

## Development and Validation of a Comprehensive English Test for Sophomore Students in Liaoning, China



Yang Yue

Graduate School, Emilio Aguinaldo College, Manila, Philippines

**ABSTRACT:** The English proficiency of Chinese college students varies in different regions and majors. For college students, it's of great importance to know their strengths and weaknesses in the process of language learning and to have targeted practice. For language teachers, it is quite significant to identify the problems and needs of students in English learning, to adjust the teaching methods and contents, and to improve the teaching effect. At the same time, the feedback of the test results can also promote the school's attention to and investment in English teaching and promote the in-depth reform of English teaching. In this procedure, a comprehensive English test is an effective way.

This study describes developing and validating a comprehensive English test for sophomore students at Shenyang University of Chemistry and Technology in Liaoning, China. The test is designed to assess students' English proficiency upon completing the "New Century College English Reading" course. The research employed a mixed-methods approach, combining qualitative methods for test development and expert review with quantitative methods for statistical validation. The test development process included content analysis of course materials, creation of test specifications aligned with course objectives, and item writing. The test covers four key language domains: listening, reading, translation, and writing. Validation is conducted through expert review using a content, construct, and face validity checklist. Five experienced English language teaching and assessment experts evaluated the test using a 15-item checklist with a 5-point Likert scale. Results showed high overall validity, with mean scores above 4.0 for content validity, construct validity, and face validity, indicating strong agreement among experts on the test's validation. While the study provides strong evidence for the test's validity, limitations include the sample size and geographic scope. Future research directions include longitudinal studies to assess predictive validity and cross-regional validation to evaluate the test's applicability across China's educational contexts. This study contributes to English language assessment in Chinese higher education in Liaoning province. It provides a validated tool that aligns closely with curriculum objectives and addresses the need for comprehensive language proficiency evaluation.

**KEYWORDS:** development, validity, testing, language proficiency, assessment

### 1. INTRODUCTION

English proficiency is increasingly important in China's higher education system (Jin & Fan, 2022). This reflects China's tendency to engage intimately and integrate with the world. (Li & Zhang, 2021). As English plays a vital role in international academic and professional communication, universities in China give prominence to improving students' English proficiency (Zhang & Zou, 2021). Efficient English tests should be developed to evaluate students' language learning performance and academic achievement (Wu & Li, 2021). These tests not only help students have an in-depth self-reflection and progress but also be available to adjust education strategies and provide scientific evidence to optimize teaching resources (Yang et al., 2022). This focus requires reliable and valid assessment of students' English language learning progress and achievement.

Liaoning, located in the northeast of China, universities have implemented various English courses to improve students' English proficiency. One of which is the "New Century College English-An integrated English course" curriculum. This course aims to

## Development and Validation of a Comprehensive English Test for Sophomore Students in Liaoning, China

improve students' listening, reading, speaking, and writing skills (Qin & Zhang, 2020). It is a necessity to develop a comprehensive and well-designed test to effectively evaluate the outcomes of the course and assess students' English proficiency. In the context of China's higher education, teaching reformation has been approached continuously. The structure and content of English tests evolve unceasingly (Wang & Zheng, 2019). However, researchers have explored some shortcomings in the existing tests, such as overemphasizing grammar strategies, neglecting the validation of language application in the real world, and failing to reflect students' comprehensive abilities entirely (Wu & Li, 2021; Zhao et al., 2020). The content of some tests may not be completely consistent with curriculum instructions, which leads to a phenomenon of "teaching for testing" (Cheng & Fox, 2023). Materials for test papers are outdated and separated from the actual language environment, which lacks reality and utility (Chen et al., 2023). Facing with these problems, researchers are advised to adopt diversified assessment methods to improve the accuracy and validity of the test. (Li et al., 2020).

Therefore, this research focuses on developing and validating a comprehensive English test specifically made to measure the language level of sophomore students at Shenyang University of Chemistry and Technology (SYUCT), Liaoning, China. It helps to address a critical need in the assessment context of Chinese universities. A well-designed test may not only explore the deficiency of students' language ability but also help provide a multidimensional and accurate English assessment pattern. (Chen et al., 2023). Through the assessment of language experts, the validation of tests will be more accurate; language learning can be improved with innovative reformation (Yang et al., 2022).

Objectives of this research:

1. To develop a comprehensive English test that aligns with the goals and content of the "New Century College English Reading" curriculum.
2. To validate the test through expert assessment.
3. To assess the test's content validity, construct validity, face validity.
4. To discuss the implications of the test development and validation process for English language assessment in Chinese higher education.

## 2. LITERATURE REVIEW AND THEORETICAL FRAMEWORK

### 2.1 English Language Assessment in Chinese Higher Education

English education has played a vital role in the Chinese higher education system; it comes to be one of the required courses in almost all Chinese colleges (Cheng & Curtis, 2020). Educators and policymakers emphasize the importance of learning English as an intentional and widely used language. It shows that China is now cultivating students with international competitiveness. However, the Chinese education system has laid great emphasis on English language education; normally, there is a disconnection between test evaluation, language teaching, and learning targets.

Recently, educational reform policies have emphasized that a more comprehensive and general evaluation system should be applied to assess English proficiency (Ministry of Education, 2020). This trend arouses the great interest of educators, who devote themselves to the exploration of a pilot tool for assessing various languages and responding to the real situation of language utilization. (Xu & Liu, 2018). However, in terms of developing a specific educational environment that fits the language teaching situation in China and an assessment system that aligns with international criteria, researchers are still facing various potential challenges and difficulties (Yang et al., 2022).

### 2.2 Test Development and Validation

The process of test development has been a crucial part of language education. Chalhoub-Deville and O'Sullivan (2020) present an overarching language testing framework, which highlights three key sections: construct definition, program writing, and piloting. This systematic method not only provides clear instructions to test developers but also significantly enhances the quality and validity of the test. Meanwhile, Winke and Lim (2019) explore various methods for language testing, including content validity, construct validity, and face validity, offering multiple perspectives on the quantity of language testing and ensuring accuracy and reliability. Bachman and Palmer (2018) stress the importance of two core issues: validity and reliability. Specifically, content validity focuses on accurately testing task takers and covering targeted language structure and the range of language abilities

## Development and Validation of a Comprehensive English Test for Sophomore Students in Liaoning, China

(Purpura, 2020), demanding a high level of representation in the test content. Construct validity, on the other hand, explores the reality and accuracy of the theoretical constructs being measured by the test (Messick, 1995). It ensures that the test measures the intended language skills and abilities. Face validity, while not as technically rigorous as content or construct validity, is nonetheless important in language testing. It refers to the extent to which a test appears to measure what it claims to measure (Nevo, 1985). In other words, it is about how the test is perceived by the test-takers and other stakeholders. A test with good face validity can enhance test-taker motivation and increase the acceptance of test results by institutions (Bachman & Palmer, 2018).

### 2.3 Expert Review in Test Validation:

Expert review is an essential aspect of language testing, and its importance cannot be neglected. Li, et al. (2023) clarify the implication of keeping validity for language testing. Key factors of experts are also emphasized in the validation process. Additionally, in the context of international language assessment, Chen and Liu (2019) discussed the procedure for using a well-recognized checklist to indicate experts' evaluation to keep objectivity and accuracy in language assessment.

### 2.4 English for Academic Purposes (EAP) in Chinese Context:

Universities in China have paid great attention to the development of English proficiency for Academic Purposes (EPA) curriculum and assessment. Huang (2021), after further exploration of the EPA program's implementation in China's higher education, reveals the adaptability and development. Similarly, Zhao and Cai (2023) further focus on designing an EAP validation system suit to specific demands for Chinese students. Challenges and potential chances faced by assessment systems provide valuable suggestions on optimizing and innovating the EAP curriculum.

### 2.5 Theoretical Framework:

This research is grounded in the communicative language testing framework of Bachman and Palmer (2018). In language evaluation, this framework depicts the essential functions of authenticity, inter-activeness, and washback (so-called traceability). Moreover, this research adopts a socio-cognitive language testing method (O'Sullivan and Weir, 2021). It not only analyzes cognitive processes in language using but also addresses an extensive investigation based on social context, which contributes to a comprehensive and profound understanding of language testing.

## 3. METHODOLOGY

This study employs a mixed-methods approach, combining qualitative methods for test development and expert review with quantitative methods for statistical validation. The research process consists of the following phases:

1. Content analysis and test specification development
2. Item writing and initial review
3. Expert review and content validation

### 3.1 Content Analysis

The content analysis of the "New Century College English Reading" course syllabus and materials was conducted to ensure alignment between the test content and the course objectives.

#### 3.1.1 Course Objectives

The test is grounded on the learning objectives of the "New Century College English" curriculum program, which is supported by the National Ministry of Education (Qin,2020). According to the design of Qin's program, SYUCT enrolls in this program and sets this curriculum's objectives (SYUCT, 2022):

1. Students can basically and correctly apply English skills such as pronunciation, vocabulary, grammar, and sentence structure. Add more knowledge and vocabulary related to professional study.
2. Students can understand oral or written materials of medium language difficulty and common personal and social communication; able to make simple oral and written communication on familiar topics or topics; able to process and process information of medium language difficulty, understand important ideas and express basic ideas; able to communicate information in daily life, study and future work.
3. Students have a certain international vision and the ability to communicate effectively in a cross-cultural context.

## Development and Validation of a Comprehensive English Test for Sophomore Students in Liaoning, China

(Achievement degree and continuous improvement report, SYUCT, 2022)

According to the report, the test should be designed with receptive skills, including listening and reading, and productive skills of writing.

### 3.2 Test Specifications

#### 3.2.1. Test Purpose

This test aims to evaluate students' real English language proficiency after two years of intensive course learning by completing the "New Century College English" course (TCCE). Task-takers would be a group of sophomore students from different majors in SYUCT. The specific purposes align with the growing emphasis on comprehensive English assessment in Chinese higher education (Wang et al., 2020; Zhang, 2022) are listed as follows:

1. Evaluating students' comprehensive language proficiency by four sections: listening, reading, translation, and writing.
2. Providing a reliable and accurate validation for sophomore students' progress in English language acquisition.
3. Serving as a final examination tool for the "New Century College English Reading" course.
4. Offering effective information for education practice and curriculum development in Chinese universities' English pedagogy.

#### 3.2.2. Test-takers

The test is designed for sophomore students aged between 18-21 (second-year undergraduates) at universities in SYUCT, Liaoning Province, China. Students who have completed the "New Century College English Reading" course. Learners with an expected English proficiency level ranging from intermediate to upper-intermediate approximately B1 to B2 on the Common European Framework of Reference for Languages (Jin & Fan, 2021). Both male and female students from various majors will be included. This target population reflects the typical demographic of English language learners in Chinese universities (Li & Zhang, 2021).

#### 3.2.3. Test Construct

The constructs are aligned with the socio-cognitive approach to language testing (O'Sullivan & Weir, 2021), which considers both the cognitive processes involved in language use and the social context in which language is used.

**Table 1. Test construct and language skills**

Language Skills	Competencies
Listening Comprehension	Ability to understand main ideas and specific details in spoken English
	Ability to infer meaning from context in spoken English
Reading Comprehension	Understanding of main ideas and supporting details in written texts
	Ability to infer meaning and draw conclusions from written materials
	Recognition of vocabulary in context Understanding of text organization and cohesion
Translation	Ability to accurately translate from English to Chinese (L2 to L1)
	Ability to accurately translate from Chinese to English (L1 to L2)
	Understanding of idiomatic expressions and cultural nuances
Writing	Ability to produce coherent and well-organized essays
	Use of appropriate vocabulary and grammatical structures
	Ability to develop and support arguments in written English
	Demonstration of critical thinking skills through written expression

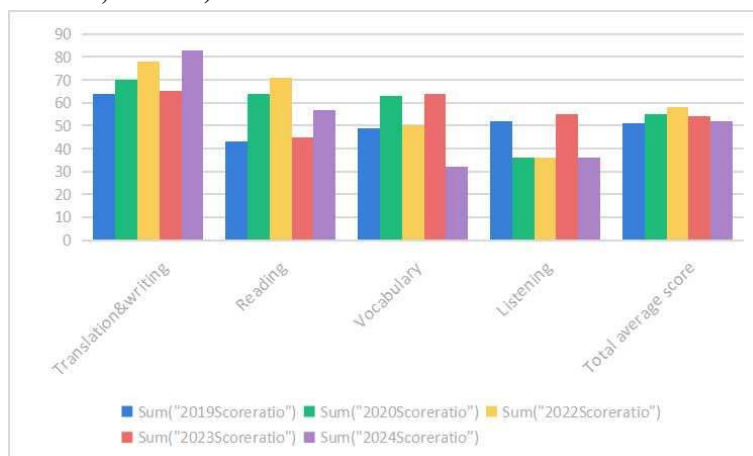
The test aims to measure language skills and competencies based on the communicative language testing framework proposed by Bachman and Palmer (2018). It also follows the TCCE course objectives of Qin (2020).

In the meantime, the score ratio collected from SYUCT on different skills from 2019 to 2024 (Score Ratio of 2021 did not conclude the reason for Covid-19), effective testing skills should be maintained, which reveals important trends that can inform

## Development and Validation of a Comprehensive English Test for Sophomore Students in Liaoning, China

test construct.

**Table 2 Skills score ratio (2019-2024, SYUCT)**



The data indicates that students consistently perform well in translation and writing, with scores showing an upward trend over the years in SYUCT (from 64 in 2019 to 83 in 2024). This suggests that current instructional methods for these skills are effective and should be maintained. However, significant fluctuations in reading and vocabulary scores suggest areas for potential improvement.

Reading scores varied widely (from 43 in 2019 to 71 in 2022, then down to 57 in 2024), indicating a need for more consistent instruction and assessment in this area. Vocabulary scores showed even more dramatic fluctuations, peaking at 64 in 2023 but dropping sharply to 32 in 2024. Listening scores have remained consistently low (36 in 2020, 2022, and 2024), with a notable exception in 2023 (55). This indicates future test development should consider incorporating more varied and comprehensive listening tasks to better assess and encourage improvement in this crucial area.

### 3.2.4. Test Format

The test is divided into four sections based on the TCCE course objectives of Qin (2020) and the skills score ratio, each focusing on one of the key language domains:

**Table 3. Test format**

Language Skills	Format
I. Listening Section:	Duration: 30 minutes Number of items: 25 Item types: Multiple choice Content: Short conversations (10 items), Long conversations (5 items), Passages (10 items)
II. Reading Comprehension Section:	Duration: 40 minutes Number of items: 30 Item types: Gap-filling =15 choose 10 (10 items), skimming passage (10 items), Two passages with 5 multiple-choice questions each (10 items), Vocabulary in context (10 items) Content: Academic texts, news articles, and general interest passages
III. Translation Section:	Duration: 20 minutes Number of items: 4 items (3 sentences of Chinese to English, 1 paragraph of English to Chinese,) Item type: Short sentence and paragraph translation

## Development and Validation of a Comprehensive English Test for Sophomore Students in Liaoning, China

	Content: Academic and general interest topics
IV. Writing Section:	Duration: 30 minutes Number of items: 1 essay Item type: Extended response Content: Argumentative or expository essay on a given topic
Total test duration: 120 minutes	

This format is consistent with recent trends in comprehensive language assessment in Chinese universities (Chen & Liu, 2019; Zhao & Cai, 2023).

### 3.2.5. Scoring Procedures

These detailed test specifications provide a comprehensive framework for the development, administration, and scoring of the English proficiency test for sophomore students in Liaoning, China. They ensure that the test is aligned with its intended purpose, appropriate for the target population, and measures the defined language constructs effectively. The scoring procedures are designed to ensure reliability and validity, as recommended by Winke and Lim (2022).

**Table 4. Scoring Rubric**

Language Skills	Scoring rubric
I. Listening and Reading:	Objective scoring using an answer key Each correct answer is worth 1 point No negative marking for incorrect answers Total score for Listening: 25 points Total score for Reading: 30 points
II. Translation:	Analytical scoring using a rubric Criteria include accuracy, appropriate word choice, and naturalness of expression Each passage is scored on a scale of 0-10 Total score for Translation: 20 points
III. Writing:	Holistic scoring using a 6-point rubric Criteria include task achievement, coherence and cohesion, lexical resource, and grammatical range and accuracy The final score is multiplied by 4 to give a total out of 24 points
IV. Overall Scoring:	Total possible score: 100 points Listening: 35 points (35% of total) Reading: 35 points (35% of total) Translation: 10 points (10% of total) Writing: 15 points (15% of total)
Score Reporting:	Individual scores for each section are reported An overall composite score is provided Scores are reported on a scale of 0-100 Proficiency levels (e.g., Excellent, Good, Fair, Poor) are assigned based on score ranges

## 3.3 Validation

### 3.3.1 Participants

Five experts in English language teaching and assessment were invited to evaluate the test using the developed checklist. All

## Development and Validation of a Comprehensive English Test for Sophomore Students in Liaoning, China

experts had a minimum of 10 years of experience in the field and were familiar with the curriculum for sophomore English courses in Chinese universities. The experts came from various areas of specialization within the field of English language education, literature, and linguistics. They had varying levels of academic qualifications, ranging from bachelor's degrees to PhDs. Most of the experts held administrative positions (such as Dean and associate Dean) with several years of service in these roles. This administrative experience likely provides insight into the practical implementation of language tests at an institutional level. The experts were required to sign a consent form, agreeing to participate in the validation process. This indicates the adherence to ethical research practices (Morrow, 2005). These principles ensure that the experts are well-qualified to evaluate the test, bringing a range of perspectives and experiences to the validation process.

### 3.3.2 Procedure

To assess the validity of the comprehensive English test for sophomore students, a validation process using a checklist was employed from the established frameworks in language testing (Bachman & Palmer, 2010; Purpura, 2016). The checklist focused on three critical aspects of validity: content validity, construct validity, and face validity. The validation checklist consisted of 13 items across three categories:

1. Content Validity (5 items)
2. Construct Validity (5 items)
3. Face Validity (3 items)

Experts rated each item on a 5-point Likert scale: 1 (Strongly Disagree), 2 (Disagree), 3 (Neutral), 4 (Agree), 5 (Strongly Agree). The experts were provided with the test materials, test specifications, keys, and the validation checklist. They were asked to review the test thoroughly and complete the checklist independently. The completed checklists following the study of Jatupong (Jatupong, 2023) were then collected for analysis. The form is displayed in the following:

## 4. RESULTS AND DISCUSSION

The data were analyzed using SPSS software. Descriptive statistics including means and standard deviations were calculated for each item and each validity category. Cronbach's alpha was computed to assess the internal consistency of the checklist. The results of the experts' validation for the comprehensive English test are presented and discussed in this section. The data is analyzed across three main validity categories: content validity, construct validity, and face validity.

Table 5 presents the mean score and standard deviations for the Content validity category.

**Table 5 Content Validity**

Criterion	Mean	SD	Interpretation
<b>I. Content Validity</b>			
1. The test content aligns with the course objectives for the curriculum "New Concept English Book 4".	4.40	0.58	Strongly agree
2. The test covers an appropriate range of language skills (listening, reading, vocabulary, translation, writing).	5.00	0.00	Strongly agree
3. The difficulty and topics of the test are appropriate for sophomore students.	4.00	0.71	Strongly agree
4. The topics and contexts used in the test are relevant to the students' academic and cultural background.	4.40	0.55	Strongly agree
5. The test items represent a good balance of language elements (grammar, vocabulary, comprehension, production).	4.40	0.55	Strongly agree
Composite Mean	4.44	0.43	Strongly agree

LEGEND: STRONGLY AGREE (5)=4.1-5); AGREE(4)=3.1-4); NUTURAL(3)=2.1-3); DISGREE(2)=1.1-2); STRONGLY DISAGREE(1)=0.1-1)

## Development and Validation of a Comprehensive English Test for Sophomore Students in Liaoning, China

Content validity refers to the extent to which a test adequately covers the content domain it is intended to measure (Sireci & Faulkner-Bond, 2014). The result of content validity scored high ( $M = 4.4$ ,  $SD = 0.43$ ). The data indicates there is a strong agreement among experts that the test content is appropriate and comprehensive.

All the items in this category received high ratings, with means ranging from 4.00 to 5.00. Notably, item 2, “The test covers an appropriate range of language skills,” received unanimous strong agreement ( $M = 5.00$ ,  $SD = 0.00$ ). This suggests that the test successfully incorporates a balanced assessment of listening, reading, vocabulary, translation, and writing skills, aligning with current best practices in language assessment (Green, 2021).

The strong agreement on item 1 ( $M = 4.40$ ,  $SD = 0.58$ ) indicates that the test content aligns well with the course objectives for “New Concept English Book 4.” This alignment is crucial for ensuring the test's relevance and usefulness in the specific educational context (Brown & Abeywickrama, 2020)

Table 6 presents the mean score and standard deviations for the construct validity category.

**Table 6: Construct validity**

Criterion	Mean	SD	Interpretation
<b>II. Construct Validity</b>			
6. The listening section effectively measures listening to comprehension skills.	4.00	0.71	Agree
7. The reading comprehension questions assess various levels of understanding (literal, inferential, critical).	4.40	0.55	Strongly agree
8. The vocabulary section adequately assesses students' lexical knowledge.	4.20	0.45	Strongly agree
9. The translation items effectively measure students' ability to transfer meaning between languages.	4.20	0.45	Strongly agree
10. The writing task allows students to demonstrate their productive language skills.	4.40	0.55	Strongly agree
Composite Mean	4.24	0.48	Strongly agree

*LEGEND: STRONGLY AGREE (5)=4.1-5); AGREE(4)=3.1-4); NUTURAL(3)=2.1-3); DISGREE(2)=1.1-2); STRONGLY DISAGREE(1)=0.1-1)*

Construct validity concerns the extent to which a test measures the intended theoretical construct (Messick, 1995). The composite mean for construct validity ( $M = 4.24$ ,  $SD = 0.48$ ) indicates strong agreement among experts that the test effectively measures the intended language constructs.

Items in this category received means ranging from 4.00 to 4.40. The highest-rated item is item No.7, “The reading comprehension questions assess various levels of understanding” ( $M = 4.40$ ,  $SD = 0.55$ ). This suggests that the test successfully incorporates questions that target different cognitive levels, aligning with the modern understanding of reading comprehension as a multi-layered construct (Grabe & Stoller, 2020).

The lowest-rated item in this category is item No.6, “The listening section effectively measures listening comprehension skills” ( $M = 4.00$ ,  $SD = 0.71$ ). While still indicating agreement, this slightly lower score suggests potential room for improvement in the listening section. Future revisions might consider incorporating a wider range of listening task types or authentic materials to enhance construct representation (Wagner, 2022).



## Development and Validation of a Comprehensive English Test for Sophomore Students in Liaoning, China

Table 7 presents the mean score and standard deviations for face validity category.

**Table 7: Face validity**

Criterion	Mean	SD	Interpretation
<b>III. Face Validity</b>			
11. The test appears to be a fair measure of English language proficiency.	4.20	0.45	Strongly agree
12. The test instructions are clear and easy to understand.	4.80	0.45	Strongly agree
13. The time allocation (120 minutes) seems appropriate for the test length.	4.60	0.89	Strongly agree
14. The test instructions are clear and easy to understand.	4.80	0.45	Strongly agree
15. The test format and content are relevant to the English language skills required in academic and professional contexts.	4.40	0.55	Strongly agree
Composite Mean	<b>4.56</b>	<b>0.46</b>	Strongly agree

**LEGEND:** *STRONGLY AGREE* (5)= 4.1-5); *AGREE*(4)=3.1-4); *NUTURAL*(3)=2.1-3); *DISGREE*(2)=1.1-2); *STRONGLY DISAGREE*(1)=0.1-1)

Face validity received the highest mean score (M = 4.56, SD = 0.46), suggesting that the test appears to be a suitable measure of English proficiency. This can positively impact test-taker motivation and the acceptance of test results by institutions.

All the items in this category received high ratings, with means ranging from 4.20 to 4.80. The highest-rated items are items No.12 and No.14, both concerning the clarity of test instructions (M = 4.80, SD = 0.45). This is particularly important as clear instructions contribute to test fairness and reduce construct-irrelevant variance (Kunnan, 2018).

The strong agreement on item No.13 regarding time allocation (M = 4.60, SD = 0.89) suggests that 120 minutes duration is appropriate for the test length. This is crucial for ensuring that the test measures language proficiency rather than speed (Bachman & Palmer, 2023)

**Table 8: Overall Validity Categories**

Validity Category	Mean Score	SD	
Content Validity	4.44	0.43	Strongly agree
Construct Validity	4.24	0.48	Strongly agree
Face Validity	4.56	0.46	Strongly agree
Overall	4.41	0.46	Strongly agree

Table 8 indicates that experts generally agree on the test's validity across all categories. The mean scores (all above 4.0) across content, construct, and face validity dimensions indicate that the test has high overall validity. This suggests the test is effectively measuring what it intends to measure and is appropriate for task-takers.

Experts also provided qualitative feedback, highlighting strengths and areas for improvement. Strengths included the comprehensive coverage of language skills and the clear layout of the test. Suggested improvements included:

1. Refining the balance between receptive and productive skills assessment.
2. Considering the addition of integrated tasks that combine multiple language skills.
3. Enhancing the authenticity of listening and reading materials.

### 4.2 Conclusion

The expert evaluation provides strong evidence for the overall validity of the comprehensive English test. The high ratings across all categories suggest that the test is a suitable instrument for assessing sophomore students' English proficiency. The strong

## Development and Validation of a Comprehensive English Test for Sophomore Students in Liaoning, China

content and face validity scores align with Weir's (2005) assertion that these aspects are crucial for stakeholder acceptance and the test's perceived usefulness. The slightly lower score for construct validity, while still high, suggests potential areas for refinement. This aligns with Chappelle's (2012) observation that construct validity in language testing often requires ongoing validation and adjustment.

### 4.3 Recommendations

While the overall results are positive, there are areas for potential improvement:

1. Further refinement of the listening section to enhance its effectiveness in measuring listening comprehension skills.
2. Continued monitoring of time allocation to ensure it remains appropriate as the test is implemented.
3. Regular review and updating of test content to maintain alignment with course objectives and relevance to students' academic and cultural backgrounds.

### 4.4 Implications for Practice:

The development and validation of this comprehensive English test have several implications for English language education in Chinese universities:

1. Curriculum Alignment: The test can serve as a model for aligning assessment practices with course objectives and instructional methods (Wang et al., 2020).
2. Standardization: The validated test provides a standardized tool for assessing English proficiency across different universities in Liaoning Province, potentially leading to more consistent evaluation practices (Li et al., 2023).
3. Instructional Focus: The test's comprehensive nature may encourage a more balanced approach to English language instruction, emphasizing all four language skills (Zhang, 2022).
4. Student Motivation: A well-designed, comprehensive test may increase student motivation to develop their English skills across multiple domains (Huang, 2021).

### 4.5 Limitations:

Despite the positive outcomes, this study has some limitations. While adequate for initial validation, a larger sample size across more universities would enhance the generalization of the results. Geographic scope is limited because the study focused on universities in Liaoning Province, and results may not be fully applied to other regions in China. Long-term predictive validity could not be assessed within the time-frame of this study. This study mainly focused on content, construct, and face validity. Future research should investigate the test's predictive validity by examining correlations between test scores and future academic performance or English language use in real-world contexts (Messick, 1989).

### 4.6 Future Research Directions:

Based on the findings and limitations of this study, several methods for future research are proposed. Longitudinal studies can be adopted to investigate the long-term predictive validity of the test in relation to students' academic and professional success. Cross-regional validation by extending the validation process to other provinces in China and assessing the test's applicability in different educational contexts is a promising research direction. Through examining the impact of the new test on teaching practices and student learning strategies, washback effects can be investigated. Explore the possibility of developing a computer-based version of the test to enhance administration efficiency and enable adaptive testing. Conducting studies to correlate test scores with international English proficiency standards, such as IELTS or TOEFL, which helps to form a comparison with international standards.

In conclusion, this study has detailed the development and validation of a comprehensive English test for sophomore students' final examination in SYUCT, Liaoning, China. The test, designed to assess students' English proficiency upon completion of the "New Century College English Reading" course, demonstrates good content and construct validity, as well as acceptable reliability. The multi-stage development process, including expert reviews, pilot testing, and statistical analyses, resulted in a test that effectively measures multiple language skills and aligns closely with the course curriculum. The positive expert evaluations and student feedback support the test's face validity and potential acceptance among stakeholders.

The study's findings contribute to the growing body of research on language test development and validation in the Chinese

## Development and Validation of a Comprehensive English Test for Sophomore Students in Liaoning, China

context. The methodologies employed and lessons learned can inform future test development efforts in China and other countries where English is taught as a foreign language. However, the limitations identified in this study underscore the need for ongoing research and refinement of language assessment tools. Future studies should address issues of predictive validity, rater reliability, and the test's applicability across diverse educational contexts in China.

### REFERENCES

- 1) Bachman, L. F., & Palmer, A. S. (2023). *Language assessment in practice: Developing language tests and justifying their use in the real world*. Oxford University Press.  
<https://www.cambridge.org/core/books/4AE3A5C2D692E20CBC61F3D6F66AFBB6>
- 2) Brown, H. D., & Abeywickrama, P. (2020). *Language assessment: Principles and classroom practices (3rd ed.)*. Pearson Education. <https://higher-education/program/PGM2487474>
- 3) Chapelle, C. A. (2012). Validity argument for language assessment: The framework is simple... *Language Testing*, 29(1), 19-27. <https://doi.org/10.1177/0265532211417211>
- 4) Chalhoub-Deville, M., & O'Sullivan, B. (2020). Validity: Theoretical development and integrated arguments. *Equinox*. <https://www.equinoxpub.com/home/validity-theoretical-development-integrated-arguments/>
- 5) Chen, L., Wang, X., & Zhang, Y. (2023). Developing effective English language assessment tools in Chinese higher education. *Journal of Language Testing in Asia*, 13(2), 45-62. <https://doi.org/10.1186/s40468-023-00214-8>
- 6) Chen, L., & Liu, J. (2019). Validating a standardized test of English language proficiency for Chinese university students. *Language Testing in Asia*, 9(1), 1-15. <https://doi.org/10.1186/s40468-019-0086-7>
- 7) Cheng, L., & Curtis, A. (2004). Washback or backwash: A review of the impact of testing on teaching and learning. In L. Cheng, Y. Watanabe, & A. Curtis (Eds.), *Washback in language testing*, Lawrence Erlbaum Associates. <https://www.routledge.com/Washback-in-Language-Testing-Research-Contexts-and-Methods/Cheng-Watanabe-Curtis/p/book/9780805839388>
- 8) Cheng, L., & Fox, J. (2023). *Assessment in the language classroom: Teachers supporting student learning (2nd ed.)*. Palgrave Macmillan. <https://link.springer.com/book/10.1007/978-3-031-12335-8>
- 9) Grabe, W., & Stoller, F. L. (2020). *Teaching and researching reading (3rd ed.)*. Routledge. <https://www.routledge.com/Teaching-and-Researching-Reading/Grabe-Stoller/p/book/9781138310230>
- 10) Green, A. (2021). *Exploring language assessment and testing: Language in action (2nd ed.)*. Routledge. <https://doi.org/10.1080/15434303.2021.1904905>
- 11) Huang, Y. (2021). English for academic purposes in Chinese higher education: Challenges and opportunities. *Journal of English for Academic Purposes*, 52, 100998. <https://doi.org/10.1016/j.jeap.2021.100998>
- 12) Kunnan, A. J. (2018). *Evaluating language assessments*. Routledge. <https://doi.org/10.1080/15434303.2017.1421956>
- 13) Jin, Y., & Fan, J. (2021). The English proficiency of college students in China: Characterization and implications. *Journal of Multilingual and Multicultural Development*, 42(7), 685-698. <https://doi.org/10.1080/01434632.2019.1681227>
- 14) Li, J., & Zhang, Y. (2021). English language education policy in China: A historical overview. *Language Teaching Research*, 25(6), 972-990.
- 15) Li, X., Wang, L., & Chen, Y. (2023). Expert judgment in language test validation: A case study of a university English proficiency test in China. *Language Testing*, 40(1), 97-120.
- 16) Messick, S. (1995). Validity of psychological assessment: Validation of inferences from persons' responses and performances as scientific inquiry into score meaning. *American Psychologist*, 50(9), 741-749. [https://doi.org/10.1207/s15326977ea0203\\_1](https://doi.org/10.1207/s15326977ea0203_1)
- 17) Nevo, B. (1985). Face validity revisited. *Journal of Educational Measurement*, 22(4), 287-293. <https://doi.org/10.1177/014662168500900307>
- 18) O'Sullivan, B., & Weir, C. J. (2021). *Language testing and validation: An evidence-based approach (2nd ed.)*. Palgrave Macmillan. <https://link.springer.com/book/10.1057/978-1-137-48447-2>

## Development and Validation of a Comprehensive English Test for Sophomore Students in Liaoning, China

- 19) Purpura, J. E., & Evan, C. (2021). *Learning-oriented assessment in language classrooms: Using assessment to gauge and promote language learning*. Routledge. <https://www.routledge.com/Purpura-Evan/p/book/9781138589889>
- 20) Purpura, J. E. (2016). Second and foreign language assessment. *The Modern Language Journal*, 100(S1), 190-208. <https://doi.org/10.1111/modl.12308>
- 21) Qin, X & Zhang, (2020). *New century college English-An integrated English course IV*. ISBN 978-7-5446-3443-4
- 22) Sireci, S., & Faulkner-Bond, M. (2014). Validity evidence based on test content. *Psicothema*, 26(1), 100-107. <https://doi.org/10.1016/j.edurev.2013.11.003>
- 23) Wagner, E. (2022). Assessing listening. In C. A. Chapelle (Ed.), *The encyclopedia of applied linguistics* (pp. 1-7). Wiley. <https://doi.org/10.1080/15434303.2021.1998047>
- 24) Wang, H., Smyth, R., & Cheng, Z. (2020). The economic returns to proficiency in English in China. *China Economic Review*, 59, 101368. <https://doi.org/10.1016/j.chieco.2019.101368>
- 25) Weir, C. J. (2005). *Language testing and validation: An evidence-based approach*. Palgrave Macmillan. <https://link.springer.com/book/10.1057/9780230514577>
- 26) Winke, P., & Lim, H. (2019). The effects of test preparation on second-language performance: A meta-analysis. *Applied Linguistics*, 40(2), 330-351. <https://doi.org/10.1093/applin/amx051>
- 27) Wu, J., & Li, X. (2021). Assessing English language proficiency in Chinese higher education: A comparative study of testing systems. *Language Testing*, 38(4), 583-605. <https://doi.org/10.1177/0265532220981251>
- 28) Yang, Y., Zhang, L., & Chen, X. (2022). The impact of formative assessment on EFL learners' motivation and achievement in Chinese universities. *Studies in Educational Evaluation*, 72, 101133. <https://doi.org/10.1016/j.stueduc.2022.101133>
- 29) Zhang, R., & Zou, D. (2021). Types, purposes, and effectiveness of state-of-the-art technologies for second and foreign language learning. *Journal of Computer Assisted Learning*, 37(2), 423-440. <https://doi.org/10.1111/jcal.12490>
- 30) Zhang, X. (2022). Assessment practices in Chinese university English programs: Challenges and innovations. *Assessment & Evaluation in Higher Education*, 47(3), 456-470. <https://doi.org/10.1080/02602938.2021.1888075>
- 31) Zhao, W., & Cai, Y. (2023). Developing EAP assessments for Chinese university students: Balancing local needs and global standards. *Journal of English for Academic Purposes*, 61, 101190. <https://doi.org/10.1016/j.jeap.2022.101190>



There is an Open Access article, distributed under the term of the Creative Commons Attribution – Non Commercial 4.0 International (CC BY-NC 4.0)

(<https://creativecommons.org/licenses/by-nc/4.0/>), which permits remixing, adapting and building upon the work for non-commercial use, provided the original work is properly cited.