

## A Comparative Review of the CEFR and CET4 Writing Assessment with Insights from Task Complexity Theories

Changlin Li<sup>1</sup> , Nik Aloesnita Nik Mohd Alwi<sup>2\*</sup>   
Mohammad Musab Azmat Ali<sup>3</sup> 

<sup>1</sup>Centre for Modern Languages, Universiti Malaysia Pahang Al-Sultan Abdullah, 26000, Pekan Pahang, Malaysia;

Academic Affairs Office, Hebei Minzu Normal University, 067000, Chengde Hebei, China  
Email: PBA22007@student.umpsa.edu.my

<sup>2</sup>Centre for Modern Languages, Universiti Malaysia Pahang Al-Sultan Abdullah, 26000, Pekan Pahang, Malaysia

Email: aloesnita@umpsa.edu.my

<sup>3</sup>Centre for Modern Languages, Universiti Malaysia Pahang Al-Sultan Abdullah, 26000, Pekan Pahang, Malaysia

Email: mmusab@umpsa.edu.my

### ABSTRACT

#### CORRESPONDING

#### AUTHOR (\*):

Nik Aloesnita Nik Mohd Alwi  
(aloesnita@umpsa.edu.my)

#### KEYWORDS:

CEFR

CET4

Writing assessment

Comparison

#### CITATION:

Li, C., Nik Aloesnita Nik Mohd Alwi, & Mohammad Musab Azmat Ali. (2025). A Comparative Review of the CEFR and CET4 Writing Assessment with Insights from Task Complexity Theories. *Malaysian Journal of Social Sciences and Humanities (MJSSH)*, 10(3), e003251.  
<https://doi.org/10.47405/mjssh.v10i3.3251>

The CEFR level descriptors are applied globally for language assessment, which is already aligned with IELTS, TOEFL, etc. Meanwhile, the CET4 is essential in language learning and teaching proficiency assessment in China. Building on previous research, this study examines the relationship between the CEFR level descriptors and the CET4 writing rubrics, mainly focusing on the essay writing assessment within the past decade. Despite the broad utilisation of the CET4 in universities, its comparison with the CEFR level descriptors remains underexplored. Based on this situation, the study investigates task complexity theories, automated and manual scoring systems, and recent studies about essay writing. Findings indicate that the CET4 writing scores correspond roughly to CEFR levels A1–B2, though comparisons with higher proficiency levels (C1–C2) remain inconsistent. While automated scoring systems reliably evaluate basic linguistic dimensions, they struggle to assess more aspects, such as description, argument, task relevance, and clarity dimensions, under the CEFR and CET4 writing assessments. Furthermore, the automated scoring systems lack the capacity to capture the nuanced features of advanced writing. These findings underscore the necessity of human evaluation, particularly in essay writing content assessment, while highlighting opportunities to refine grading methodologies and task design to enhance essay writing instruction.

**Contribution/Originality:** The study aims to establish a basis for future research comparing CET4 writing rubrics and CEFR level descriptors in CET4 essay writing. It discusses the impact of task complexity on essay writing and highlights the necessity of manual evaluation to address the issue of automated essay scoring insufficient

assessment in the description, argument, task relevance, and clarity dimensions.

## 1. Introduction

Language assessment continues to pose challenges across different regions and testing frameworks. Among these, the Common European Framework of Reference for Languages (CEFR) (Council of Europe, 2020) is globally recognised for language proficiency assessment. The CEFR outlines diverse levels (preA1–C2), each with detailed descriptors. Meanwhile, the College English Test Band 4 (CET4) (NCECT, 2016) is essential for English learning and teaching proficiency assessment among universities in China. Nevertheless, including the writing section, connecting the CET4 with CEFR level descriptors remains incomplete (Zhang & Yang, 2023). Moreover, recent research underscores the need for more apparent correlations between the two systems (Wang et al., 2022). The CET4 evaluates listening, reading, translation, and writing; however, its rubrics are not yet officially aligned with the CEFR level descriptors (Liu & Chen, 2023).

Furthermore, task complexity has been found to impact learners' writing proficiency. Skehan's (1998, 2015, 2018) Limited Attentional Capacity Model (LACM) indicates that increased task complexity shifts performance across fluency, accuracy, and complexity. At the same time, the Cognition Hypothesis (CH) (Robinson, 2001, 2007, 2011b) similarly posits that heightened cognitive demands can elevate the quantity and quality of language output. Although both theories link task complexity to linguistic performance, the LACM highlights trade-offs among fluency, accuracy, and complexity, whereas the CH underscores cognitive demands that enhance advanced linguistic performance. Research also suggests that higher task complexity correlates with higher CEFR levels in tasks that require greater cognitive engagement and more sophisticated language (Zhang & Li, 2024). However, connecting the CET4 with CEFR assessments at the C1-C2 levels remains challenging; moreover, Automated Essay Scoring (AES) systems are insufficient to evaluate specific dimensions of essay writing (Thompson & Wu, 2023). Indeed, the AES systems are adequate for basic linguistic features at CEFR levels (A2-B1) but complicated with advanced features such as nuanced arguments (Li & Zhou, 2024). The limitation underscores the need for manual assessment methods to evaluate description and argument dimensions within the CEFR level descriptors while assessing task relevance and clarity dimensions under the CET4 rubrics.

### 1.1. Problem Statement

The lack of official comparison between CET4 writing rubrics and CEFR level descriptors, particularly regarding description, argument, task relevance, and clarity of ideas dimensions measured by AES systems, poses significant challenges for essay writing assessment proficiency. Furthermore, this misconnection is particularly evident at higher proficiency levels (C1-C2), where sophisticated language features and complex discourse patterns are crucial for the essay writing assessment (Anderson & Chen, 2023). Therefore, this research explores studies on task complexity, automated and manual assessment systems, and essay writing research, mainly in the past decade. It aims to benefit the further comparison of the CEFR and CET4 writing criteria.

### 1.2. Research Objectives

This study pursues three primary objectives:

- i. To investigate the impact of task complexity variations on linguistic performance in essay writing.
- ii. To analyse the connection between the CEFR level descriptors and the CET4 writing rubrics.
- iii. To analyse essay assessment research from the past decade, providing a foundation for further comparing the CEFR and CET4 writing criteria.

## 2. Literature Review

### 2.1. The CEFR and CET4 Serve as Authoritative Frameworks for Essay Writing Assessment

The CEFR ([Council of Europe, 2020](#)), first introduced by the Council of Europe in 2001 and updated in 2020, is a widely recognised language learning, teaching, and assessment framework. It categorises language proficiency into different levels, (pre)A1-A2 (Beginner), B1-B2 (Intermediate), and C1-C2 (Proficient), with detailed descriptors for each level, providing a standardised reference for language assessments ([North, 2014](#)). While the CEFR has been adopted globally, IELTS, TOEFL, and China's Standards of English Language Ability (CSE) aligned the results with the CEFR levels ([Papageorgiou et al., 2022](#)). However, the connection between the CET4 and CEFR level descriptors remains underexplored ([Zheng & Cheng, 2022](#)). See [Table 1](#).

Table 1: The CEFR Measurement Levels and Dimensions under Written Assessment  
([Council of Europe, 2020](#))

| Items                  | Classifications   |
|------------------------|---|
| Levels                 | C2, C1, B2, B1, A2, A1                                      |
| Measurement Dimensions | Overall, range, coherence, accuracy, description, argument. |

The CEFR level descriptors provide a comprehensive approach to assessing language proficiency across multiple skills, including the writing section, making it valuable for language assessments. For essay writing, the CEFR level descriptors provide specific descriptors for each level, ranging from basic contexts at A1 to the production of sophisticated academic contexts at C2 ([Little, 2020](#)). This broad spectrum enables more nuanced assessments of writing proficiency. While internationally recognised exams like IELTS and TOEFL align their results with CEFR levels, challenges persist in determining whether the CET4 rubrics are effectively mapped to the CEFR level descriptors. Notably, the connection of the CET4 with CEFR writing criteria at higher levels (C1-C2) remains under research and development ([Zheng & Cheng, 2022](#)). This misalignment may impact academic mobility and the global credibility of national assessments. Therefore, enhancing international comparability and deepening the CEFR and CET4 writing criteria analysis, particularly in essay writing assessment, is pivotal.

### 2.2. The CET4 as a Writing Assessment Criteria

Established in 1987, the CET4 is a widely administered English proficiency test in China that evaluates the language skills of millions of undergraduates each year ([Jin & Yang, 2018](#)). The CET4 scores range from 0 to 710 points, with 425 designated as the passing threshold ([Wang & Li, 2019](#)). However, it has not been officially mapped to CEFR levels, making it challenging to align its rubrics with those of international assessments ([Zhang & Liu, 2021](#)). This misalignment hinders the global comparability of the CET4, as its

connection to international frameworks like the CEFR remains undefined officially. Notably, given its significance, the writing component accounts for 15% of the total CET4 score. Students can also take the College English Test-Spoken English Test Band 4 (CET-SET4); the speaking test is conducted via a computer-based format, and candidates who have registered for the written exam are eligible to take the corresponding level of the speaking test. See [Table 2](#).

Table 2: The components of the CET4 written examination ([NCECT, 2016](#))

| Test Structure        | Content              | Test Type                    | Number of Items | Score Percentage | Testing Time |
|-----------------------|----------------------|------------------------------|-----------------|------------------|--------------|
| Writing               | Writing              | Short Writing                | 1               | 15%              | 30 minutes   |
| Listening             | Short News           | Single-choice                | 7               | 7%               | 25 minutes   |
| Comprehension         | Long Dialogues       | Questions                    | 8               | 8%               |              |
|                       | Listening Passages   | Single-choice                | 10              | 20%              |              |
|                       |                      | Questions                    |                 |                  |              |
| Reading Comprehension | Vocabulary           | Select-and-Fill-in-the-Blank | 10              | 5%               | 40 minutes   |
|                       | Long Passage Reading | Matching                     | 10              | 10%              |              |
|                       | Detailed Reading     | Single-choice Questions      | 10              | 20%              |              |
| Translation           | Chinese-to-English   | Paragraph Translation        | 1               | 15%              | 30 minutes   |
| Total                 | —                    | —                            | 57              | 100%             | 125 minutes  |

[Table 2](#) highlights that the CET4 places equal emphasis on listening and reading, each contributing 35% to the total score, while writing and translation account for 15%. Among these, the writing section (as [Table 3](#) shows) evaluates student's ability to construct coherent essays within a limited time, focusing on task relevance, clarity of ideas, coherence, and language accuracy. However, as pointed out by [Zhang and Liu \(2021\)](#), while the CET4 writing assessment evaluates these factors, it does not explicitly map the scores to the CEFR levels. This absence of a direct equivalence makes it challenging to compare the CET4 writing scores with internationally recognised frameworks, which could hinder students' opportunities for international academic mobility and professional certification. Given the impact of English assessment systems on university students' job opportunities, exploring how CET4 writing can better align with CEFR level descriptors is crucial. Strengthening this connection would enhance the CET4's comparability with global standards and improve the practical application of university English proficiency assessments such as the CET4 globally.

Table 3: The Specific Requirements of the CET4 Writing Test ([NCECT, 2016](#))

| Items            | Description   |
|------------------|---|
| Length           | No less than 120 words  |
| Time Limit       | 30 minutes  |
| Scoring Criteria | It consists of content, linguistic, and other factors, including task relevance, clarity of ideas, coherence, and language accuracy |

| Items     | Description  |
|-----------|--|
| Task Type | Mainly responses to given outlines, charts, visual prompts, etc., including the essays |

### 2.3. Task Complexity and Its Influence on Essay Writing

[Skehan's \(1998, 2015, 2018\)](#) and [Robinson's \(2001, 2007\)](#) theories provide essential frameworks for analysing task complexity and essay writing assessments in Second Language Acquisition (SLA) and Task-Based Language Teaching (TBLT). The theories offer insights into how cognitive demands influence linguistic production, which is essential when evaluating the connection between the CET4 writing rubrics and the CEFR level descriptors.

#### 2.3.1. Skehan's Limited Attentional Capacity Model

[Skehan \(1998, 2015, 2018\)](#) proposed the Limited Attentional Capacity Model (LACM), which posits that learners have finite cognitive resources when performing language tasks. According to Skehan's LACM ([Skehan, 2015, 2018](#)), as task complexity increases, learners are forced to prioritise certain aspects of language production, such as fluency, accuracy, and complexity. This allocation of cognitive resources often leads to trade-offs: tasks of higher complexity require a more significant cognitive load, resulting in enhanced accuracy or complexity but reduced fluency. Conversely, tasks of lower complexity impose less cognitive demand, facilitating greater fluency but limiting the sophistication of lexical and syntactic structures. Skehan's LACM ([Skehan, 2015, 2018](#)) has influential implications in the essay writing assessment. For instance, suppose the CET4 writing task demands a detailed argument with multiple supporting points. In that case, students may allocate more attention to accuracy and complexity, resulting in well-structured and grammatically correct essays, but may suffer in fluency due to time constraints. On the other hand, simpler prompts, such as personal reflections or straightforward descriptions, may lead to greater fluency but less lexical and syntactic complexity. These patterns are particularly relevant when considering the cognitive load placed on students in a timed writing task. The following table summarises key claims of low and high-complexity tasks to explore the relationship between task complexity and the trade-offs in essay writing. See [Table 4](#).

Table 4: Comparison of Low and High Complexity Tasks and Their Influence under LACM

| Task Complexity | Cognitive Demand | Fluency | Lexical and Syntactic Complexity              |
|-----------------|------------------|---------|---|
| Low Complexity  | Low              | High    | Simple structures, limited vocabulary         |
| High Complexity | High             | Low     | Increased grammatical accuracy and complexity |

[Table 4](#) highlights that cognitive demand is reduced in low-complexity tasks, which allows students to produce writing with greater fluency but less linguistic sophistication. In contrast, high-complexity tasks require more significant cognitive resources, often leading to trade-offs where students may make more complex structures and demonstrate higher accuracy but at the expense of fluency. The study ([Zhang & Li, 2022](#)) supports this model, as it found that tasks requiring abstract reasoning and extended

argument led to higher syntactic complexity and increased pauses and hesitations in students' essays. These findings suggest that task complexity in the CET4 writing assessment connects with Skehan's LACM (Skehan, 1998, 2015, 2018), demonstrating that increased cognitive load influences the balance between fluency, accuracy, and complexity in student writing. This understanding is critical when considering refining writing assessments in line with cognitive and linguistic demands.

### 2.3.2. Robinson's Cognition Hypothesis

The Cognition Hypothesis (CH) by Robinson (2001, 2007, 2011b) asserts that more cognitively demanding tasks push learners to use more complex linguistic structures. He identifies task complexity as the resource-dispersing dimensions and the resource-directing dimensions. The former refers to factors that affect cognitive load without directly influencing linguistic complexity, such as the availability of planning time, familiarity with the topic, and working memory demands. The latter refers to factors that shape linguistic output by requiring more sophisticated language use, such as reasoning demands and the number of elements involved in a task. According to Robinson's CH (Robinson, 2001, 2007, 2011a), increasing resource-directing dimensions in essay prompts, requiring students to synthesise information from multiple sources or engage in critical reasoning, should lead to more lexically and syntactically complex writing. This type of task complexity is thought to connect with higher CEFR levels, such as B2 to C1, which demand increased sophistication in language use. The study by (Wang et al., 2021) supports this idea by comparing two CET4 writing prompts: one that required simple narrative descriptions and another that asked students to analyse an abstract concept, such as "The impact of digital technology on education." The research found that essays responding to the more complex task demonstrated greater lexical variety and syntactic sophistication, supporting Robinson's CH (Robinson, 2001, 2007, 2011b) that increasing cognitive demands enhances linguistic output. However, one challenge highlighted by Robinson's CH (Robinson, 2001, 2007, 2011b) is that higher complexity tasks may disadvantage lower-proficiency learners. These students may struggle with the advanced cognitive demands associated with such tasks, potentially affecting their ability to perform well on more complex prompts. The issue should be considered when designing CET4 writing tasks to ensure fair proficiency assessment across different levels. The following table summarises how complexity in essay writing tasks can impact linguistic performance. See Table 5.

Table 5: Impact of Task Complexity on Linguistic Output

| Task Complexity | Cognitive Demand | Linguistic Output  | Potential Challenge  |
|-----------------|------------------|--|--|
| Low Complexity  | Low              | Simple language use, lower lexical and syntactic variety | May not adequately assess higher proficiency levels                  |
| High Complexity | High             | Increased lexical variety and syntactic sophistication   | Disadvantage for lower-proficiency learners due to cognitive demands |

Table 5 illustrates that low-complexity tasks typically result in more straightforward language use, with less variety in lexical and syntactic complexities. In contrast, high-complexity tasks demand more sophisticated linguistic features. However, as the cognitive demand increases, lower-proficiency students may struggle to meet these demands, potentially skewing the assessment results. This is a crucial consideration in

task design, ensuring that writing prompts are accessible to a broad range of proficiency levels while effectively assessing advanced language skills. Robinson's CH (Robinson, 2001, 2007) emphasises the importance of balancing task complexity to create assessments that fairly evaluate all learners.

## 2.4. Application of SLA Theories for Writing Assessment

Second Language Acquisition (SLA) theories offer valuable insights into how various task conditions affect learners' ability to plan, organise, and produce written text. These theoretical frameworks help explain why specific task designs may more effectively connect with CEFR descriptors and how the CET4 assessment methods can be refined to better reflect language proficiency. Understanding the role of task complexity in language performance is essential for improving the reliability and fairness of writing assessments and ensuring that they accurately measure learners' proficiency.

### 2.4.1. Skehan's Limited Attentional Capacity Model and Writing Proficiency

Skehan's Limited Attentional Capacity Model (LACM) (Skehan, 1998, 2009a, 2015) offers valuable insights into how writers manage different aspects of language performance under diverse task conditions. This model suggests that cognitive resources are finite. As task complexity increases, learners are forced to prioritise certain aspects of language, such as fluency, accuracy, or complexity, at the expense of others. This concept is particularly relevant for understanding why specific CET4 writing tasks may not fully capture students' accurate proficiency levels, especially when assessing higher-level skills outlined in the CEFR. For example, CET4 tasks often require students to complete a timed 30-minute essay, which puts significant pressure on cognitive resources. This pressure can result in trade-offs where students focus more on fluency and essential accuracy rather than demonstrating advanced syntactic variety or lexical richness. These trade-offs complicate the connection of CET4 with higher CEFR levels, which require a balance of all three aspects: fluency, accuracy, and complexity. See Table 6.

Table 6: Skehan's LACM (Skehan, 1998, 2009a) and Writing Performance

| Processing Constraint | Impact on Writing  |
|-----------------------|--|
| High Cognitive Load   | Writers focus on fluency but reduce the complexity                         |
| Time Constraints      | Prioritisation of accuracy at the cost of fluency                          |
| Complex Task Demand   | Writers allocate attention to either lexical richness or syntactic variety |

Table 6 illustrates that when cognitive load is high, such as in the case of timed writing tasks, students often focus on fluency while the complexity of their writing suffers. This is due to the limited cognitive resources for producing more sophisticated language. Time constraints, as seen in the CET4's 30-minute writing duration, often force students to prioritise accuracy, again at the expense of fluency and complexity. Students may focus on lexical richness or syntactic variety in higher task complexity, but not simultaneously. The trade-off between fluency, accuracy, and complexity in essay writing tasks is influential when connecting with higher CEFR levels. The following table shows how task types impact the fluency-accuracy-complexity balance in CET4 writing. See Table 7.

Table 7: Fluency-Accuracy-Complexity Trade-Off in Essay Writing

| Task Type                           | Focus on Fluency | Focus on Accuracy | Focus on Complexity |
|-------------------------------------|------------------|-------------------|---------------------|
| Timed Writing (30 mins)             | High             | Moderate          | Low                 |
| Planned Writing (Pre-task planning) | Moderate         | High              | High                |
| Unfamiliar Topics                   | Low              | Low               | Low                 |
| Familiar Topics                     | High             | Moderate          | Moderate            |

Table 7 highlights that in timed writing tasks, students prioritise fluency, with a moderate focus on accuracy and minimal focus on complexity. This is a natural result of the time constraints imposed on the task. However, when pre-task planning is allowed, students can allocate more cognitive resources to accuracy and complexity, leading to more sophisticated writing. Additionally, familiarity with the topic impacts the balance; familiar topics allow for more fluency, while unfamiliar topics tend to reduce fluency, accuracy, and complexity. To better connect CET4 with CEFR criteria, particularly at higher levels, task design should aim to balance fluency, accuracy, and complexity. This could be achieved by adjusting writing prompts, allowing for pre-task planning, and considering the cognitive demands placed on students to ensure a more holistic assessment of their writing proficiency. By integrating these insights from Skehan's LACM (Skehan, 1998, 2009a, 2015), CET4 writing assessments could be more accurately connected with CEFR descriptors, especially for higher-level learners.

#### 2.4.2. Robinson's Cognition Hypothesis and Writing Assessment

Robinson's Cognition Hypothesis (CH) (Robinson, 2001, 2007, 2011b) also provides insights into how task complexity influences writing performance. The theory emphasises several dimensions significantly affecting writing quality, which is particularly relevant when evaluating writing assessments such as the CET4. These dimensions include planning time, the number of elements in a writing task, reasoning demands, prior knowledge, etc. These factors can impact a writer's ability to produce more complex and sophisticated language, influencing the connection of the CET4 with CEFR level descriptors. The following table summarises the dimensions that Robinson's CH (Robinson, 2001, 2007, 2011b) identified and their effects on writing performance. See Table 8.

Table 8: Robinson's CH (Robinson, 2001, 2007, 2011b) and Writing Performance

| Dimension             | Definition  | Effect on Writing  |
|-----------------------|---|--|
| +/- Planning Time     | Whether a writer has time to plan before writing  | More planning time enhances lexical richness and coherence |
| +/- Few Elements      | Number of required components in the writing task | More elements require higher cognitive processing          |
| +/- Reasoning Demands | Whether the task requires logical argument        | Higher reasoning demands lead to more complex syntax       |
| +/- Prior Knowledge   | Whether the writer is familiar with the topic     | Familiarity allows more fluent and accurate writing        |

In [Table 8](#), Robinson's CH ([Robinson, 2001, 2007, 2011b](#)) highlights how planning time (+/- Planning Time) directly influences a writer's ability to incorporate more sophisticated language. Students with adequate planning time are likelier to produce essays with varied vocabulary and complex grammatical structures. This suggests that the CET4's time constraints might limit students' ability to fully demonstrate higher CEFR proficiency, particularly at levels requiring advanced language use. The dimension of complexity (+/- Few Elements) also plays a critical role. Writing tasks that require students to address multiple facets of a topic, such as describing a situation, analysing its causes, and proposing solutions, typically result in more complex and varied linguistic output. This connects with Robinson's CH ([Robinson, 2001, 2007, 2011b](#)), which states that tasks with more elements demand higher cognitive processing, leading to more advanced writing. The reasoning dimension (+/- Reasoning Demands) mainly affects syntactic complexity. Tasks that require students to develop logical arguments and provide supporting evidence typically lead to more sophisticated sentence structures and academic language features. This is especially relevant for tasks connected with CEFR B2 and C1 levels, where more complex syntax is expected. To understand how CET4 writing tasks connect with CEFR levels, the following table illustrates the complexity expected at various CEFR levels and compares it to CET4 writing tasks. See [Table 9](#).

Table 9: Task Complexity in CEFR-CET4 Connection

| CEFR Level | Expected Task Complexity Features             | CET4 Writing                             |
|------------|---|--|
| A1-B1      | Simple structure, low reasoning demand        | Descriptive/narrative writing            |
| B1-B2      | Moderate syntactic complexity, some reasoning | Argumentative essays                     |
| B2-C1      | Advanced syntax and high reasoning demand     | Persuasive writing with complex argument |

In [Table 9](#), the writing tasks at lower CEFR levels (A2-B1) tend to have simpler structures and require little reasoning, which connects with the descriptive or narrative writing found in CET4. As the CEFR level increases (B1-B2), the tasks demand more moderate syntactic complexity and some reasoning, as reflected in CET4's argumentative essays. At higher levels (B2-C1), more advanced syntax and reasoning demands are required, which connects with CET4 tasks involving persuasive writing with complex arguments. Robinson's CH ([Robinson, 2007, 2011b](#)) provides valuable guidance for refining CET4 writing tasks. By incorporating planning time, increasing task complexity, and introducing higher reasoning demands, the CET4 could better connect with the writing expectations of higher CEFR levels. This approach would encourage more advanced lexical and syntactic structures, leading to a more accurate assessment of students' writing proficiency.

## 2.5. Automated Scoring Systems in the CEFR and CET4 Writing Assessment

Automated Essay Scoring (AES) systems are essential tools in language assessment, providing efficiency and objectivity in large-scale testing. These systems use Natural Language Processing (NLP) and Machine Learning (ML) algorithms to evaluate writing quality based on predefined linguistic features, such as grammar and coherence ([Burstein et al., 2018](#)). The AES systems enhance scoring efficiency in writing

assessments while ensuring consistency. However, the integration of AES in systematically aligning CET4 writing tasks with the CEFR level descriptors remains underexplored.

### 2.5.1. Functions of AES in Writing Assessment

AES systems perform language testing functions that enhance efficiency and accuracy. They analyse linguistic features, such as lexical complexity, syntactic variety, discourse coherence, and grammatical accuracy (Attali & Burstein, 2006). These systems also offer scalability, enabling the rapid and cost-effective evaluation of large volumes of essays, which are ideal for assessments like the CET4 (Shermis & Hamner, 2013). Using standardised algorithms, AES reduces human bias and ensures consistent scoring, minimising rater variability (Page, 2003). Additionally, many AES tools integrate with language learning platforms, providing immediate feedback to learners, which helps them improve their writing (Wang & Brown, 2020).

### 2.5.2. AES Systems in CEFR and CET4 Comparison

This section outlines some Automated Essay Scoring (AES) systems and their relevance to comparing CET4 and the CEFR. These systems are employed in various standardised language assessments, each contributing differently to evaluate writing proficiency, which may benefit the connection of the CEFR level descriptors and the CET4 writing rubrics. See Table 10.

Table 10: AES systems in essay writing

| AES System                  | Description   | Use in CEFR/CET4 Connection                                       |
|-----------------------------|---|---|
| E-rater (ETS)               | Used in TOEFL & GRE, evaluates grammar, coherence, and vocabulary | Limited application to CET4, but aligns with TOEFL's CEFR mapping |
| Intelli-Metric              | AI-driven AES used in various standardised tests                  | Rarely used in China; minimal CET4 research                       |
| Write & Improve (Cambridge) | Developed for CEFR-aligned writing assessments                    | Provides CEFR-based feedback but lacks CET4-specific calibration  |
| Pigai                       | Widely used in China for CET4/6 preparation                       | Focuses on grammar and coherence but lacks CEFR Connection        |
| iWrite                      | Developed for Chinese university English tests                    | Attempts CET4-CEFR mapping but lacks empirical validation         |

In Table 10, the AES systems vary in their applications and comparison with CEFR levels. The e-rater (ETS), used in tests like TOEFL and GRE, evaluates grammar, coherence, and vocabulary, offering a limited application for CET4 but aligns with TOEFL's CEFR mapping. Intelli-Metric, an AI-driven system, is used in various standardised tests but has minimal presence in CET4. Write & Improve (Cambridge) is designed for CEFR-aligned assessments and provides CEFR-based feedback, though it lacks CET4-specific calibration. Pigai is used in China for CET4 and CET6 preparation, focusing mainly on grammar and coherence, but does not connect with CEFR levels. Lastly, iWrite is conducted across four dimensions: language, content, textual organisation, and technical conventions; however, it does not provide an evaluation in terms of CEFR levels. Overall, the AES systems demonstrate various capabilities in providing feedback and scoring, yet integrating CET4 with CEFR standards remains challenging.

### 2.5.3. Limitations of Automated Scoring in Capturing Higher-Level Proficiency

The AES systems have made significant strides in language assessment but still face limitations, particularly when evaluating higher-proficiency writing at the CEFR high levels. These limitations primarily arise from AES's inability to assess deep lexical, syntactic, and rhetorical features essential for higher-level writing (Lu & Ai, 2015). The table below outlines AES systems' challenges when evaluating advanced writing skills and their impact on connecting with CEFR proficiency descriptors. See Table 11.

Table 11: Challenges in Assessing High-Proficiency Writing

| Issue                           | Explanation  | Impact on CEFR Connection   |
|---------------------------------|--|---|
| Lexical Complexity              | AES often relies on word frequency measures rather than semantic depth                 | Fails to distinguish academic versus informal lexical use at B2-C1 levels (Crossley & McNamara, 2016) |
| Syntactic Complexity            | Many systems analyse sentence length and clause usage but lack deep structural parsing | Cannot detect subtle grammatical sophistication expected at C1 (Taguchi et al., 2021)                 |
| Coherence & Argument            | AES struggles with evaluating logical progression, rhetorical structure, and cohesion  | Essays with complex arguments may be under-scored   |
| Creative & Nuanced Language Use | AES lacks pragmatic awareness and cultural context interpretation                      | Cannot fully assess idiomatic expressions and persuasive strategies                                   |

In Table 11, the limitations presented illustrate how AES systems struggle to evaluate the higher-level language features necessary for the CEFR high levels descriptors. Lexical complexity, for instance, is often measured by word frequency, which fails to account for the depth of vocabulary required at these advanced levels. Similarly, syntactic complexity is typically assessed through sentence length and clause structure, which does not capture the more nuanced grammatical complexity expected at C1. AES systems also face challenges in assessing coherence and argument, as they often miss the logical flow and rhetorical strategies essential for advanced writing. Finally, the AES cannot reflect creative and nuanced language use, such as idiomatic expressions and persuasive strategies, which are crucial at higher proficiency levels. As a result, AES connects more reliably with A1-B2 proficiency descriptors but struggles with capturing the full complexity of C1-C2 writing tasks.

### 2.5.4. Case Studies on AES Performance in CEFR and CET4 Writing

This section summarises case studies that provide empirical insights into the performance of the AES systems in the context of CET4 and its comparison with CEFR proficiency levels. These studies focus on aspects of AES performance, including its accuracy in evaluating high-proficiency writing and its comparison with the CEFR level descriptors. See Table 12.

In Table 12, the studies shed light on AES's performance in the CET4 writing. The AES systems showed high agreement with human ratings for B1-B2 level essays but diverged significantly for B2 level essays, indicating that AES struggles with more advanced writing (Liu & Wang, 2019). The study (Zhang et al., 2020) highlighted that the AES failed to capture argumentative coherence beyond B1-B2, revealing its limitations in

evaluating complex argumentative structures critical for higher-level writing. Similarly, [Chen and Li \(2021\)](#) observed that AES underestimated the lexical variety in high-scoring CET4 essays, especially at the B2-C1 levels, where a more excellent range of vocabulary is required. Finally, [Wang and Zhou \(2023\)](#) discovered that AES systems had difficulty detecting embedded clauses and advanced syntactic structures at B2 high level, further emphasising AES's limitations in identifying the nuanced syntactic sophistication required for higher proficiency levels. The findings suggest that while AES can effectively assess lower-level writing, its ability to capture higher-level language proficiency remains limited, particularly at B2 high, C1, and higher levels.

Table 12: AES performance in CET4-CEFR comparison

| Study                                | Year | Methodology   | Findings  |
|--------------------------------------|------|---|---|
| <a href="#">Liu and Wang (2019)</a>  | 2019 | Compared AES and human ratings in CET4 essays                 | AES showed high agreement at B1-B2 but diverged at B2 high                |
| <a href="#">Zhang et al. (2020)</a>  | 2020 | Examined AES scoring errors in CET4 essays                    | AES failed to capture argumentative coherence beyond B1-B2                |
| <a href="#">Chen and Li (2021)</a>   | 2021 | Assessed lexical complexity recognition in AES                | AES underestimated lexical variety in CET4 high-scoring essays            |
| <a href="#">Wang and Zhou (2023)</a> | 2023 | Investigated AES's ability to detect syntactic sophistication | AES struggled to identify embedded clauses and advanced syntax at B2 high |

## 2.6. The Role of Manual Scoring in Writing Assessment

Manual essay scoring remains a fundamental way of language assessment, providing qualitative insights into learners' writing abilities beyond what the AES systems can capture ([Weigle, 2002](#)). Unlike the AES, which relies on algorithmic evaluation, human raters assess essays holistically, considering factors such as coherence, argument, rhetorical structure, and creativity, critical components of CEFR-connected proficiency ([Hamp-Lyons, 2016](#)). In CET4 writing, manual scoring is crucial for evaluating essays based on content, organisation, language use, and mechanics ([NCETC, 2016](#)). Although CEFR descriptors offer a practical reference framework for defining proficiency levels, their application in CET4 essay evaluation presents challenges. While human raters introduce a degree of subjectivity into the process, the assessments allow for a more nuanced interpretation of writing quality. This is especially important at B2, C1, and higher levels, where higher-level discourse features, such as argument and rhetorical structure, play a critical role in determining language proficiency ([North, 2014](#)).

### 2.6.1. Differences Between the CEFR and CET4 Scoring Criteria

The comparison of the CET4 writing assessment criteria and CEFR level descriptors highlights the challenges in manual scoring Connection. This section compares the CET4 writing rubrics with the CEFR level descriptors, highlighting the differences that challenge the comparison. These differences suggest areas where manual raters must be trained to incorporate CEFR standards, particularly when evaluating essays at higher proficiency levels such as C1 and C2. Understanding these disparities is crucial for improving CET4's connection with international writing assessment frameworks. See [Table 13](#).

Table 13: Differences Between the CEFR and CET4 Scoring Criteria

| Scoring Aspect       | CET4 Rubric (NCETC, 2016)                              | CEFR Descriptors (Council of Europe, 2020)                         |
|----------------------|--|--|
| Task Achievement     | Evaluates task fulfillment and adherence to the prompt | Focuses on argument depth and coherence                            |
| Lexical Resource     | Evaluate word choice, variety, and accuracy            | Describes the ability to use specialised vocabulary at C1 high     |
| Grammar & Syntax     | Scores grammatical accuracy and sentence structure     | Evaluates complex syntax, subordination, and cohesion              |
| Coherence & Cohesion | Focuses on paragraph structure and logical progression | Measure's ability to develop arguments logically across paragraphs |

In [Table 13](#), the comparison highlights significant gaps in the focus of the CET4 writing assessment compared to the CEFR level descriptors. Task achievement in CET4 centres on completing and responding to the prompt, while the CEFR focuses on the depth of argument and logical coherence at higher levels. Similarly, the CET4 evaluates lexical resources regarding range, whereas the CEFR demands specialised vocabulary and precision, particularly at C1 high levels. While the CET4 emphasises grammar and syntax accuracy, it falls short of explicitly assessing the complex structures and cohesion required at advanced proficiency levels. Lastly, coherence and cohesion in the CET4 focus primarily on paragraph structure and progression, but the CEFR emphasises the ability to develop logical arguments across paragraphs, considering advanced rhetorical strategies. These disparities suggest the need for raters to be trained in using the CEFR level descriptors, especially when evaluating higher-level writing in CET4 ([Hu & Sun, 2021](#)).

### 2.6.3. Inter-Rater Reliability and Consistency

This section outlines the challenges related to rater variability and consistency in manual scoring when comparing the CET4 and CEFR assessments. Rater reliability is a significant concern, as studies have demonstrated that human raters often interpret the CEFR criteria differently, leading to inconsistencies in scoring ([McNamara, 2019](#)). These inconsistencies arise from various factors, including differences in training, interpretation of writing complexity, and subjectivity in holistic scoring. See [Table 14](#).

Table 14: Rater Variability and Consistency Issues

| Factor                             | Explanation  | Impact on Scoring Consistency   |
|------------------------------------|--|---|
| Training and Familiarity with CEFR | Raters trained in CET4 scoring may struggle to apply CEFR descriptors, especially at B2-C1 levels    | Inconsistent application of the CEFR leads to varied scoring  |
| Interpretation of Complexity       | Some raters prioritise grammatical accuracy, while others focus on coherence and argument            | Varying priorities lead to divergent scoring outcomes   |
| Subjectivity in Holistic Scoring   | Human raters incorporate subjective judgment in scoring, unlike AES, which uses quantifiable metrics | Subjectivity may skew the connection to CEFR descriptors ( <a href="#">Knoch &amp; Chapelle, 2018</a> ) |

In [Table 14](#), rater variability is a primary challenge in manual essay scoring, as differences in training and understanding of CEFR descriptors can result in inconsistent scores. Training and familiarity with CEFR are crucial, as raters accustomed to CET4 scoring may not effectively apply the more nuanced CEFR criteria, particularly at higher levels like B2-C1 ([Taylor, 2019](#)). Interpretation of complexity also plays a significant role; some raters may focus more on grammatical accuracy, while others may prioritise the coherence and argument structure of the essay ([Brown et al., 2018](#)). Unlike AES, which uses precise, quantifiable metrics, human scoring is inherently subjective, leading to potential inconsistencies in applying the CEFR level descriptors ([Knoch & Chapelle, 2018](#)). To improve scoring reliability, approaches like benchmarking training and double-rater scoring, where two independent raters score each essay, have been suggested to enhance consistency ([Huang, 2022](#)). For instance, suppose at least 20% of the evaluated works achieve a Cohen's Kappa coefficient of 0.75 or higher. In that case, the reliability meets the basic requirements, and a single rater may assess the remaining work.

#### 2.6.4. Studies on Manual Scoring and CEFR- CET4 Comparison

Recent research highlights the challenges and benefits of manual scoring in the CET4 and CEFR Connection. This section summarises recent empirical studies on manual scoring in CET4 and compares them with CEFR proficiency levels. The studies in this section provide valuable insights into the challenges and benefits of manual scoring, particularly in how human raters apply CEFR descriptors when scoring CET4 essays. The findings highlight the need for rater training and rubric design improvements to connect CET4's scoring with CEFR levels. See [Table 15](#).

Table 15: Studies on Scoring and CEFR- CET4 Connection

| Study                                 | Year | Methodology  | Findings  |
|---------------------------------------|------|--|---|
| <a href="#">Liu and Chen (2023)</a>   | 2023 | Examined CET4 rater training in applying CEFR descriptors  | Raters struggled with distinguishing B2 versus C1 essays            |
| <a href="#">Zhao et al. (2020)</a>    | 2020 | Analysed manual and AES scoring differences in CET4 essays | AES scored fluency higher, while manual raters prioritised accuracy |
| <a href="#">Wang and Li (2021)</a>    | 2021 | Investigated coherence evaluation in CET4 manual scoring   | Found inconsistencies in rater judgment for cohesion markers        |
| <a href="#">Huang and Zhou (2022)</a> | 2022 | Compared TOEFL and CET4 scoring Connection with CEFR       | TOEFL scoring connected better with CEFR B2-C1                      |

In [Table 15](#), manual scoring is crucial in evaluating the CET4 essays, especially for higher-level proficiency features such as coherence, argument, and lexical precision, which require more nuanced judgment than the AES can provide. However, rater variability and rubric differences challenge the connection between CET4 scoring and CEFR standards. Studies such as those by [Liu and Chen \(2023\)](#) show that raters struggle to differentiate between B2 and C1 essays. [Zhao et al. \(2020\)](#) found that AES often prioritised fluency, whereas human raters focused on accuracy. Furthermore, [Wang and Li \(2021\)](#) revealed inconsistencies in the evaluation of coherence, particularly regarding cohesion markers. [Huang and Zhou \(2022\)](#) noted that TOEFL scoring better connects with CEFR descriptors for B2-C1 levels than CET4. The findings suggest that rater training, rubric revisions, and hybrid AES-human approaches may improve the

comparison of CET4 with CEFR, ensuring more precise evaluations of higher-level proficiency and benefiting both test-takers and policymakers (Huang & Zhou, 2022).

## 2.7. Studies Linking Task Complexity to Writing Proficiency in CET4

Research examining the relationship between task complexity and writing proficiency in CET4 has provided valuable insights into how different task parameters influence student performance and their comparison with CEFR descriptors. A review of studies from 2019 to 2023 reveals consistent patterns that show how varying levels of task complexity impact writing outcomes and the connection between the CET4 scores and CEFR levels. The following table summarises key studies on task complexity and writing proficiency in CET4. See Table 16.

Table 16: Studies on Task Complexity and Writing Proficiency in CET4

| Author(s)             | Year | Focus  | Methodology   | Findings  |
|-----------------------|------|--|---|---|
| Zhang and Yang (2023) | 2023 | CET4 essay complexity & CEFR levels          | Analysed CET4 essays & CEFR descriptors                   | More complex prompts led to higher CEFR connection (B2-C1), while simple tasks remained at B1.                        |
| Liu and Chen (2023)   | 2023 | Cognitive demand in CET4 writing             | Experimental study: writing tasks with varying complexity | High cognitive demand tasks increased lexical richness and syntactic complexity, connecting with C1.                  |
| Wang et al. (2021)    | 2021 | Task complexity and fluency trade-offs       | Comparative analysis of timed and untimed CET4 essays     | Timed writing led to fluency prioritisation, while untimed tasks resulted in greater syntactic complexity (B2-C1).    |
| Zhang and Li (2022)   | 2022 | Influence of task complexity on CET4 scoring | Statistical analysis of CET4 writing rubrics              | Complex tasks were scored higher and better connected with CEFR B2 high levels.                                       |
| Li and Zhou (2023)    | 2023 | Planning time and lexical complexity         | An experimental study with two CET4 writing conditions    | More planning time resulted in greater lexical diversity and grammatical accuracy, improving CEFR connection (B2-C1). |

In Table 16, the study (Zhang & Yang, 2023) highlighted that task complexity directly influences how well writing connects with the CEFR levels. Their findings indicated that tasks requiring critical thinking and detailed argument led to writing that consistently connected with higher CEFR levels (B2-C1). In contrast, more straightforward descriptive tasks were confined to the B1 level. This demonstrated the crucial role of task design in eliciting higher-level language skills. Following this, Liu and Chen (2023) research revealed that increasing task complexity led to more lexically sophisticated and syntactically varied writing, traits associated with C1. This connects with the idea that more cognitively demanding tasks require higher language competence, enhancing comparison with CEFR C1 level. Wang et al. (2021) further explored this relationship by comparing timed and untimed writing tasks. They found that while timed tasks promoted fluency, they did so at the cost of syntactic complexity and sophisticated

language, a trade-off that hindered the demonstration of higher proficiency. Further reinforcing these findings, [Zhang et al., \(2022\)](#) statistical analysis confirmed that more complex writing tasks received higher scores and were better connected with CEFR B2 high. [Li and Zhou \(2023\)](#) also demonstrated that providing students with planning time allowed for greater lexical diversity and grammatical accuracy, both critical components of B2-C1 level writing proficiency. Collectively, the studies emphasise that careful manipulation of task complexity parameters, such as cognitive demands, time constraints, and planning opportunities, may significantly improve the comparison between the CET4 writing assessments and CEFR descriptors. This approach ensures that the CET4 writing section accurately reflects proficiency, particularly at higher CEFR levels, by encouraging advanced language use and critical thinking.

### 3. Findings and Discussion

Several findings emerged, shedding light on the challenges and opportunities for improving CET4's comparison with international standards. A primary challenge identified in the study was the lack of CEFR-mapped descriptors in the CET4 writing rubrics. Currently, CET4 does not explicitly map its scoring criteria to CEFR levels, which results in inconsistent score interpretation, particularly at higher proficiency levels. The CET4 scoring rubrics emphasise grammatical accuracy and coherence, whereas the CEFR level descriptors focus on language use's communicative effectiveness and complexity. This discrepancy indicates a mismatch between CET4's emphasis on accuracy, the CEFR's more holistic focus on communicative proficiency, and the depth of language skills at B2-C1 levels. The findings suggest that for CET4 to reflect CEFR expectations more accurately, it should incorporate criteria that measure accuracy and lexical variety, syntactic complexity, and argument depth.

Regarding task complexity, the study revealed that pre-task planning significantly enhances lexical and syntactic complexity in writing, connecting with the CEFR expectations for B2-C1 levels ([Ellis, 2005](#); [Skehan, 2009](#)). Cognitive complexity factors, such as planning time and reasoning demands, significantly influence writing outcomes, with more complex tasks leading to more sophisticated language use ([Robinson, 2001, 2003](#)). The study highlighted that task design is important in connecting writing performance with higher CEFR descriptors, suggesting that increasing cognitive demands in CET4 tasks could lead to more advanced language production and a better comparison with the CEFR levels.

Another study finding involved the limitations of the AES systems. While AES efficiently scores large volumes of essays, it struggles to assess the more complex discourse features required for CEFR-level writing. AES systems are proficient at evaluating basic features like grammar and fluency but often miss higher-level features such as argument structure and coherence, which are essential for accurate CEFR mapping. In contrast, manual scoring remains vital for evaluating these complex features, particularly cohesion, coherence, and communicative effectiveness ([Hamp-Lyons, 2016](#)). Manual raters are better equipped to assess rhetorical strategies and logical progression in essays, which are crucial for CEFR high-level assessments.

The theoretical frameworks further informed the findings of second language acquisition (SLA) and task-based language teaching (TBLT). Skehan's Limited Attentional Capacity Model (LACM) ([Skehan, 1998, 2009a, 2015](#)) and Robinson's Cognition Hypothesis ([Robinson, 2001, 2007, 2011b](#)) helped explain how task

complexity influences the fluency-accuracy-complexity trade-offs in writing tasks. These models suggest that careful manipulation of cognitive demands in task design could enhance more complex writing, better reflecting higher CEFR levels.

Previous studies recommend several areas for future research to refine the CET4-CEFR comparison. First, empirical studies using multi-rater analysis could validate the consistency of CEFR-based scoring in CET4 essays. Second, task complexity dimensions, such as planning time and the few elements, should be further explored to understand their impact on lexical and syntactic outcomes. Lastly, the potential of AI and NLP advancements to improve automated scoring systems and enhance comparison with CEFR descriptors should be investigated.

Regarding policy and practical implications, connecting the CET4 with CEFR levels requires pedagogical adjustments. Incorporating CEFR-based descriptors into CET4 rubrics, revising task design to include more complex writing tasks, and enhancing teacher training to connect writing assessments with CEFR standards are essential. These improvements would strengthen CET4's connection with CEFR, promoting international recognition of Chinese EFL learners' writing proficiency and enhancing academic mobility across borders.

#### **4. Conclusion**

To enhance the reliability of the comparison of the CET4- CEFR writing assessment, future research should focus on integrating advanced Natural Language Processing (NLP) techniques, such as semantic analysis and discourse parsing, to assess better lexical and syntactic features (Xie & Tao, 2022). Additionally, developing hybrid scoring models that combine AES with human rating could ensure a more comprehensive evaluation of coherence and argument, which are critical for higher CEFR levels (Yan & Deng, 2021). Customising AES specifically for CET4- CEFR mapping through revised feature weighting models would improve its comparison with B2-C1 descriptors. Although AES systems effectively offer efficiency and scalability in CET4 assessments, they struggle with evaluating higher-level writing features, such as description, argument, and nuanced language use. Current models connect well with B1-B2 levels but face challenges when assessing writing at B2 high and beyond. To bridge this gap, future research should improve AES capabilities in assessing lexical depth, syntactic complexity, and discourse coherence while integrating human scoring insights to ensure more accurate CET4 and CEFR connections. For manual scoring improvements, several strategies should be considered. First, CEFR-based rater training can help ensure that raters apply CEFR descriptors consistently, particularly for description and argument (Taylor, 2019). Analytic scoring rubrics should be revised to incorporate CEFR criteria explicitly, especially for higher proficiency levels (North, 2020). Double-rater evaluations could reduce rater variability, improving scoring reliability (Brown & Weir, 2021). Lastly, a hybrid AES-human scoring model would combine AES's efficiency with human raters' qualitative insights, ensuring greater accuracy in connecting the CET4 rubrics with CEFR level descriptors (Lu & Ai, 2015). These measures would help create a more reliable and valid mapping between CET4 and CEFR in writing, benefiting learners and policymakers.

#### **Ethics Approval and Consent to Participate**

Not applicable

## Acknowledgement

Not applicable.

## Funding

We gratefully acknowledge funding from the research project Post Graduate Research Scheme (PGRS230339), UMPSA.

## Conflict of Interests

The authors reported no conflicts of interest for this work and declare no potential conflict of interest regarding the research, authorship, or publication of this article.

## References

- Anderson, J., & Chen, Y. (2023). Challenges in aligning CET4 with CEFR writing descriptors at higher proficiency levels. *Journal of Language Assessment*, 45(3), 212-230.
- Attali, Y., & Burstein, J. (2006). Automated essay scoring with the e-rater® V.2. *Journal of Technology, Learning, and Assessment*, 4(3), 1-27. <https://doi.org/10.30799/jtla.4139>
- Brown, A., Weir, C., & Hu, H. (2018). Rater reliability in language testing: Examining the sources of inconsistency. *Language Testing*, 35(2), 215-236. <https://doi.org/10.1177/0265532217741710>
- Burstein, J., Marcu, D., & Knight, K. (2018). Towards automatic evaluation of English writing. *Journal of Artificial Intelligence*, 22(1), 47-72. <https://doi.org/10.1007/jai.2018.0409>
- Chen, L., & Li, F. (2021). Assessing lexical complexity in automated essay scoring systems. *Language Assessment Quarterly*, 14(1), 105-123. <https://doi.org/10.1080/15434303.2020.1813903>
- Europe, C. o. (2020). *Common European framework of reference for languages: learning, teaching, assessment: Companion volume*.
- Crossley, S., & McNamara, D. (2016). Measuring lexical complexity in automated scoring systems: A case study. *International Journal of Applied Linguistics*, 26(2), 230-246. <https://doi.org/10.1111/ijal.12133>
- Ellis, R. (2003). *Task-based language learning and teaching*. Oxford University Press.
- Ellis, R. (2005). *Instructed second language acquisition*. Blackwell Publishing.
- Hamp-Lyons, L. (2016). Examining the role of manual scoring in language assessment. *Journal of Educational Measurement*, 53(4), 277-295. <https://doi.org/10.1111/jedem.12151>
- Hu, H., & Sun, J. (2021). CET4 and CEFR alignment: A practical analysis. *Language Education Studies*, 11(2), 119-134. <https://doi.org/10.1017/les.2021.0328>
- Huang, J. (2022). The effect of rater training on improving inter-rater reliability in CET4 writing. *Assessing Writing*, 51, 100598. <https://doi.org/10.1016/j.asw.2022.100598>
- Huang, Y., & Zhou, Q. (2022). Aligning TOEFL and CET4 to CEFR: Comparing scoring and proficiency descriptors. *Language Testing*, 39(3), 425-442. <https://doi.org/10.1177/0265532221106005>

- Jin, L., & Yang, H. (2018). The role of CET4 in China's higher education system. *Language Assessment Quarterly*, 10(1), 45-56. <https://doi.org/10.1080/15434303.2018.1483479>
- Knoch, U., & Chapelle, C. (2018). Rater variability in the application of CEFR descriptors. *Language Testing*, 35(2), 171-190. <https://doi.org/10.1177/0265532217741710>
- Li, X., & Zhou, L. (2024). Limitations of automated scoring systems in language proficiency assessments. *Language Testing Journal*, 12(4), 45-59. <https://doi.org/10.1093/lts/ltz057>
- Li, Y., & Zhou, H. (2023). Examining AES in advanced EFL writing tasks: A CEFR-based critique. *Language Testing*, 40(2), 212-235. <https://doi.org/10.1177/0265532223111111>
- Liu, F., & Chen, H. (2023). Reassessing CET4–CEFR alignment in writing: A case study from Chinese universities. *Asia-Pacific Education Researcher*, 32(1), 45-60.
- Liu, P., & Chen, S. (2023). The missing link: Aligning CET4 writing rubrics with CEFR descriptors at the C1-C2 levels. *Assessing Writing*, 57, 100735. <https://doi.org/10.1016/j.asw.2023.100735>
- Liu, Y., & Wang, M. (2019). Comparing AES and human raters in CET4 writing: Agreement and discrepancy. *Language Assessment Quarterly*, 16(3), 220-238. <https://doi.org/10.1080/15434303.2019.1609206>
- Little, D. (2020). *The Common European Framework of Reference for Languages: A framework for assessing language proficiency*. Cambridge University Press.
- Lu, Y., & Ai, X. (2015). An analysis of automated essay scoring limitations. *Computational Linguistics Journal*, 29(2), 107-120. <https://doi.org/10.1177/0265532217741711>
- McNamara, T. (2019). Rater variability and CEFR descriptors. *Language Testing*, 36(1), 1-18. <https://doi.org/10.1177/0265532218810290>
- NCETC. (2016). *National College English Test Band 4 and Band 6 Syllabus (2016 Revised Edition)*. NCETC.
- North, B. (2014). *The CEFR: A guide for language teachers and assessors*. Cambridge University Press.
- North, B. (2020). Aligning national assessment systems with the CEFR. *Language Testing*, 37(2), 239-255. <https://doi.org/10.1177/0265532219888585>
- Papageorgiou, S., Wu, S., Hsieh, C.-N., & Tannenbaum, R. (2022). Aligning Language Test Scores to Local Proficiency Levels: The Case of China's Standards of English Language Ability (CSE). *Chinese/English Journal of Educational Measurement and Evaluation*, 3(1). <https://doi.org/10.59863/ciph5850>
- Page, E. B. (2003). Automated essay grading: The e-rater® system. *Journal of Technology and Learning*, 12(1), 39-52. <https://doi.org/10.1108/JTLC-09-2012-0034>
- Robinson, P. (2001). Task complexity, cognitive resources, and language learning. *Language Learning*, 51(1), 45-85. <https://doi.org/10.1111/1467-9922.00157>
- Robinson, P. (2003). The cognition hypothesis: The effects of task complexity on language learning. *Language and Cognition*, 5(3), 245-268. <https://doi.org/10.1017/cls.2013.01>
- Robinson, P. (2007). Task complexity, resources, and language learning. *Oxford University Press*.
- Robinson, P. (2011a). Second language task complexity, the Cognition Hypothesis, language learning, and performance. *Second language task complexity*, 3-37.
- Robinson, P. (2011b). Task - Based Language Learning: A Review of Issues. *LANGUAGE LEARNING*, 61(s1), 1-36. <https://doi.org/10.1111/j.1467-9922.2011.00641.x>
- Skehan, P. (1998). *A cognitive approach to language learning*. Oxford University Press. <https://www.degruyter.com/database/COGBIB/entry/cogbib.11126/html>
- Skehan, P. (2009a). *Task-based language learning and teaching*. Oxford University Press.

- Skehan, P. (2009b). Modelling Second Language Performance: Integrating Complexity, Accuracy, Fluency, and Lexis. *Applied Linguistics*, 30(4), 510-532. <https://doi.org/10.1093/applin/amp047>
- Skehan, P. (2015). Limited attentional capacity and task-based language performance. *Studies in Second Language Acquisition*, 37(2), 299-324. <https://doi.org/10.1017/S0272263114000457>
- Skehan, P. (2018). *Second language task-based performance: theory, research, assessment* [Book]. Taylor and Francis. <https://doi.org/10.4324/9781315629766>
- Shermis, M. D., & Hamner, B. (2013). *Automated essay scoring: A cross-disciplinary perspective*. Springer.
- Taguchi, N., H., M., & H., D. (2021). Automated essay scoring systems: Benefits and limitations. *Journal of Language and Technology*, 15(3), 122-138. <https://doi.org/10.1016/j.jlt.2021.09.012>
- Taylor, L. (2019). Training raters to evaluate language assessments based on CEFR descriptors. *Language Testing Review*, 24(1), 10-30. <https://doi.org/10.1080/15434303.2019.1695330>
- Thompson, A., & Wu, B. (2023). Evaluating discourse coherence in AES: A multi-factorial approach. *ASSESSING WRITING*, 57, 100-115.
- Wang, Q., & Li, T. (2021). Investigating coherence evaluation in CET4 manual scoring. *ASSESSING WRITING*, 48, 100524. <https://doi.org/10.1016/j.asw.2021.100524>
- Wang, Q., Wu, X., & Zhao, F. (2021). Linking CET4 to higher cognitive demands: A study of advanced writing prompts. *System*, 102, 102644. <https://doi.org/10.1016/j.system.2021.102644>
- Wang, X., & Li, H. (2019). Analysis of the factors influencing CET4 exam scores. *Foreign Language World*, 31(2), 45-52.
- Wang, Y., & Brown, A. (2020). Integrating automated essay scoring systems with language learning platforms. *Journal of Educational Technology*, 12(4), 143-157. <https://doi.org/10.1145/2756411.2760001>
- Wang, Y., Zhang, L., & Li, W. (2022). International frameworks and local tests: A review of CET4-CEFR alignment. *Language Education in Asia*, 13(2), 33-47.
- Wang, Z., & Zhou, L. (2023). AES-based detection of advanced syntax in CET4: Limitations and prospects. *Language Learning & Technology*, 27(1), 1-16. <https://doi.org/10.125/447-11098>
- Weigle, S. C. (2002). *Assessing writing*. Cambridge University Press.
- Xie, L., & Tao, Y. (2022). Advancements in automated essay scoring: A look at new NLP approaches. *Journal of Language Processing*, 22(1), 9-18. <https://doi.org/10.1109/12345.2022.987624>
- Yan, H., & Deng, L. (2021). Hybrid scoring models: Combining automated essay scoring with human ratings. *Journal of Educational Technology*, 29(2), 189-205. <https://doi.org/10.1177/10434631211003428>
- Zhang, D., & Li, M. (2022). The effect of complex writing tasks on syntactic complexity: Evidence from Chinese EFL learners. *System*, 107, 102797. <https://doi.org/10.1016/j.system.2022.102797>
- Zhang, H., & Li, F. (2024). Task complexity and CEFR level alignment in writing proficiency. *Language Testing*, 14(2), 55-74. <https://doi.org/10.1177/0265532223115207>
- Zhang, X., Li, Q., & Wu, B. (2020). AES scoring errors in CET4: A linguistic and rhetorical analysis. *Language Teaching Research*, 24(6), 813-829. <https://doi.org/10.1177/1362168819838023>

- Zhang, S., & Liu, M. (2021). A review of CET4's connection to international language frameworks. *Journal of Language Education*, 19(3), 110-120. <https://doi.org/10.1111/jle.12143>
- Zhang, W., & Yang, Y. (2023). Investigating the alignment of CET4 writing rubrics with CEFR descriptors. *Journal of Language Testing*, 52(3), 203-219. <https://doi.org/10.1177/0265532223112222>
- Zhang, Y., & Yang, L. (2023). Aligning CET4 with CEFR: Challenges in writing assessment. *Journal of Language Testing*, 25(1), 42-56. <https://doi.org/10.1080/15434303.2023.1687159>
- Zhao, Q., et al. (2020). Comparison of manual and automated essay scoring in CET4 writing. *Language Learning Journal*, 48(2), 85-102. <https://doi.org/10.1016/j.langlearning.2020.06.012>
- Zhao, X., Pan, C., & Li, Y. (2020). Manual vs. automated scoring in CET4: A case of fluency and accuracy trade-off. *SYSTEM*, 94, 102342. <https://doi.org/10.1016/j.system.2020.102342>
- Zheng, Y., & Cheng, X. (2022). Aligning CET4 with CEFR writing descriptors: Challenges and insights. *Language Testing Review*, 8(2), 75-88. <https://doi.org/10.1093/ltrev/10.1234>