_level] $\leftarrow$ bloom_scores[bloom_level] / bloom_max_scores[bloom_level] <b>ENDIF</b> <b>ENDFOR</b> 3. <b>PREPEND TO</b> feedback: "Student: {student_profile.name}, Score: {performance_metrics.total_score}/{test.score}\n" <b>FOR EACH</b> topic <b>IN</b> topic_scores <b>DO</b> Pt $\leftarrow$ topic_scores[topic] strength $\leftarrow$ Pt > 0.75 ? "Strong" : Pt $\geq$ 0.5 ? "Moderate" : "Weak" advice $\leftarrow$ Pt $\geq$ 0.5 <b>AND</b> Pt $\leq$ 0.75 ? " review concepts" : Pt < 0.5 ? " focus on fundamentals" : "" <b>APPEND TO</b> feedback: "Topic {topic} ({Pt:.2f}): {strength} at {test.L} level{advice}\n" <b>ENDFOR</b> <b>FOR EACH</b> clo <b>IN</b> clo_scores <b>DO</b> Pclo $\leftarrow$ clo_scores[clo] status $\leftarrow$ Pclo $\geq$ 0.5 ? "Meets" : "Below" advice $\leftarrow$ Pclo < 0.5 ? " focus on related chapters" : "" <b>APPEND TO</b> feedback: "CLO{clo} ({Pclo:.2f}):{status} at {test.L} level{advice}\n" <b>ENDFOR</b> <b>IF ANY</b> (question.question_type = PE <b>IN</b> test.question_sets) <b>THEN</b> <b>APPEND TO</b> feedback: "Programming Recommendations:\n" <b>IF</b> issues > 0 <b>THEN APPEND TO</b> feedback: "- Review coding style\n" <b>IF</b> failed > 0 <b>THEN APPEND TO</b> feedback: "- Debug failed test cases\n" <b>ENDIF</b> 4. <b>RETURN</b> (performance_metrics, feedback)

Figure 4. The algorithm of Per-test knowledge assessment and feedback

### 3.5. Knowledge progression across tests

**Definition 3.5:** The knowledge progression across tests problem evaluates a student's performance trends over multiple tests and provides feedback on their learning trajectory. It is modeled as:  $(student\_profile, test\_list, R) \rightarrow (progression\_metrics, progression\_feedback)$ , where the inputs follow the

user model (Definition 2.2) and  $R$  maps questions to topics and CLOs (Definition 2.1).

Algorithm 5 aggregates results from Algorithm 4 across tests, computes average performance, classifies levels, and identifies trends using a slope-based method to generate actionable recommendations.

Algorithm 5: ASSESS_ACROSS_TESTS
<b>Input:</b> <i>student_profile, test_list, R</i>
<b>Output:</b> <i>progression_metrics, progression_feedback</i>
<b>Procedure:</b> 1. <b>INITIALIZE</b> <i>progression_metrics</i> $\leftarrow$ { <i>topic_progression</i> $\leftarrow$ {}, <i>clo_progression</i> $\leftarrow$ {}, <i>bloom_progression</i> $\leftarrow$ {}, <i>overall_progression</i> $\leftarrow$ []}, <i>progression_feedback</i> $\leftarrow$ "" 2. <b>FOR EACH</b> <i>test_entry</i> <b>IN SORT</b> ( <i>test_list</i> <b>BY</b> timestamp) <b>DO</b> <i>metrics</i> $\leftarrow$ ASSESS_PER_TEST( <i>test_entry.test</i> , <i>student_profile</i> , <i>test_entry.question_scores</i> , <i>test_entry.student_responses</i> , <i>R</i> ) <b>APPEND</b> ( <i>metrics.total_score</i> / <i>test_entry.test.score</i> ) <b>TO</b> <i>progression_metrics.overall_progression</i> <b>FOR EACH</b> ( <i>dim, score</i> ) <b>IN</b> <i>metrics.topic_scores</i> <b>DO</b> <b>APPEND</b> <i>score</i> <b>TO</b> <i>topic_progression</i> [ <i>dim</i> , []] <b>ENDFOR</b> <b>FOR EACH</b> ( <i>dim, score</i> ) <b>IN</b> <i>metrics.clo_scores</i> <b>DO</b> <b>APPEND</b> <i>score</i> <b>TO</b> <i>clo_progression</i> [ <i>dim</i> , []] <b>ENDFOR</b> <b>FOR EACH</b> ( <i>dim, score</i> ) <b>IN</b> <i>metrics.bloom_scores</i> <b>DO</b> <b>APPEND</b> <i>score</i> <b>TO</b> <i>bloom_progression</i> [ <i>dim</i> , []] <b>ENDFOR</b> <b>ENDFOR</b> 3. <b>FOR EACH</b> ( <i>key, data_name</i> ) <b>IN</b> [( <i>topic_progression</i> , "Topic"), ( <i>clo_progression</i> , "CLO"), ( <i>bloom_progression</i> , "Bloom")] <b>DO</b> <b>FOR EACH</b> ( <i>dim, s_list</i> ) <b>IN</b> <i>key</i> <b>DO</b> <i>mu</i> $\leftarrow$ MEAN( <i>s_list</i> ) <i>level</i> $\leftarrow$ <i>mu</i> < 0.3 ? "Failed" : <i>mu</i> < 0.5 ? "Weak" : <i>mu</i> < 0.7 ? "Average" : <i>mu</i> < 0.9 ? "Good" : "Excellent" <i>slope</i> $\leftarrow$ LINEAR_SLOPE( <i>s_list</i> ) <b>IF</b> LENGTH( <i>s_list</i> ) > 1 <b>ELSE</b> 0 <i>trend</i> $\leftarrow$ <i>slope</i> > 0 ? "Improving" : <i>slope</i> < 0 ? "Declining" : "Stable" <i>key</i> [ <i>dim</i> ] $\leftarrow$ ( <i>s_list</i> , <i>level</i> , <i>trend</i> ) <b>APPEND TO</b> <i>progression_feedback</i> : "{ <i>data_name</i> } { <i>dim</i> }: { <i>level</i> }" "({{s:.2f FOR s IN s_list}}, { <i>trend</i> })\n" <b>ENDFOR</b> <b>ENDFOR</b> 4. <i>mu</i> $\leftarrow$ MEAN( <i>progression_metrics.overall_progression</i> ) <i>level</i> $\leftarrow$ <i>mu</i> < 0.3 ? "Failed" : <i>mu</i> < 0.5 ? "Weak" : <i>mu</i> < 0.7 ? "Average" : <i>mu</i> < 0.9 ? "Good" : "Excellent" <i>slope</i> $\leftarrow$ LINEAR_SLOPE( <i>progression_metrics.overall_progression</i> ) <b>IF</b> LENGTH( <i>progression_metrics.overall_progression</i> ) > 1 <b>ELSE</b> 0 <i>trend</i> $\leftarrow$ <i>slope</i> > 1 ? "Improving" : <i>slope</i> < -1 ? "Declining" : "Stable" <i>progression_metrics.overall_progression</i> $\leftarrow$ ( <i>progression_metrics.overall_progression</i> , <i>level</i> , <i>trend</i> ) <b>PREPEND TO</b> <i>progression_feedback</i> : "Overall: { <i>level</i> } ({{s:.2f FOR s IN <i>progression_metrics.overall_progression</i> [0]}}, { <i>trend</i> })\n" 5. <b>IF ANY</b> ( <i>l</i> <b>IN</b> ["Failed", "Weak"] <b>FOR</b> _, (_, l, _) <b>IN</b> [ <i>topic_progression</i> , <i>clo_progression</i> ]) <b>THEN</b> <b>APPEND TO</b> <i>progression_feedback</i> : "Recommendation: Focus on topics/CLOs with Failed or Weak performance.\n" <b>IF</b> <i>progression_metrics.overall_progression</i> [1] <b>IN</b> ["Failed", "Weak"] <b>THEN</b> <b>APPEND TO</b> <i>progression_feedback</i> : "Recommendation: Review foundational concepts across all areas.\n" <b>ENDIF</b> <b>ENDIF</b> 6. <b>RETURN</b> ( <i>progression_metrics, progression_feedback</i> )

Figure 5. The algorithm of progression tracking across multiple tests

#### 4. IMPLEMENTATION AND EVALUATION

##### 4.1. System architecture and knowledge base

The architecture of the knowledge assessment system for DSA, as illustrated in Figure 6, supports two user types: students and knowledge engineers. Students engage in testing and progress tracking, while knowledge engineers manage the knowledge base via the User Interface (1). The system comprises a User Interface (1) and seven other interconnected modules: (2) Request Classification routes user requests; (3) Knowledge Management handles insert, update, and delete operations; (4) Knowledge Evaluation assesses general per-

formance; (5) PE Evaluation processes Programming Exercises using *cppcheck*, and unit test cases; (6) Test Generation creates adaptive tests via Algorithm 2; (7) Knowledge Base stores data in an SQL database; and (8) Explanation generates detailed feedback.

The Knowledge Base (7) comprises 1250 questions spanning six chapters, categorized into eight question types and four cognitive levels based on Bloom's Taxonomy. These questions are systematically mapped to specific topics and CLOs, as defined in Definition 2.1 and detailed in Table 4. Data flows from the User Interface (1) through Request

Classification (2) to specialized modules (3-6), with PE Evaluation (5) validating accuracy for coding tasks.

Explanation (8) integrates results, ensuring comprehensive DSA assessment and feedback.

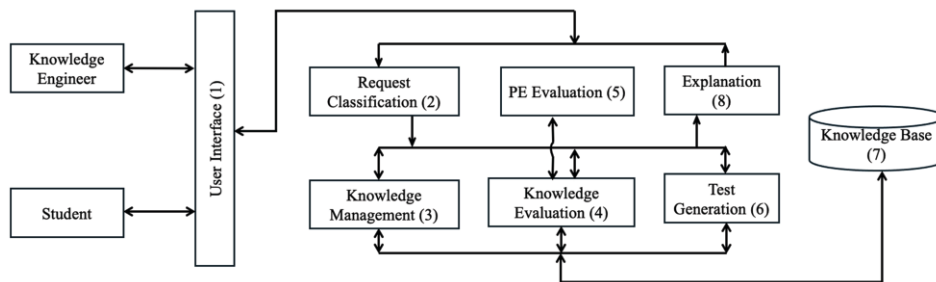


Figure 6. The architecture of knowledge assessment system

Table 4. Question Bank Distribution by Chapter and Bloom's Level

Chapter	Remembering (R)	Understanding (U)	Applying (AP)	Analyzing (AN)	Total
C1: Overview	60	55	45	45	205
C2: Searching & Sorting	80	75	65	60	280
C3: Linked Lists	60	55	45	40	200
C4: Stacks & Queues	55	45	45	45	190
C5: Tree	45	45	55	55	200
C6: Hash Table	45	45	45	40	175
<b>Total</b>	<b>345</b>	<b>320</b>	<b>300</b>	<b>285</b>	<b>1250</b>

4.2. Implementation and testing

The system was implemented using Spring Boot for backend logic and React for the frontend interface, with a MySQL database. It provides two main views: (1) a post-test feedback page showing individual question results, PE-specific analysis, and performance summaries aligned with CLOs and Bloom's levels; and (2) an overall progress dashboard displaying trend graphs, performance levels, and personalized recommendations. Pilot testing was conducted with 20 IT students at Hong Bang International University using the 1250-question knowledge base. Three test types were evaluated: quizzes (10 questions, Level 1 - 2), midterms (30 questions, Level 3 - 4), and finals (50 questions, Level 5). Algorithm 3 achieved 90% scoring accuracy. Progression tracking (Algorithms 4 - 5) revealed that only 10% of students showed positive learning trends, indicating a need for stronger teacher support and reevaluation of question difficulty. User feedback highlighted the interface clarity but suggested optimizing PE evaluation speed.

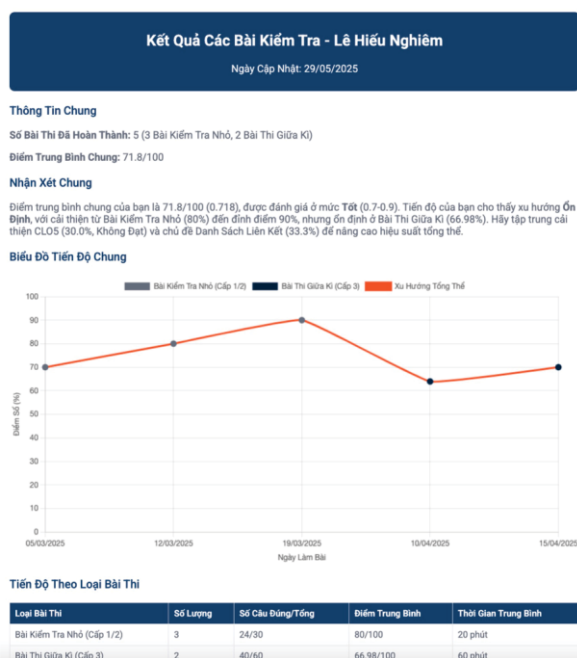
4.3. Discussion of results

Pilot testing with 20 IT students demonstrated the system's practical value. The architecture-with eight modules including PE Evaluation (5) using GoogleTest, cppcheck-supports diverse question types and difficulty levels, leveraging classical test theory for adaptive test generation aligned with CLOs, utilizing a 1250-question knowledge base (Table 3). Algorithm 3 achieved 90% scoring accuracy, and the system effectively tracks knowledge across topics and CLOs. The post-test

feedback page and progress dashboard (Figure 7–8) were well-received for providing clear, personalized feedback. It automatically assesses knowledge per test or across tests by topic and CLOs, using data mapping (e.g., question\_type, bloom\_level) for objective selection, with stored results enabling progress tracking and recommendations. However, only 10% of students showed positive learning trends, highlighting engagement challenges. The small sample size and single-institution context also limit the generalizability of the results. Future work should focus on optimizing PE evaluation performance and conducting larger-scale trials with statistical validation to strengthen the findings.



Figure 7. A part of test result user interface (summary of the test, topic and CLOs alignment, and recommendations)



**Figure 8.** Overall progress dashboard (summary of tests, trend graph of scores accross tests, performance levels, and recommendations)

**5. CONCLUSION**

This study presents an automated assessment system for the DSA course that supports eight

**REFERENCES**

[1]C.-C. Lin, A. Y. Q. Huang and O. H. T. Lu, "Artificial intelligence in intelligent tutoring systems toward sustainable education: a systematic review," *Smart Learning Environments*, vol. 10, no. 41, 2023.

[2]M. Kurni, M. S. Mohammed and S. K. G, A Beginner's Guide to Introduce Artificial Intelligence in Teaching and Learning, Springer Cham, 2023.

[3]E. Mousavinasab, N. Zarifasanaiey, S. R. N. Kalhori, M. Rakhshan, L. Keikha and M. G. Saeedi, "Intelligent tutoring systems: A systematic review of characteristics, applications, and evaluation methods," *Interactive Learning Environments*, vol. 29, no. 1, pp. 142-163, 2021.

[4]H. N. Long, M. T. Thanh and D. V. Nhon, "Designing a knowledge assessment system for the Data structures and Algorithms course," *Hong Bang International University Journal of Science*, vol. 6, pp. 65-74, 2024.

[5]H. D. Nguyen, "Intelligent System in Education: Requirements and Design Method". *Journal of Electronic Voltage and Application*, vol. 4, no. 2, pp. 12-19, 2023.

[6]International English Test, "International English

diverse question types and integrates knowledge matrices aligned with CLOs, Bloom's Taxonomy, and psychometric parameters. The system enables adaptive test generation, automated scoring through exact matching, pattern recognition, symbolic evaluation, and test-case evaluation, while providing detailed feedback and tracking student progress across topics and CLOs. Pilot deployments demonstrated its effectiveness in improving assessment quality, reducing grading workload, and supporting data-driven teaching. The architecture is scalable and suitable for larger classes or broader curricular integration. Future work will focus on two main directions: (1) developing intelligent question recommendation mechanisms using reinforcement learning or Hidden Markov Models, and (2) extending the system to other foundational computing courses with enhanced natural language processing capabilities. This research contributes to the development of intelligent, scalable assessment solutions in computer science education.

**ACKNOWLEDGEMENT**

This work is funded by Hong Bang International University under grant code GVTC18.02.

Test," [Online]. Available: <https://international-englishtest.com/>. [Accessed 23 May 2025].

[7]FPT, "VioEdu," FPT, [Online]. Available: <https://vio.edu.vn/>. [Accessed 23 May 2025].

[8]Hong Bang International University, "Elearning," Hong Bang International University, [Online]. Available: <https://elearning.hiu.vn/>. [Accessed 23 May 2025].

[9]A. L. C. Barczak, A. Mathrani, B. Han and N. H. Reyes, "Automated assessment system for programming courses: a case study for teaching data structures and algorithms," *Education Tech Research Dev*, no. 71, p. 2365-2388, 2023.

[10]R. Lobb and J. Harlow, "Coderunner: a tool for assessing computer programming skills," *ACM Inroads*, vol. 7, no. 1, pp. 47-51, 2016.

[11]LeetCode, "LeetCode," LeetCode, [Online]. Available: <https://leetcode.com/>. [Accessed 23 May 2025]. [12]M. West, N. Walters, M. Silva, T. Bretl and C. Zilles, "Integrating Diverse Learning Tools using the PrairieLearn Platform," in *Proceedings of SPLICE 2021 workshop CS Education Infrastructure for All III: From Ideas to Practice*, Virtual Event, 2021.

# Hệ thống đánh giá kiến thức tự động cho cấu trúc dữ liệu và giải thuật với nhiều loại câu hỏi

Hoàng Ngọc Long, Đỗ Văn Nhơn  
Trường Đại học Quốc tế Hồng Bàng

## TÓM TẮT

Bài báo đề xuất một hệ thống đánh giá tự động cho môn Cấu trúc Dữ liệu và Giải Thuật (DSA) hỗ trợ tám loại câu hỏi đa dạng-Hỏi Đáp Trắc Nghiệm (MCQ), Điền Số (FN), Điền Dãy Số (FNS), Điền Cú Pháp Lập Trình (FPS), Điền Chuỗi Ngắn (FSS), Điền Biểu Thức/Công Thức (FE), Ghép Đôi (MP), và Bài Tập Lập Trình (PE)-trên bảy cấp độ độ khó (Dễ đến Chuyên Gia). Sử dụng Phân loại Bloom và các tham số tâm lý trắc nghiệm (độ khó:  $p$ , khả năng phân biệt:  $d$ ), hệ thống áp dụng ma trận tri thức để tích hợp Kết Quả Học Tập Khóa Học (CLOs), sáu chương nội dung, và các cấp độ nhận thức, cho phép tạo bài kiểm tra thích ứng cho bài tập, kỳ thi giữa kỳ, và kỳ thi cuối kỳ. Đánh giá tự động sử dụng phương pháp khớp chính xác cho MCQ/MP/FN/FNS, khớp mẫu cho FPS/FSS, đánh giá biểu thức cho FE và đánh giá dựa trên bộ kiểm thử cho PE. Hệ thống theo dõi tiến trình tri thức qua các bài kiểm tra, theo chủ đề, và theo CLO, cung cấp phản hồi toàn diện. Các triển khai ban đầu đã cho thấy khả năng của hệ thống trong việc nâng cao chất lượng đánh giá, giảm bớt công sức chấm điểm và hỗ trợ các phương pháp giảng dạy trong các lĩnh vực kỹ thuật như Cấu trúc Dữ liệu và Giải Thuật.

**Từ khoá:** đánh giá tự động, cấu trúc dữ liệu và giải thuật, học trực tuyến, hệ thống hướng dẫn thông minh, biểu diễn tri thức.

---

Received: 29/5/2025

Revised: 09/10/2025

Accepted for publication: 16/10/2025