

# iTEP Academic

# 2016

TECHNICAL REPORT



# TABLE OF CONTENTS

About this Report	4
Chapter 1 Introduction	5
Approach and Rationale for the Development of iTep Academic	6
Chapter 2 Detailed Description of iTep Academic	7
Theoretical Model for Language Assessment	7
Description of iTep Academic Scales	8
Grammar Section	8
Listening Section	9
Reading Section	11
Speaking Section	12
Writing Section	13
Assessment Administration	14
Delivery Method	14
Examinee Experience	14
Scoring/Grading	15
Proficiency Levels	15
Chapter 3 iTep Academic Development Process, Reliability, and Validity <sup>1</sup>	16
Development Process	16
Reliability	17
Internal Consistency Reliability	17
Test-Retest Reliability	19
Rater Agreement	20
Validity	23
Content Validity	23
Convergent and Discriminant Validity	24
References	26
Appendix A: Examinee Pre-Assessment Modules and Instructions	27
Appendix B: Speaking Scale Rater Scoring Rubric	30
Appendix C: Writing Scale Rater Scoring Rubric	32
Appendix D: iTep Ability Guide	34
Appendix E: Summary of Steps to Minimize Content-Irrelevant Test Material	36

# LIST OF FIGURES

Figure 1. Continuous Cycle of Item Development	16
--	----

# LIST OF TABLES

Table 1. Internal Consistency Reliability Estimates for Relevant iTEP Academic Scales	18
Table 2. Test-Retest Reliability Estimates for iTEP Academic	19
Table 3. Raw Rater Agreement Analysis – Speaking Scale	20
Table 4. Raw Rater Agreement Analysis – Writing Scal	21
Table 5. $r_{WG}$ Rater Agreement Statistics – Speaking Scale	22
Table 6. $r_{WG}$ Rater Agreement Statistics – Writing Scale	22
Table 7. iTEP Academic Scale Intercorrelations	24

# ABOUT THIS REPORT

## Purpose

The purpose of this report is to describe the technical details and rationale for the development of iTEP Academic and to summarize the reliability and validity of the assessment.

## Acknowledgements

The data analyzed for this report were provided in raw form by Boston Educational Services to Dr. Stephanie Seiler, Independent Consultant. The data were not manipulated in any way, unless the data manipulation was done according to accepted statistical procedures and/or was done according to other rationale as described in this report. Dr. Seiler conducted all analyses and wrote the report. Dr. Jay Verkuillen and Dr. Rob Stilson provided statistical consultation.

## About the Author

Stephanie Seiler holds a PhD in Industrial and Organizational Psychology from the University of Illinois at Urbana-Champaign. Stephanie worked in the personnel selection industry for 5 years as a Selection and Assessment Manager and then as the Director of Research and Development for a leading personnel selection company. In these roles, she was responsible for customer research, consulting, implementation of assessments, and development of new and innovative assessments. Stephanie has published and presented in the areas of educational and personnel assessment, innovative assessment methodologies, and research ethics. She is currently operating as an independent assessment consultant and is a Lead Research Associate and product developer with the University of Illinois at Urbana Champaign with the National Center for Professional and Research Ethics (NCPRE).

## About Boston Educational Services

Boston Education Services (BES), founded by career international educators, developed the iTEP suite of examinations to provide institutions and individual test-takers with an efficient, secure, accurate, and affordable on-demand language proficiency assessment.

There are currently seven versions of iTEP available: iTEP Academic, iTEP SLATE (Secondary Level Assessment Test of English), iTEP Business, iTEP Au Pair, iTEP Intern, iTEP Hospitality, and iTEP Conversation. All seven exams have the same basic structure, standardized rubric scoring, and administration procedures. The exams assess all or some of the following five components of English language proficiency: Grammar, Listening, Reading, Speaking, and/or Writing.

# CHAPTER 1

## INTRODUCTION

---

The International Test of English Proficiency - Academic (iTEP Academic), developed and published by Boston Educational Services (BES) is a multimedia assessment that evaluates the English language proficiency of English as a Second Language (ESL) college applicants and students.

iTEP Academic is commonly used for:

- Making admissions decisions
- Placing students within language programs
- Guiding course instruction and curriculum development
- Evaluating pre- and post-course progress
- Determining eligibility for scholarships

iTEP Academic is also used to assess the proficiency of English language teachers.

In order to target the level and type of English proficiency needed to be a successful college student, the content of iTEP Academic is tailored to reflect the academic and life experiences of individuals who attend or plan to attend college. iTEP Academic does not require any specialized academic or cultural knowledge, so it is well-suited for testing in any academic discipline. The assessment evaluates examinees' ability to apply their English knowledge and skill to process, learn from, and respond appropriately to new information that is presented in English. iTEP Academic is delivered over the Internet at secure Certified iTEP Test Centers around the world. Examinees can schedule a testing date within three business days of contacting the test center.

There are two versions of iTEP Academic:

- **iTEP Academic-Core:** assesses Grammar, Listening, and Reading and is 50 minutes in length, with an additional 10 minutes for pre-test preparation. Results are available immediately.
- **iTEP Academic-Plus:** assesses Grammar, Listening, Reading, Speaking, and Writing and is 80 minutes in length, with an additional 10 minutes for pre-test preparation. Results are available within 24 hours.

iTEP automatically emails the examinee's official score report to the client. An online iTEP client account provides a variety of tools for managing results.

## Approach and Rationale for the Development of iTEP Academic

College is a social and communicative experience. Whether the student is listening to a lecture, writing a paper, reading exam instructions, working on a group project, or making a purchase at the bookstore, the ability to understand and use the college's primary language is a fundamental prerequisite for the student to succeed. Though success in college can also depend on factors that have little direct link to language, such as intelligence, motivation, self-discipline, and physical and emotional health, these will have little use for the student if he/she is unable to process, learn from, and respond to information.

iTEP Academic was designed and developed to provide English language proficiency scores that are valid for many types of educational decision making. The developers of iTEP Academic recognized that in order to thoroughly evaluate English proficiency, the assessment needed to include items that evaluated both written and spoken language, as well as the examinee's grasp of English grammar. In addition, iTEP developers made the distinction between receptive language skills (i.e., listening and reading) and expressive language skills (i.e., writing and speaking). Assessment items that measure an examinee's ability to express ideas in English were developed for inclusion in iTEP Academic – Plus.

When language proficiency is measured accurately, reliably, and comprehensively, educators or administrators can use examinees' scores on the assessment to make more rigorous, evidence-based decisions. iTEP Academic was developed with these goals in mind. Furthermore, iTEP Academic uses the best technology available and on-demand support to help ensure an engaging, user-friendly examinee and administrator experience.

iTEP is recognized by the Academic Credentials Evaluation Institute (ACEI) and Accrediting Council for Continuing Education and Training (ACCET), as an approved internationally regarded English proficiency exam that meets institutional standards. In addition, BES is committed to actively engaging with the international education community through memberships and affiliations with NAFSA, EnglishUSA, TESOL, ACEI, ACCET, and AISAP.

## CHAPTER 2

# DETAILED DESCRIPTION OF iTEP ACADEMIC

---

### Theoretical Model for Language Assessment

Traditionally, language researchers and educators have grouped language skills into four distinct categories (Listening, Reading, Speaking, and Writing), and from a commonsense perspective this categorization is no surprise, as each of these elements of communication refers to a distinct set of activities and knowledge used for distinct purposes. In addition, it is common for a distinction to be made between language *skills* and language *knowledge* (e.g., grammar and vocabulary) (Bachman, 1990).

On the surface, the Listening, Reading, Speaking, and Writing sections of iTEP Academic align with the traditional categorization of language *skills*, and the Grammar section aligns with the notion of language *knowledge*. Listening scores reflect the ability to comprehend spoken language, Grammar scores reflect the knowledge of correct grammar, and so on. Additionally, practical considerations clearly warrant testing across multiple competency areas. In the case of admissions or certification, use of multiple measures helps ensure content coverage (measurement breadth) across the most critical elements of language; in the case of placement or program evaluation, multiple measures help pinpoint different areas of examinee strengths and weaknesses.

The traditional categorization of language into skills and knowledge domains may seem to suggest that each iTEP Academic scale measures an isolated language capability; however, modern theories of language emphasize the interrelatedness of language knowledge and skill and the practical fact that any attempt to measure a single component of language will likely be confounded by other language skills that are necessary to answer the question (for example, an evaluation of reading proficiency requires knowledge of grammar, sentence structure, vocabulary, etc.). In addition, these theories emphasize that one must consider the context in which the communication occurs; communication in a casual setting is likely to involve a different set of competencies—and a different judgment of effectiveness—than communication in an academic or business setting. These modern theories suggest that in practice, language effectiveness must be evaluated in the situational context for which the assessment is to be used (Association of Language Testers in Europe (ALTE), 2011; Bachman, 1990). Plainly stated, a language assessment should represent the real-world use of language. The component parts of language are still relevant to language assessment, but they must be interpreted in context.

iTEP Academic aligns with best practices in language assessment by evaluating one's ability to communicate effectively in the context of common scenarios that are encountered in college.

## Description of iTEP Academic Scales

### Grammar Section

The ability to understand and use a language's grammar rules correctly is an important component of effective communication. Grammar does not need to be perfect in order for someone to comprehend the meaning of a statement, yet as the number of grammatical errors increases, the likelihood that the information will be conveyed incorrectly also increases. Still higher standards for grammatical correctness are present within most academic settings.

The iTEP Academic: Grammar section evaluates an examinee's understanding of and ability to use proper English grammar. It is comprised of twenty-five multiple-choice questions, each of which tests the examinee's familiarity with a key feature of English structure (e.g., use of the correct article, verb tense, modifier, or conjunction; identifying the correct sentence structure, pronoun, or part of speech). The Grammar section includes a range of sentence structures from simple to more complex, as well as both beginning and advanced vocabulary. The first 13 questions require the examinee to select the word or phrase that correctly completes a sentence, and the next 12 questions require the examinee to identify the word or phrase in a sentence that is grammatically incorrect. Each of the two question types is preceded by an on-screen example.

The Grammar section takes 10 minutes to complete.

### Sample Grammar Item

**DIRECTIONS:** Click on the word or phrase that is NOT CORRECT in this sentence.

The dog was played outside, and now he is covered with mud.

- The
- played
- is
- covered

## Listening Section

The ability to comprehend spoken information is of central importance within an academic setting—as well as for navigating the social aspects of college life. The typical model for a college course, particularly during the first two years of coursework, involves students attending lectures. The iTEP Academic: Listening section evaluates an examinee’s proficiency in understanding spoken English information. In this section, the examinee listens to two types of spoken information: (1) a short conversation between two speakers; and (2) a brief lecture on an academic topic. After listening to the conversation or lecture, the examinee is presented with a question (orally and in writing) that measures several key indicators of whether the information was understood. These indicators include: identifying the primary subject of the conversation or lecture (Main Idea), recalling important points (Catching Details), understanding why a particular statement was made (Determining the Purpose), inferring information based on contextual information (Making Implications), and determining the relationship between key pieces of information (Connecting Content).

To ensure realism in the Listening section, item writers take steps to ensure that the content reflects a conversational tone. In addition, while the examinee listens to each audio file, a static image of the speaker(s) is presented onscreen.

The Listening section takes 20 minutes to complete and consists of three parts:

**Part 1:** Four high-beginning to low-intermediate-difficulty level conversations of 2-3 sentences, each followed by 1 multiple-choice question

**Part 2:** One 2- to 3-minute intermediate- difficulty level conversation followed by 4 multiple-choice questions

**Part 3:** One 4-minute lecture followed by 6 multiple-choice questions

## Sample Listening Item:



Transcript of audio played to examinee [text is for demonstration in this report and is not presented to the examinee]

Male Student

“Hi Tara. Did you hear that Professor Johnson’s biology class was cancelled? He moved the quiz to next week.”

Tara

“No, I didn’t. Thanks for telling me. That will give me more time to write my history report and finish my math homework.”

**In what class does Tara have a quiz?**

- Math.
- History.
- Writing.
- Biology.

## Reading Section

Along with Listening, the ability to comprehend written information is critical for effective learning in an academic setting—as well as for navigating college life in general. Course lectures are typically paired with required textbooks or other reading materials, and students are frequently evaluated on their recall and understanding of both the lectures and the readings. Additionally, the typical examination in lower-level college courses involves written materials such as multiple-choice questions.

The iTEP Academic: Reading section evaluates an examinee’s level of reading comprehension by measuring several key indicators of whether a written passage was understood. These indicators include: identifying the significant points and main focus of the written passage (Catching Details and Main Idea, respectively), determining what a word means based on its context (Vocabulary), and understanding why a particular statement within a larger passage was written by connecting together relevant information (Synthesis). In addition, the Reading section evaluates the examinee’s understanding of how a paragraph should be constructed in order to properly convey information (Sequencing). Sequencing items require the examinee to read a paragraph and determine where a new target sentence should be placed based on the surrounding content.

The Reading section takes 20 minutes to complete and consists of two parts:

**Part 1:** One intermediate reading level passage about 250 words in length, followed by 4 multiple-choice questions

**Part 2:** One upper reading level paragraph about 450 words in length, followed by 6 multiple-choice questions

### Sample Reading Item:

times their body length. If a human had similar ability, he could jump 90 meters.

▶ Not only can spiders jump far, but they can also walk upside down on smooth surfaces. Their feet are covered with tiny hairs that enable them to hold 170 times their body weight before coming unstuck. That is equivalent to a children's super-hero carrying 170 people from danger while clinging to the side of a building with his fingers and toes.

Spiders can also spin as many as seven different kinds of silk. Some of the silk is so strong that it rivals the strength of steel. Spiders use the silk for many different purposes, such as catching insects in webs and then wrapping them up so that they cannot escape. They also use silk to travel from place to place and to form egg sacs.

Spiders come in a wide variety of sizes. The largest known spider is the Goliath bird eater tarantula. This South American spider can be as big as a dinner plate. The smallest known spider is the mygalomorph spider from Borneo. Its body is the size of a pinhead.

**DIRECTIONS:** Answer the question below. To view the next question, click on the "Next" button.

According to Paragraph 2, what permits spiders to walk upside-down on smooth surfaces?

- The shape of their body
- A sticky substance that comes out of their feet
- Small hairs on their feet
- Toes that can grasp tiny irregularities in the surface

Note: an arrow [▶] points to Paragraph .

## Speaking Section

The speaking and writing in a new language are often considered more advanced skills, developed after the individual has acquired a basic grasp of the language’s grammar and vocabulary and learned to apply this knowledge to comprehend written and spoken information. The longer version of iTEP Academic, iTEP Academic – Plus, evaluates the examinee’s English Speaking ability (along with Writing ability as described next).

During the Speaking section of the assessment, the examinee listens to and reads a prompt (either a question or a brief lecture), and then prepares an oral response. The examinee then records his/her response for later evaluation by a trained iTEP rater.

The Speaking section takes 5 minutes to complete and consists of two parts:

**Part 1:** The examinee hears and reads a short question geared at low-intermediate level, then has 30 seconds to prepare a spoken response, and 45 seconds to speak.

**Part 2:** The examinee hears a brief upper-level statement presenting two sides of an issue, then is asked to express his or her thoughts on the topic, with 45 seconds to prepare, and 60 seconds to speak.

### Sample Speaking Item:

**DIRECTIONS:** You will both hear and read a question about school life. Answer the question giving specific reasons and examples that support your answer. After you hear the question, you will have 30 seconds to prepare your answer, and 45 seconds to speak.

**Topic:** After you complete your studies, what kind of work do you want to do?

<b>PREPARE</b>	<b>SPEAK</b>
29	45
SECONDS	SECONDS

## Writing Section

In addition to the Speaking section, iTEP Academic – Plus evaluates the examinee’s English Writing ability.

During the Writing section of the assessment, the examinee reads a question and then writes a response. The responses are submitted for later evaluation by a trained iTEP rater.

The Writing section takes 25 minutes to complete and consists of two parts:

**Part 1:** The examinee is given five minutes to write a 50-75 word note, geared at the low intermediate level, on a supplied topic

**Part 2:** The examinee is given 20 minutes to write a 175-225 word piece expressing and supporting his or her opinion on an upper-level written topic

### Sample Writing Item:

**Topic:** In some societies, families expect their children to leave home and live on their own as soon as they finish high school. In other societies, children live with their families until they get married – and sometimes even after they are married. Explain which approach you prefer, and why.

copy

cut

paste

undo

Word count: 100 left.

# Assessment Administration

## Delivery Method

iTEP Academic is administered via the Internet. Items are administered to examinees at random from a larger item bank, according to programming logic and test development procedures that ensure each examinee receives an overall examination of comparable content and difficulty to other examinees.

A static paper-and-pencil version of iTEP Academic is also available.

iTEP Academic must be administered at a secure location or a Certified iTEP Test Center.

The examinee inputs responses to the test in the following manner:

- During the Reading, Listening, and Grammar sections, the examinee selects from a list of multiple choice options for each question
- Writing samples are keyboarded directly into a text entry field
- Speaking samples are recorded with a headset and microphone at the examinee's computer

## Examinee Experience

Prior to the start of the test, the examinee logs in and completes a registration form. The system guides the examinee through a series of steps to ensure technical compatibility and to prepare him/her for the format of the assessment.

Each section/scale has a fixed time allotted to it. In the Reading and Grammar sections, examinees can advance to the next section if there is time remaining, or they are free to use any extra time to review and revise their answers. In the Listening section, the prompts each play only once and once submitted, an item response cannot be reviewed or changed. In the Writing section, there are fixed time limits for each part, but examinees may advance to the next section before time expires. In the Speaking section, there are fixed time limits for each part and examinees cannot advance until time expires.

The directions for each section are displayed for a set amount of time, and are also read aloud. The amount of time instructions are displayed varies according to the amount of text to be read. If an examinee needs more time to read a particular section's directions, he or she can access them by clicking the Help button, which displays a complete menu of directions for all test sections.

Following each section of the test, examinees see a transition screen indicating which section will be completed next. These transition screen provides a 15-second break between sections, and displays a progress bar showing completed and remaining test sections. After the last test section is completed, examinees see a final screen telling them to wait for further directions from the administrator.

Screenshots of the examinee experience, including pre-assessment modules and instructions, are shown in Appendix A.

## Scoring/Grading

iTEP Academic computes an overall proficiency level from 0 (Beginner) to 6 (Mastery), as well as individual proficiency levels from 0 to 6 for each scale. Sub-scale scores are also computed (e.g. parts of speech, synthesis, main idea), in order to give a more detailed picture of the examinee's skill level. The Overall score represents the combination of scores across each scale; for greater accuracy, Overall scores are reported to one decimal point (e.g., 0.0, 0.1, 0.2, ... , 5.9, 6.0).

iTEP Academic is graded as follows:

- The Grammar, Listening, and Reading scales are scored automatically by the computer. Each response is worth 1 point. There is no penalty for guessing.
- The Speaking and Writing scales are evaluated by native English-speaking, ESL-trained professionals, according to a standardized scoring rubric (see Appendix B and Appendix C). Raters attend refresher training sessions throughout the year to ensure continued adherence to the rubric.
- For computing the Overall score, each test scale is weighed equally.
- The official score report presents an individual's scoring information in both tabular and graphical formats. The graphical format, or skill profile, is particularly useful for displaying an examinee's strengths and weaknesses in each of the skills evaluated.

## Proficiency Levels

The seven iTEP Academic proficiency levels may be expressed briefly as follows:

**Level 0:** Beginning

**Level 1:** Elementary

**Level 2:** Low Intermediate

**Level 3:** Intermediate

**Level 4:** High Intermediate

**Level 5:** Low Advanced

**Level 6:** Advanced

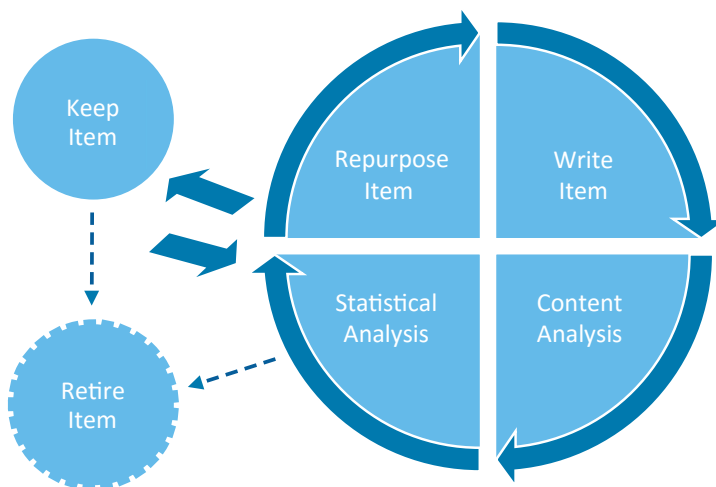
BES has mapped iTEP Academic Proficiency Levels to the levels described in the Common European Framework of Reference for Languages (CEFR; See Appendix D).

# CHAPTER 3

## ITEP ACADEMIC DEVELOPMENT PROCESS, RELIABILITY, AND VALIDITY<sup>1</sup>

### Development Process

Boston Educational Services adheres to a continuous cycle of item analysis (see Figure 1) to ensure the content of the assessment adheres to the reliability and validity goals of the assessment. The cycle begins with item writing, enters an expert review and content analysis stage, and then works through a number of statistical analyses to evaluate the difficulty level and other psychometric properties of the item. Items that do not meet quality standards during the content analysis and/or statistical analysis phase are either removed from further consideration, or repurposed if it is determined that minor adjustments will improve the item. Items that meet quality standards during the content analysis and statistical analysis phases are retained in the assessment; in order to maintain a secure assessment and minimize the likelihood of an item being shared among examinees over time, all items used in the assessment are retired after a certain length of time. Items may also be identified as having “drifted” in difficulty over time, indicating that the item may have been compromised; these items are retired immediately upon identification.



**Figure 1. Continuous Cycle of Item Development**

1 All analysis and evaluation of iTEP Academic as described in Chapter 3 was conducted in accordance with the Standards for Educational and Psychological Testing (hereafter *Standards*; American Educational Research Association, American Psychological Association, & National Council on Measurement in Education, 2014), Uniform Guidelines on Employee Selection Procedures (Equal Employment Opportunity Commission, 1978), and the Principles for the Validation and Use of Personnel Selection Procedures (Society for Industrial and Organizational Psychology, 2003)

## Reliability

The reliability of an assessment refers to the degree to which the assessment provides stable, consistent information about an examinee. Demonstrating reliability is important because if a test is not stable and consistent—whether across the items in the assessment, across repeated administrations of the assessment, or based on performance scores provided by trained raters—then the results cannot be relied upon as accurate. Moreover, the reliability of an assessment theoretically sets a maximum limit for its validity; when an assessment is not consistent, it is less effective as an indicator of a person’s true ability and will therefore demonstrate lower correlations with relevant outcomes (such as grades, academic adjustment, or attrition).

### Internal Consistency Reliability

Internal consistency reliability refers to the stability of the items within a particular assessment, or in this case, within each assessment scale. When it can be shown that the items are statistically related to each other, the case can be made that the assessment is consistent in its measurement. Cronbach’s alpha (Cronbach, 1951) is a commonly-used and accepted classical test theory (CTT) statistic that is used to estimate internal consistency reliability. The statistic reflects the average correlation between all items within an assessment or assessment section. Values of .70 or above have traditionally been considered desirable, with some scholars stating that test developers should aim to develop tests with values of at least .80 or even .90 and higher. These benchmarks are general rules and do not take into account other desirable characteristics of an assessment, such as assessment brevity to minimize testing time (Gatewood & Field, 2001), the breadth of content coverage within the assessment (to ensure a large domain of the characteristic being measured is represented) (Loevinger, 1954), and the validity of the assessment (Nunnally & Bernstein, 1994). Test developers must think critically about the interrelated factors influencing test reliability and validity and use their best judgment when deciding what should be considered acceptable (Gatewood & Field, 2001).

Because the calculation of internal consistency reliability requires that the assessment scale contain multiple items, this class of statistics is appropriate for the Grammar, Listening, and Reading scales of iTEP Academic; calculation of internal consistency reliability is not possible for the Speaking and Writing scales, as trained raters provide only one summary score for each of these sections based on the examinee’s overall Speaking or Writing performance.

Within the Grammar, Listening, and Reading sections of iTEP Academic, the set of items administered to each examinee are selected at random from a larger item bank; therefore, the traditional CTT calculation of Cronbach’s alpha it is not possible. In order to compute an internal consistency reliability estimate for each scale, the following procedure was used to derive an estimate that can be interpreted in a manner similar to Cronbach’s alpha. The procedure relies on statistics derived from item response theory (IRT), a class of statistical models that are particularly suited to handling randomly-administered items.

- 1 For each scale, compute the IRT common discrimination parameter using the 1-Parameter Logistic model (1PL). The common  $a$  parameter reflects the average extent to which each item provides statistical information that distinguishes lower-performing examinees from higher-performing examinees. The  $a$  parameter is in concept most similar to an item-total correlation from classical test theory.
- 2 Use the  $a$  parameter estimate to compute an intraclass correlation coefficient (ICC). This formula is:

$$ICC = \frac{a^2}{a^2 + \pi^2/3}$$

- 3 The resulting value of the ICC reflects the average internal consistency reliability for any one item in the scale, and therefore the final internal reliability estimate ( $\alpha$ ) must be “stepped up” using the Spearman-Brown prophecy formula to reflect the reliability of the total scale. The Spearman-Brown prophecy method is the same method that would be used to examine the impact of shortening or lengthening a test (for example, cutting a 50-item test in half). The Spearman-Brown prophecy formula is:

$$\alpha = \frac{K \cdot ICC}{1 + (K - 1) \cdot K \cdot ICC}$$

Where K is a scaling factor reflecting the proportional increase or decrease in the number of test items. In the current case, K is the number of items in the scale.

The internal consistency reliability results, which can be interpreted as conceptually similar to Cronbach’s alpha estimates, were computed for a sample of over 17,000 examinees who completed iTEP Academic between 2014 and 2016.<sup>2</sup> The results are provided in Table 1. As shown, all values exceed the .70 benchmark, and the Grammar estimate exceeds the .80 benchmark.

**Table 1. Internal Consistency Reliability Estimates for Relevant iTEP Academic Scales**

Scale	Number of Items	Discrimination ( $a$ )	Intraclass Correlation (ICC)	Internal Consistency Reliability ( $\alpha$ )
Grammar	25	1.08	.25	.89
Listening	14	0.85	.21	.78
Reading	10	0.93	.22	.74

Note: The sample size for the analysis was N = 17,731. The Internal consistency reliability estimates are not Cronbach’s alpha values, but can be interpreted in a similar manner to Cronbach’s alpha.

- 2 All examinee data provided by BES was included in the analysis, with the exception of the following: (1) when a unique identifier indicated the data was for an examinee re-testing, only the examinee’s first testing occasion was included; or (2) if the examinee timed-out on any scale without seeing one or more of the items, the examinee was removed; or (3) examinees younger than 14 years of age were removed. Examinee non-responses to items that were seen but not answered were scored as incorrect.

## Test-Retest Reliability

Test-retest reliability refers to the stability of test scores across repeated administrations of the test. A high level of test-retest reliability indicates that the examinee is likely to receive a similar score every time he or she takes it—assuming the examinee’s actual skill in the domain being measured has not changed. Test-retest reliability estimates for all iTEP Academic scales, and the Overall score, were computed using a sample of 198 examinees who took iTEP Academic twice in an operational environment (i.e., at a testing center for college admissions purposes). Analyses were restricted to examinees with at least 5 days and less than 2 months between testing occasions (average time elapsed: 24 days).

The test-retest values shown in Table 2 reflect the correlation between the Time 1 and Time 2 scores for the sample. Values can range from -1.0 to 1.0, with values at or exceeding .70 typically considered desirable. As can be seen, only the Overall score exceeds this threshold. However, it should be noted that the sample used to compute the test-retest correlations was an operational sample, and it could reasonably be assumed that at least some of the sample had worked diligently to improve their performance between Time 1 and Time 2 testing occasions; given the number of days between test administrations for the sample (up to 2 months; 24 days on average), this seems very likely. Had the test-retest estimates been computed on a research sample and/or if the sample size of available data allowed for the analysis of a shorter time period between testing occasions, the correlations would likely be higher. Therefore, the values given in Table 2 can be considered lower-bound estimates of the true test-retest reliability of iTEP Academic.

**Table 2. Test-Retest Reliability Estimates for iTEP Academic**

Scale	Test-Retest Reliability
Grammar	.63
Listening	.49
Reading	.48
Speaking	.64
Writing	.62
Overall	.77
Overall – Core	.71

Note: The sample size for the analysis was N = 198. The OVERALL – Core score was approximated by removing the Speaking and Writing section scores from the Overall scores of examinees who completed the longer iTEP Academic – Plus.

## Rater Agreement

The iTEP Academic Speaking and Writing sections are evaluated by a trained rater and as such, it is necessary to estimate the accuracy of these judgments—specifically, the extent to which the scores given by rater are interchangeable with the scores of another. Evaluations of rater *agreement*, as opposed to rater *reliability*, are more appropriate in cases where the examinee’s absolute score is of interest rather than the examinee’s rank order position relative to other examinees (LeBreton & Senter, 2008).

Tables 3 and 4 summarize a raw investigation of rater agreement using a sample of Speaking and Writing ratings from six examinees obtained from eight raters during a training exercise. The examinees completed either iTEP Academic or iTEP SLATE.

It should be noted that the results in Tables 3-6 likely reflect a lower-bound estimate of rater agreement, as the cases used for the training exercise were purposely selected to be more challenging to rate than a typical case.

**Table 3. Raw Rater Agreement Analysis – Speaking Scale**

Examinee	Average Score	Rater Deviations from Average Score								Average Deviation	Max. Deviation
		R1	R2	R3	R4	R5	R6	R7	R8		
E1	1.79	.04	.54	.29	.04	.29	.46	.71	-	.34	.71
E2	4.88	.63	.38	.63	.38	.13	.38	.63	.88	.50	.88
E3	2.53	.28	.97	.03	.22	1.53	.22	.72	.28	.53	1.53
E5	4.03	.22	.72	.53	.53	.22	.03	.53	.47	.41	.72
E6	3.50	.50	.00	1.50	-	-	-	.25	.75	.60	1.50
<b>Average</b>	<b>3.34</b>	<b>.33</b>	<b>.52</b>	<b>.59</b>	<b>.29</b>	<b>.54</b>	<b>.27</b>	<b>.57</b>	<b>.59</b>	<b>.47</b>	<b>1.07</b>

Note: No Speaking scale ratings were provided for Examinee 4 due to a technical issue with the audio recording. The missing values occurred because the rater(s) did not provide a rating. Average Score: the examinee’s average rating across all eight raters. Rater Deviations from Average Score: the absolute value of the difference between each rater’s score and the Average Score for each examinee. Average Deviation: average Rater Deviation for each examinee. Max Deviation: highest Rater Deviation value that was observed across all eight raters.

**Table 4. Raw Rater Agreement Analysis – Writing Scale**

Examinee	Average Score	Rater Deviations from Average Score								Average Deviation	Max. Deviation
		R1	R2	R3	R4	R5	R6	R7	R8		
E1	1.79	.04	.71	.29	.79	.04	.29	.71	-	.41	.79
E2	3.81	.56	.69	.06	.31	.44	.56	.06	.44	.39	.69
E4	3.43	.32	-	.07	.18	.18	.32	.18	.18	.20	.32
E5	4.41	.16	.66	.59	.09	.09	.09	.16	.09	.24	.66
E6	3.60	.15	.40	.85	-	-	-	.15	.15	.34	.85
<b>Average</b>	<b>3.41</b>	<b>.25</b>	<b>.61</b>	<b>.37</b>	<b>.34</b>	<b>.19</b>	<b>.32</b>	<b>.25</b>	<b>.21</b>	<b>.32</b>	<b>.66</b>

Note: No Writing scale ratings were provided for Examinee 3. The missing values occurred because the rater(s) did not provide a rating. Average Score: the examinee’s average rating across all eight raters. Rater Deviations from Average Score: the absolute value of the difference between each rater’s score and the Average Score for each examinee. Average Deviation: average Rater Deviation for each examinee. Max Deviation: highest Rater Deviation value that was observed across all eight raters.

As seen in Table 3, in all but 2 instances the raters’ Speaking scores for each examinee deviated less than 1 point from the average rating across all raters (as a reminder, scale scores can range from 0 to 6). Across all raters and examinees, the average deviation was .47 points, and the average maximum deviation was 1.07 points. These results suggest a moderately strong agreement across raters.

As seen in Table 4, all of the raters’ Writing scores for each examinee deviated less than 1 point from the average rating across all raters. Across all raters and examinees, the average deviation was .32 points, and the average maximum deviation was .66 points. These results suggest a strong agreement across raters.

Using the same data that were used for Tables 3 and 4, rater agreement was also estimated using a version of the  $r_{WG}$  agreement statistic (James, Demaree, and Wolf, 1984). The value of  $r_{WG}$  can theoretically range from 0 to 1, and represents the observed variability in scores among raters relative to the amount of variability that would be present if all raters had assigned scores completely at random. The formula for  $r_{WG}$  is:

$$r_{WG} = 1 - \frac{S^2_x}{\sigma^2_E}$$

Where  $S^2_x$  is the observed variance of ratings on the variable across raters and  $\sigma^2_E$  is the variance expected if the ratings were completely random.

The specific version of  $r_{WG}$  chosen for the analysis uses a value for  $\sigma^2_E$  that would occur if the raters’ completely random scores came from a triangular (approximation of normal) distribution (see LeBreton & Senter, 2008).

The closer an  $r_{WG}$  value is to 1, the higher the agreement. There is no agreed-upon minimum value that is considered acceptable for  $r_{WG}$ , but as a benchmark, test developers might consider .80 or .90 to be a minimally acceptable value for an application such as assigning ratings based on a score rubric. To put these values in perspective, an  $r_{WG}$  of .80 would suggest that 20% (1 – .80) of an average rater’s score across examinees was due to error, or factors other than the examinee’s “true score” on the exercise.

The  $r_{WG}$  agreement statistics are presented in Tables 5 and 6.

**Table 5.  $r_{WG}$  Rater Agreement Statistics – Speaking Scale**

Examinee	Observed Variance	Error Variance	$r_{WG}$
E1	.20	2.1	.91
E2	.34	2.1	.84
E3	.58	2.1	.72
E5	.24	2.1	.89
E6	.78	2.1	.63
<b>Average</b>			<b>.80</b>

Note: No Speaking scale ratings were provided for Examinee 4 due to a technical issue with the audio recording.

**Table 6.  $r_{WG}$  Rater Agreement Statistics – Writing Scale**

Examinee	Observed Variance	Error Variance	$r_{WG}$
E1	.30	2.1	.86
E2	.23	2.1	.89
E4	.06	2.1	.97
E5	.12	2.1	.94
E6	.24	2.1	.89
<b>Average</b>			<b>.91</b>

Note: No Writing scale ratings were provided for Examinee 3.

The results in Table 5 indicate moderately strong agreement amongst the raters. The minimum  $r_{WG}$  was observed for Examinee 6, with a value of .63. The average  $r_{WG}$  across all examinees was .80, indicating that 20% of the average rater’s score across examinees was due to factors other than the examinee’s “true score” on the exercise.

The results in Table 6 indicate strong agreement amongst the raters. The minimum  $r_{WG}$  was observed for Examinee 1, with a value of .86. The average  $r_{WG}$  across all examinees was .91,

indicating that only 9% of the average rater’s score across examinees was due to factors other than the examinee’s “true score” on the exercise.

Overall, the results of the rater agreement analyses suggest that ratings provided by any one iTEP rater are likely to be a reliable indication of an examinee’s actual proficiency on the Speaking and Writing scales.

## Validity

The iTEP Academic examination was designed and developed to provide English language proficiency scores that are valid for many types of educational decision making. The *Standards* define validity as “the degree to which accumulated evidence and theory support specific interpretations of test scores entailed by proposed uses of a test” (AERA et al., 2014, p. 184). In other words, the term validity refers to the extent to which an assessment measures what it is intended to measure. Evidence for validity can, and should, come from multiple lines of investigation that together converge to form a conclusion regarding the relative validity of the assessment, including:

- 1 Expert judgments regarding the extent to which the content of the assessment reflects the real-world knowledge, skills, characteristics, or behaviors the assessment is designed to measure (Content Validity)
- 2 An examination of the degree to which the assessment (or assessment scale) is correlated with theoretically similar measures and un-correlated with theoretically unrelated measures (Convergent and Discriminant Validity; traditionally conceived of as the main contributors to Construct Validity<sup>3</sup>)
- 3 An examination of the degree to which the assessment is correlated with the real-world outcomes it is intended to measure, for example: adjustment to college, grades, or improvement in language proficiency (Criterion Validity)

## Content Validity

Content validity, or content validation, refers to the process of obtaining expert judgments on the extent to which the content of the assessment corresponds to the real-world knowledge, skill, or behavior the assessment is intended to measure. For example, an assessment that asks questions about an examinee’s knowledge of cooking techniques may be judged by experts to be content valid for measuring that aspect of cooking skill, but it would not be content valid for measuring the examinee’s athletic ability—even if it turned out that cooking assessment scores were correlated with athletic ability.

According to the *Standards* (AERA et al., 2014), evidence for assessment validity based on test content can be both logical and empirical and can include scrutiny of both the items/prompts themselves as well as the assessment’s delivery method(s) and scoring.

<sup>3</sup> The modern conception of *construct validity* refers not just to convergent and discriminant validity, but to the accumulation of all forms of evidence in support of an assessment’s validity (AERA et al., 2014).

Content-related validity evidence for iTEP Academic, for the purposes of academic decision-making, can be demonstrated via a correspondence between the assessment’s content and relevant college educational and social experiences. To ensure correspondence, developers conducted a comprehensive curriculum review and met with educational experts to determine common educational goals and the knowledge and skills emphasized in curricula across the country. This information guided all phases of the design and development of iTEP Academic.

Content validity evidence for iTEP Academic is also demonstrated through the use of trained item writers who are experts in the field of education and language assessment and who have substantial experience in item-writing. The content and quality of items submitted by item-writers is continually supervised, and feedback is provided in order to ensure ongoing adherence to the content goals of the assessment and to avoid content-irrelevant test material. Some of the critical steps taken to achieve this objective are summarized in Appendix E.

Finally, content validity evidence for iTEP Academic is shown via its correspondence with the CEFR framework. BES mapped iTEP Academic to the CEFR framework through a process of expert evaluation and judgment on the content of the assessment and associated scores.

## Convergent and Discriminant Validity

Convergent and discriminant validity evidence is demonstrated through a pattern of high correlations among scales that measure concepts that are known to be closely related, and lower correlations among scales measuring unrelated concepts (AERA et al., 2014). The intercorrelations among iTEP Academic scales are shown in Table 3. The examinee data analyzed are the same as described in the Reliability section.

**Table 7. iTEP Academic Scale Intercorrelations**

Scale	Listening	Reading	Speaking	Writing	Overall
Grammar	.59	.57	.57	.66	.83
Listening	–	.55	.54	.57	.80
Reading	–	–	.50	.56	.79
Speaking	–	–	–	.82	.82
Writing	–	–	–	–	.86

Note: N = 16,425 for correlations involving Speaking or Writing; N = 17,760 for all other correlations.

The pattern of correlations within iTEP Academic provides preliminary evidence for the convergent and discriminant validity of the assessment. Overall, the relatively strong correlations between the majority of scales (i.e., in the .50-.60 range) indicates that each scale is likely measuring related components of language proficiency, and the fact that the correlations do not approach 1.0 indicates that each scale likely measures a distinct element of proficiency. Compared with the Grammar/Speaking correlation, the higher correlation between Grammar

and Writing is conceptually logical given more weight is placed on grammar, by design, when iTEP raters evaluate examinees' writing ability than when evaluating their spoken ability. The strong correlation between Speaking and Writing is also to be expected, given these skills are considered more advanced demonstrations of language proficiency that require expressive, as opposed to receptive, language skills.

In addition to the internal examination of convergent and discriminant validity within the iTEP Academic scales, preliminary analyses conducted by a BES partner suggested a .93 correlation between iTEP scores and TOEFL® scores. The correlation indicates that iTEP scores are closely aligned with those of other language proficiency tests.

# REFERENCES

- American Educational Research Association, American Psychological Association, & National Council on Measurement in Education (1999). *Standards for educational and psychological testing*. Washington, DC: American Psychological Association.
- Association of Language Testers in Europe (2011). *Manual for language test development and examining*. Cambridge: ALTE.
- Bachman, L. (1990). *Fundamental considerations in language testing*. Oxford: Oxford.
- Cronbach, L. J. (1951). Coefficient alpha and the internal structure of tests. *Psychometrika*, *16*, 297-334.
- Equal Employment Opportunity Commission, Civil Service Commission, Department of Labor, and Department of Justice. (1978). *Uniform guidelines on employee selection procedures*.
- Gatewood, R.D. & Field, H.S. (2001). *Human resource selection* (5th ed.). Ohio: South-Westin.
- James, L. R., Demaree, R. G., & Wolf, G. (1984). Estimating within-group interrater reliability with and without response bias. *Journal of Applied Psychology*, *69*, 85-98.
- LeBreton, J. M., & Senter, J. L. (2008). Answers to 20 questions about interrater reliability and interrater agreement. *Organizational Research Methods*, *11*, 815-852.
- Loevinger, J. (1954). The attenuation paradox in test theory. *Psychological Bulletin*, *51*, 493—504.
- Nunnally, J. C., & Bernstein, I. H. (1994). *Psychometric theory* (3rd ed.). New York, NY: McGraw-Hill, Inc.
- Society for Industrial and Organizational Psychology, Inc. (2003). *Principles for the validation and use of personnel selection procedures* (4th ed.). Bowling Green, OH: SIOP.

# APPENDIX A: EXAMINEE PRE-ASSESSMENT MODULES AND INSTRUCTIONS

**\* = Required Information (scroll down to see all items)**

* Last (Family or Surname) Name:	<input type="text"/>	* Country of Residence:	<input type="text"/>
* First Name:	<input type="text"/>	* Nationality:	<input type="text"/>
Middle Name (Optional):	<input type="text"/>	* Government Identification Number:	<input type="text"/>
* Date of Birth:	Month <input type="text"/> Day <input type="text"/> Year <input type="text"/>	* Country Issuing Identification:	<input type="text"/>
* E-mail Address:	<input type="text"/>	* Type of Identification:	<input type="text"/>
* Phone (with city/country code):	<input type="text"/>	* Native Language:	<input type="text"/>
* Gender:	<input type="radio"/> Male <input type="radio"/> Female	* Highest Level of Education Attained:	Please Select One <input type="text"/>
		Field of Study:	Please Select One <input type="text"/>

School Level (if used):

Referral:

\* Have you ever taken the iTEP Test?  Yes  No

\* Are you applying to a school?  Yes  No

Where will you take the test?

Please read the following terms and conditions for taking this test:

1. Candidate's government-issued photo ID is required and will be verified before beginning the test.
2. The iTEP Administrator will verify that all information provided on the Registration Form is identical to the Candidate's official ID document(s).
3. Reference materials/tools and other personal effects (e.g. dictionaries, mobile phones, audio recording devices, pagers, notepaper, etc.) are not permitted in the room during the test.
4. Smoking, eating, or drinking is not permitted in the room during the test.
5. The iTEP Administrator reserves the right to dismiss a Candidate from the test or declare a Candidate's test results void if the Candidate violates any of the above conditions or fails to follow the Administrator's instructions during the test.
6. If for technical or any other reasons a given test is not able to be completed or results cannot be provided, Boston Educational Services' and the iTEP Administrator's liability shall be limited to providing a refund of fees received for said test and, at the Candidate's request, rescheduling a replacement test.

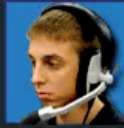
Note: Webcam will be used to save images during this test. Do not cover webcam during test.

\* I pledge on my honor that I will not give or receive any unauthorized assistance on this test. I am also aware that the penalty for cheating is severe and may include disqualification from any academic program.

Type "I Agree" into the text box:

\*  By checking this box, I agree to the above terms and conditions.

Please put your headphones on now. Then, click on the "Next" button.



Whenever you see this picture, you must have your headphones on.

While listening to these directions, click on the "Volume Slider Bar" found below. Using this "Volume Slider Bar", you can increase (right) or decrease (left) the volume to a comfortable level.

These directions will be repeated as you make your adjustments. If you have any sound problems that cannot be corrected by adjusting the volume, please ask your test proctor to help you.

When you are satisfied with the volume level, please click on "Next".

Note: You will also be able to adjust the volume during the test whenever you see the "Volume Slider Bar".



< Back

Next >

Please test your computer's voice recorder:

**STEP 1**

Click on the "Record" button below, then speak into your microphone for 6 seconds:



Say 2 times: "This is an English test."

**STEP 2**

Listen to the recording as it automatically plays.

**STEP 3**

Did you hear your voice clearly?  Yes |  No

You can always review the directions later by clicking the "Help" button. Now, click "Next" to continue.

Help

When you are answering the test questions, the numbers at the lower left of your screen will tell you:

**1/25**    **14:45**  
QUESTION    TIME LEFT

- A. The number of the current question.
- B. The total number of questions in the section.
- C. The amount of time remaining in the section.

At the beginning of each test section, you will have a limited amount of time to read the directions. This time will be displayed for you both visually and in numbers of seconds remaining, as in this example:

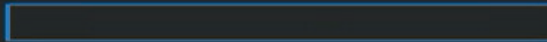
○○○○○○○○○ **56** seconds

The test has 5 sections proceeding in the following order:

- Grammar
- Listening
- Reading
- Writing
- Speaking

The test will take approximately 85 minutes to complete.

Test is loading ... please wait



# APPENDIX B: SPEAKING SCALE RATER SCORING RUBRIC

Rating criteria for each Level are discussed in terms of the following:

- General statement of ability and control of language, fluency
- Syntax and grammar
- Lexicon: sophistication of vocabulary
- Degree of elaboration in content, cultural/stylistic appropriateness
- Intelligibility and required listener/reader effort (includes mechanics such as spelling, punctuation, and capitalization)

Level	Rating Criteria
<b>6</b> <b>ADVANCED</b>	<ul style="list-style-type: none"> <li>• Highly effective control of the language; high degree of fluency; pauses and self-correction highly similar to those of native speakers</li> <li>• High degree of syntactic variety and sophistication; rare minor errors in grammatical usage</li> <li>• High degree of variety and sophistication in vocabulary, including idiomatic expressions; rare minor errors in word usage</li> <li>• Content elaboration is highly detailed and relevant to the task; high degree of cultural and stylistic appropriateness; high degree of organizational markers and coherence</li> <li>• Requires only rare effort by listener to determine intended meaning; pronunciation and intonation are highly intelligible with slight non-native influence</li> </ul>
<b>5</b> <b>LOW ADVANCED</b>	<ul style="list-style-type: none"> <li>• Mostly effective control of the language; fairly strong fluency; pauses and self-correction mostly similar to those of native speakers; occasional hesitation and false-starts</li> <li>• Fairly strong degree of syntactic variety and sophistication; occasional minor errors and awkwardness in grammatical usage, especially in complex structures</li> <li>• Fairly strong degree of variety and sophistication in vocabulary, including occasional idiomatic expressions; occasional minor errors in word usage</li> <li>• Content elaboration is detailed and relevant to the task; fairly high degree of cultural and stylistic appropriateness; fairly high degree of coherence and organizational markers</li> <li>• Requires occasional effort by listener to determine intended meaning; pronunciation and intonation are mostly intelligible with some non-native influence</li> </ul>

Level	Rating Criteria
<b>4</b> <b>HIGH INTERMEDIATE</b>	<ul style="list-style-type: none"> <li>Fairly effective control of the language; adequate fluency; some hesitation and false-starts</li> <li>Some syntactic variety and sophistication; fairly frequent significant errors and awkwardness in grammatical usage, especially in complex structures</li> <li>Fair variety and sophistication in vocabulary; rare use of idiomatic expressions; fairly frequent errors in word usage</li> <li>Content elaboration is somewhat detailed and mostly relevant to the task; some cultural and stylistic appropriateness; some degree of organizational markers and coherence</li> <li>Requires fair degree of effort by listener to determine intended meaning; pronunciation and intonation are fairly intelligible with moderate non-native influence</li> </ul>
<b>3</b> <b>INTERMEDIATE</b>	<ul style="list-style-type: none"> <li>Emerging control of the language; some degree of fluency; frequent hesitation and false-starts</li> <li>Some syntactic variety and sophistication; fairly frequent errors and awkwardness in grammatical usage, even in simple structures</li> <li>Attempts at variety and sophistication in vocabulary; rare use of idiomatic expressions; frequent errors in word usage</li> <li>Content elaboration is minimally detailed and fairly relevant to the task; occasional cultural and stylistic appropriateness; attempts at organizational markers and coherence</li> <li>Requires significant degree of effort by listener to determine intended meaning; pronunciation and intonation are somewhat intelligible with considerable non-native influence</li> </ul>
<b>2</b> <b>LOW INTERMEDIATE</b>	<ul style="list-style-type: none"> <li>Weak control of the language; little fluency; considerable hesitation and false-starts</li> <li>Little syntactic variety and sophistication; significantly frequent errors and awkwardness in grammatical usage, even in simple structures</li> <li>Little variety and sophistication in vocabulary; little use of idiomatic expressions; significantly frequent errors in word usage</li> <li>Content elaboration is very minimally detailed and parts may be irrelevant to the task; little cultural and stylistic appropriateness; few attempts at organizational markers and coherence</li> <li>Requires sustained effort by listener to determine intended meaning; pronunciation and intonation are markedly non-native</li> </ul>
<b>1</b> <b>ELEMENTARY</b>	<ul style="list-style-type: none"> <li>Very little control of the language; no fluency; intended meaning is mostly obscured; significant hesitation and false-starts</li> <li>Very limited syntactic and grammatical skills</li> <li>Very limited vocabulary</li> <li>Content elaboration is neither detailed nor culturally appropriate</li> <li>Requires extreme effort by reader to determine intended meaning; pronunciation and intonation are significantly non-native</li> </ul>
<b>0</b> <b>BEGINNING</b>	<ul style="list-style-type: none"> <li>No response or able to respond only with a few prompt-related words or reading of the prompt</li> <li>Mostly unintelligible; pronunciation and intonation are extremely non-native; extreme hesitation and false-starts</li> <li>May be off-topic or canned</li> </ul>

Note: Examinee scores can be between two levels, e.g., a score of 4.5 indicates the examinee is between levels 4 and 5.

# APPENDIX C: WRITING SCALE RATER SCORING RUBRIC

Rating criteria for each Level are discussed in terms of the following:

- General statement of ability and control of language, fluency
- Syntax and grammar
- Lexicon: sophistication of vocabulary
- Degree of elaboration in content, cultural/stylistic appropriateness
- Intelligibility and required listener/reader effort (includes mechanics such as spelling, punctuation, and capitalization)

Level	Rating Criteria
<b>6</b> ADVANCED	<ul style="list-style-type: none"> <li>• Highly effective control of the language; high degree of fluency</li> <li>• High degree of syntactic variety and sophistication; rare minor errors in grammatical usage</li> <li>• High degree of variety and sophistication in vocabulary, including idiomatic expressions; rare minor errors in word usage</li> <li>• Content elaboration is highly detailed and relevant to the task; high degree of cultural and stylistic appropriateness; high degree of organizational markers and coherence</li> <li>• Requires only rare effort by reader to determine intended meaning</li> </ul>
<b>5</b> LOW ADVANCED	<ul style="list-style-type: none"> <li>• Mostly effective control of the language; fairly strong fluency</li> <li>• Fairly strong degree of syntactic variety and sophistication; occasional minor errors and awkwardness in grammatical usage, especially in complex structures</li> <li>• Fairly strong degree of variety and sophistication in vocabulary, including occasional idiomatic expressions; occasional minor errors in word usage</li> <li>• Content elaboration is detailed and relevant to the task; fairly high degree of cultural and stylistic appropriateness; fairly high degree of coherence and organizational markers</li> <li>• Requires occasional effort by reader to determine intended meaning</li> </ul>
<b>4</b> HIGH INTERMEDIATE	<ul style="list-style-type: none"> <li>• Fairly effective control of the language; adequate fluency</li> <li>• Some syntactic variety and sophistication; fairly frequent significant errors and awkwardness in grammatical usage, especially in complex structures</li> <li>• Fair variety and sophistication in vocabulary; rare use of idiomatic expressions; fairly frequent errors in word usage</li> <li>• Content elaboration is somewhat detailed and mostly relevant to the task; some cultural and stylistic appropriateness; some degree of organizational markers and coherence</li> <li>• Requires fair degree of effort by reader to determine intended meaning</li> </ul>

Level	Rating Criteria
<b>3</b> <b>INTERMEDIATE</b>	<ul style="list-style-type: none"> <li>• Emerging control of the language; some degree of fluency</li> <li>• Some syntactic variety and sophistication; fairly frequent errors and awkwardness in grammatical usage, even in simple structures</li> <li>• Attempts at variety and sophistication in vocabulary; rare use of idiomatic expressions; frequent errors in word usage</li> <li>• Content elaboration is minimally detailed and fairly relevant to the task; occasional cultural and stylistic appropriateness; attempts at organizational markers and coherence</li> <li>• Requires significant degree of effort by reader to determine intended meaning</li> </ul>
<b>2</b> <b>LOW INTERMEDIATE</b>	<ul style="list-style-type: none"> <li>• Weak control of the language; little fluency</li> <li>• Little syntactic variety and sophistication; significantly frequent errors and awkwardness in grammatical usage, even in simple structures</li> <li>• Little variety and sophistication in vocabulary; little use of idiomatic expressions; significantly frequent errors in word usage</li> <li>• Content elaboration is very minimally detailed and parts may be irrelevant to the task; little cultural and stylistic appropriateness; few attempts at organizational markers and coherence</li> <li>• Requires sustained effort by reader to determine intended meaning</li> </ul>
<b>1</b> <b>ELEMENTARY</b>	<ul style="list-style-type: none"> <li>• Very little control of the language; rare fluency; intended meaning is mostly obscured</li> <li>• Very limited syntactic and grammatical skills</li> <li>• Very limited vocabulary</li> <li>• Content elaboration is neither detailed nor culturally appropriate</li> <li>• Requires extreme effort by reader to determine intended meaning</li> </ul>
<b>0</b> <b>BEGINNING</b>	<ul style="list-style-type: none"> <li>• No response or able to respond only with a few prompt-related words or reading of the prompt</li> <li>• Mostly unintelligible</li> <li>• May be off-topic or canned</li> </ul>

Note: Examinee scores can be between two levels, e.g., a score of 4.5 indicates the examinee is between levels 4 and 5.

# APPENDIX D: ITEP ABILITY GUIDE

ITEP	CEFR	Listening	Reading	Writing	Speaking
6	<b>C2</b> MASTERY	<ul style="list-style-type: none"> <li>Comprehends overall meaning and virtually all details of lectures on diverse topics</li> <li>Understands English spoken in a variety of nonnative accents</li> </ul>	<ul style="list-style-type: none"> <li>Comprehends virtually all aspects of a wide variety of academic material for non-specialists</li> <li>Reads at near-native speed</li> <li>Rarely requires use of dictionary</li> </ul>	<ul style="list-style-type: none"> <li>Writes complex documents such as research reports using appropriate style and vocabulary</li> <li>Grammar and orthographic accuracy is at near-native level</li> <li>Expresses complex relationships between ideas</li> </ul>	<ul style="list-style-type: none"> <li>Communicates accurately and effectively on practically all academic and social topics in culturally appropriate ways</li> <li>Pronunciation is close to that of native speakers</li> </ul>
5.9	<b>C1</b> ADVANCED	<ul style="list-style-type: none"> <li>Identifies attitude and purpose of speakers</li> <li>Grasps main ideas and the majority of supporting details from academic lectures</li> <li>Is challenged by complex social and cultural references</li> </ul>	<ul style="list-style-type: none"> <li>Understands main ideas and most of the details of academic texts, journal articles, and abstracts</li> <li>Requires little extra reading time</li> </ul>	<ul style="list-style-type: none"> <li>Vocabulary is strong in specialty</li> <li>Satisfies demands of most general academic tasks with occasional grammar and style mistakes</li> <li>Exhibits fairly good organization and development</li> </ul>	<ul style="list-style-type: none"> <li>Vocabulary is strong in specialty</li> <li>Satisfies demands of most general academic tasks with occasional grammar and style mistakes</li> <li>Exhibits fairly good organization and development</li> </ul>
5.0					
4.9	<b>B2</b> UPPER INTERMEDIATE	<ul style="list-style-type: none"> <li>Identifies main ideas and details in conversation</li> <li>Occasionally needs to ask for repetition or clarification</li> <li>Begins to determine the attitudes of speakers</li> <li>Understands main ideas from academic lectures, but misses significant details</li> </ul>	<ul style="list-style-type: none"> <li>Utilizes contextual and syntactic cues to interpret meaning of complex sentences and new vocabulary</li> <li>Gathers most main ideas from textbooks and articles, but has an uneven grasp of details</li> <li>Misinterprets some abstract content and cultural references</li> </ul>	<ul style="list-style-type: none"> <li>Writes reasonably coherent essays on familiar topics, but with some grammatical weakness</li> <li>Does not have complete grasp of stylistic features</li> <li>Vocabulary frequently lacks precision and sophistication</li> </ul>	<ul style="list-style-type: none"> <li>Begins to express abstract concepts, especially on familiar topics</li> <li>Fluency is occasionally hampered by gaps in vocabulary and grammar</li> <li>Expresses viewpoints in fairly long stretches of discourse</li> <li>Sometimes is asked to repeat words or phrases</li> </ul>
4.0					

ITEP	CEFR	Listening	Reading	Writing	Speaking
3.9	<b>B1</b> INTERMEDIATE	<ul style="list-style-type: none"> <li>Grasps the general outline of topics discussed in an academic setting</li> <li>Unfamiliarity with complex structures and higher-level vocabulary leaves major gaps in understanding</li> </ul>	<ul style="list-style-type: none"> <li>Limited vocabulary impedes speed</li> <li>Grasps the gist of material on familiar subjects, and identifies some significant details</li> <li>Follows step-by-step instructions in exams, labs, and assignments</li> </ul>	<ul style="list-style-type: none"> <li>Communicates basic ideas, but with weak organizational structure and grammatical mistakes that sometimes hinder understanding</li> <li>Expresses him/herself with some circumlocution on topics such as family, hobbies, work, etc</li> </ul>	<ul style="list-style-type: none"> <li>Manages day-to-day communications with peers and instructors, marked by frequent grammar and vocabulary errors</li> <li>Pronunciation requires significant effort from listeners</li> </ul>
3.5		<ul style="list-style-type: none"> <li>Maintains comprehension during conversations on familiar topics</li> <li>Relies heavily on non-verbal cues and repetition</li> <li>Understands very basic exchanges when spoken slowly using simple vocabulary</li> </ul>	<ul style="list-style-type: none"> <li>Major vocabulary gaps lead to frequent inaccurate or incomplete comprehension, and slow pace</li> <li>Understands simplified material</li> <li>Begins to determine the meaning of words by familiar surrounding context</li> </ul>	<ul style="list-style-type: none"> <li>Limited vocabulary results in repetitive style and simple sentences</li> <li>Considerable effort required by the reader to identify intended meaning</li> <li>Uses only basic vocabulary and simple grammatical structures</li> </ul>	<ul style="list-style-type: none"> <li>Generates simple questions, greetings, expressions of needs, and preferences</li> <li>Pronunciation requires significant effort from listeners</li> <li>Pronunciation often obscures meaning</li> </ul>
2.5	<b>A2</b> ELEMENTARY	<ul style="list-style-type: none"> <li>Understands simple greetings, statements, and questions when spoken with extra clarity</li> <li>Follows simple familiar instructions</li> <li>Frequently requires repetition for comprehension</li> <li>Understands a few isolated words or phrases spoken slowly</li> </ul>	<ul style="list-style-type: none"> <li>Comprehends only highly simplified phrases or sentences</li> <li>Recognizes familiar cohesive devices and basic pronouns</li> <li>Demonstrates understanding of a few simple grammatical and lexical structures</li> <li>Recognizes the alphabet and isolated words</li> </ul>	<ul style="list-style-type: none"> <li>Writes only short simple sentences. often characterized by errors that obscure meaning</li> <li>Provides personal details with correct spelling and can copy familiar words and phrases</li> <li>Produces isolated words and phrases</li> </ul>	<ul style="list-style-type: none"> <li>Capable of short simple presentation on familiar topic</li> <li>Responds to simple statements or questions</li> <li>Speech is marked with non-native stress and intonation patterns</li> <li>Communication is understood for short utterances</li> <li>Pauses, false starts, and reformulation are common</li> <li>Communicates with single words and short phrases at “survival level”</li> <li>Intense listener effort required</li> <li>Produces a few isolated words and phrases</li> <li>Pronunciation is mostly unintelligible</li> </ul>
2.4		<ul style="list-style-type: none"> <li>Understands simple greetings, statements, and questions when spoken with extra clarity</li> <li>Follows simple familiar instructions</li> <li>Frequently requires repetition for comprehension</li> <li>Understands a few isolated words or phrases spoken slowly</li> </ul>	<ul style="list-style-type: none"> <li>Comprehends only highly simplified phrases or sentences</li> <li>Recognizes familiar cohesive devices and basic pronouns</li> <li>Demonstrates understanding of a few simple grammatical and lexical structures</li> <li>Recognizes the alphabet and isolated words</li> </ul>	<ul style="list-style-type: none"> <li>Writes only short simple sentences. often characterized by errors that obscure meaning</li> <li>Provides personal details with correct spelling and can copy familiar words and phrases</li> <li>Produces isolated words and phrases</li> </ul>	<ul style="list-style-type: none"> <li>Capable of short simple presentation on familiar topic</li> <li>Responds to simple statements or questions</li> <li>Speech is marked with non-native stress and intonation patterns</li> <li>Communication is understood for short utterances</li> <li>Pauses, false starts, and reformulation are common</li> <li>Communicates with single words and short phrases at “survival level”</li> <li>Intense listener effort required</li> <li>Produces a few isolated words and phrases</li> <li>Pronunciation is mostly unintelligible</li> </ul>
0.1	<b>A1</b> BEGINNER	<ul style="list-style-type: none"> <li>Understands simple greetings, statements, and questions when spoken with extra clarity</li> <li>Follows simple familiar instructions</li> <li>Frequently requires repetition for comprehension</li> <li>Understands a few isolated words or phrases spoken slowly</li> </ul>	<ul style="list-style-type: none"> <li>Comprehends only highly simplified phrases or sentences</li> <li>Recognizes familiar cohesive devices and basic pronouns</li> <li>Demonstrates understanding of a few simple grammatical and lexical structures</li> <li>Recognizes the alphabet and isolated words</li> </ul>	<ul style="list-style-type: none"> <li>Writes only short simple sentences. often characterized by errors that obscure meaning</li> <li>Provides personal details with correct spelling and can copy familiar words and phrases</li> <li>Produces isolated words and phrases</li> </ul>	<ul style="list-style-type: none"> <li>Capable of short simple presentation on familiar topic</li> <li>Responds to simple statements or questions</li> <li>Speech is marked with non-native stress and intonation patterns</li> <li>Communication is understood for short utterances</li> <li>Pauses, false starts, and reformulation are common</li> <li>Communicates with single words and short phrases at “survival level”</li> <li>Intense listener effort required</li> <li>Produces a few isolated words and phrases</li> <li>Pronunciation is mostly unintelligible</li> </ul>

# APPENDIX E: SUMMARY OF STEPS TO MINIMIZE CONTENT-IRRELEVANT TEST MATERIAL

- Implement best practices in item writing to reduce the likelihood that “test wise” test-takers will be able to select the best answer, through cues in the test, without needing to understand the test item itself (for example, by selecting the lengthiest option, eliminating options that are saying the same thing in different ways)
- Avoid content that may influence test-takers’ performance on the test—items respect people’s value, beliefs, identity, culture, and diversity.
- Topics on which a set of items may be based are submitted by item writers to BES; BES pre-approves topics prior to item writing
- Assessment content reflects the domain and difficulty of knowledge of someone with the educational level of a high school junior who expects to attend college. The content reflects materials that an examinee would be expected to encounter in textbooks, journals, classroom lectures, extra-curricular activities, and social situations involving students and professors. Items do not reflect specialized knowledge.
- Write items at an appropriate reading level (no higher than grade 12; lower reading level for easier items); avoid words that are used with low frequency
- Test items assess comprehension within the item, as opposed to common knowledge. Passages establish adequate context for the topic, but then go on to introduce material that is not generally known. Examinees should be able to gain sufficient new information from the passage to answer the questions.
- Content does not unduly advantage examinees from particular regions of the world.

Classified by month change

Table 1				Table 2				Table 3				Table 4			
Year	Jan	Feb	Total	Year	Jan	Feb	Total	Year	Jan	Feb	Total	Year	Jan	Feb	Total
2010	10,000	15,000	25,000	2011	12,000	18,000	30,000	2012	14,000	20,000	34,000	2013	16,000	22,000	38,000
2014	18,000	25,000	43,000	2015	20,000	28,000	48,000	2016	22,000	30,000	52,000	2017	24,000	32,000	56,000
2018	26,000	35,000	61,000	2019	28,000	38,000	66,000	2020	30,000	40,000	70,000	2021	32,000	42,000	74,000
2022	34,000	45,000	79,000	2023	36,000	48,000	84,000	2024	38,000	50,000	88,000	2025	40,000	52,000	92,000
2026	42,000	55,000	97,000	2027	44,000	58,000	102,000	2028	46,000	60,000	106,000	2029	48,000	62,000	110,000
2030	50,000	65,000	115,000	2031	52,000	68,000	120,000	2032	54,000	70,000	124,000	2033	56,000	72,000	128,000
2034	58,000	75,000	133,000	2035	60,000	78,000	138,000	2036	62,000	80,000	142,000	2037	64,000	82,000	146,000
2038	66,000	85,000	151,000	2039	68,000	88,000	156,000	2040	70,000	90,000	160,000	2041	72,000	92,000	164,000
2042	74,000	95,000	169,000	2043	76,000	98,000	174,000	2044	78,000	100,000	178,000	2045	80,000	102,000	182,000
2046	82,000	105,000	187,000	2047	84,000	108,000	192,000	2048	86,000	110,000	196,000	2049	88,000	112,000	200,000
2050	90,000	115,000	205,000	2051	92,000	118,000	210,000	2052	94,000	120,000	214,000	2053	96,000	122,000	218,000
2054	98,000	125,000	223,000	2055	100,000	128,000	228,000	2056	102,000	130,000	232,000	2057	104,000	132,000	236,000
2058	106,000	135,000	241,000	2059	108,000	138,000	246,000	2060	110,000	140,000	250,000	2061	112,000	142,000	254,000
2062	114,000	145,000	259,000	2063	116,000	148,000	264,000	2064	118,000	150,000	268,000	2065	120,000	152,000	272,000
2066	122,000	155,000	277,000	2067	124,000	158,000	282,000	2068	126,000	160,000	286,000	2069	128,000	162,000	290,000
2070	130,000	165,000	295,000	2071	132,000	168,000	300,000	2072	134,000	170,000	304,000	2073	136,000	172,000	308,000
2074	138,000	175,000	313,000	2075	140,000	178,000	318,000	2076	142,000	180,000	322,000	2077	144,000	182,000	326,000
2078	146,000	185,000	331,000	2079	148,000	188,000	336,000	2080	150,000	190,000	340,000	2081	152,000	192,000	344,000
2082	154,000	195,000	349,000	2083	156,000	198,000	354,000	2084	158,000	200,000	358,000	2085	160,000	202,000	362,000
2086	162,000	205,000	367,000	2087	164,000	208,000	372,000	2088	166,000	210,000	376,000	2089	168,000	212,000	380,000
2090	170,000	215,000	385,000	2091	172,000	218,000	390,000	2092	174,000	220,000	394,000	2093	176,000	222,000	398,000
2094	178,000	225,000	403,000	2095	180,000	228,000	408,000	2096	182,000	230,000	412,000	2097	184,000	232,000	416,000
2098	186,000	235,000	421,000	2099	188,000	238,000	426,000	2100	190,000	240,000	430,000	2101	192,000	242,000	434,000
2102	194,000	245,000	439,000	2103	196,000	248,000	444,000	2104	198,000	250,000	448,000	2105	200,000	252,000	452,000
2106	202,000	255,000	457,000	2107	204,000	258,000	462,000	2108	206,000	260,000	466,000	2109	208,000	262,000	470,000
2110	210,000	265,000	475,000	2111	212,000	268,000	480,000	2112	214,000	270,000	484,000	2113	216,000	272,000	488,000
2114	218,000	275,000	493,000	2115	220,000	278,000	498,000	2116	222,000	280,000	502,000	2117	224,000	282,000	506,000
2118	226,000	285,000	511,000	2119	228,000	288,000	516,000	2120	230,000	290,000	520,000	2121	232,000	292,000	524,000
2122	234,000	295,000	529,000	2123	236,000	298,000	534,000	2124	238,000	300,000	538,000	2125	240,000	302,000	542,000
2126	242,000	305,000	547,000	2127	244,000	308,000	552,000	2128	246,000	310,000	556,000	2129	248,000	312,000	560,000
2130	250,000	315,000	565,000	2131	252,000	318,000	568,000	2132	254,000	320,000	572,000	2133	256,000	322,000	576,000
2134	258,000	325,000	583,000	2135	260,000	328,000	588,000	2136	262,000	330,000	592,000	2137	264,000	332,000	596,000
2138	266,000	335,000	601,000	2139	268,000	338,000	606,000	2140	270,000	340,000	610,000	2141	272,000	342,000	614,000
2142	274,000	345,000	623,000	2143	276,000	348,000	628,000	2144	278,000	350,000	632,000	2145	280,000	352,000	636,000
2146	282,000	355,000	641,000	2147	284,000	358,000	646,000	2148	286,000	360,000	650,000	2149	288,000	362,000	654,000
2150	290,000	365,000	663,000	2151	292,000	368,000	668,000	2152	294,000	370,000	672,000	2153	296,000	372,000	676,000
2154	298,000	375,000	681,000	2155	300,000	378,000	686,000	2156	302,000	380,000	690,000	2157	304,000	382,000	694,000
2158	306,000	385,000	703,000	2159	308,000	388,000	708,000	2160	310,000	390,000	712,000	2161	312,000	392,000	716,000
2162	314,000	395,000	721,000	2163	316,000	398,000	726,000	2164	318,000	400,000	730,000	2165	320,000	402,000	734,000
2166	322,000	405,000	743,000	2167	324,000	408,000	748,000	2168	326,000	410,000	752,000	2169	328,000	412,000	756,000
2170	330,000	415,000	765,000	2171	332,000	418,000	768,000	2172	334,000	420,000	772,000	2173	336,000	422,000	776,000
2174	338,000	425,000	783,000	2175	340,000	428,000	788,000	2176	342,000	430,000	792,000	2177	344,000	432,000	796,000
2178	346,000	435,000	805,000	2179	348,000	438,000	810,000	2180	350,000	440,000	814,000	2181	352,000	442,000	818,000
2182	354,000	445,000	831,000	2183	356,000	448,000	836,000	2184	358,000	450,000	840,000	2185	360,000	452,000	844,000
2186	362,000	455,000	853,000	2187	364,000	458,000	858,000	2188	366,000	460,000	862,000	2189	368,000	462,000	866,000
2190	370,000	465,000	875,000	2191	372,000	468,000	880,000	2192	374,000	470,000	884,000	2193	376,000	472,000	888,000
2194	378,000	475,000	893,000	2195	380,000	478,000	898,000	2196	382,000	480,000	902,000	2197	384,000	482,000	906,000
2198	386,000	485,000	911,000	2199	388,000	488,000	916,000	2200	390,000	490,000	920,000	2201	392,000	492,000	924,000
2202	394,000	495,000	931,000	2203	396,000	498,000	936,000	2204	398,000	500,000	940,000	2205	400,000	502,000	944,000
2206	402,000	505,000	953,000	2207	404,000	508,000	958,000	2208	406,000	510,000	962,000	2209	408,000	512,000	966,000
2210	410,000	515,000	975,000	2211	412,000	518,000	980,000	2212	414,000	520,000	984,000	2213	416,000	522,000	988,000
2214	418,000	525,000	993,000	2215	420,000	528,000	998,000	2216	422,000	530,000	1002,000	2217	424,000	532,000	1006,000
2218	426,000	535,000	1011,000	2219	428,000	538,000	1016,000	2220	430,000	540,000	1020,000	2221	432,000	542,000	1024,000
2222	434,000	545,000	1031,000	2223	436,000	548,000	1036,000	2224	438,000	550,000	1040,000	2225	440,000	552,000	1044,000
2226	442,000	555,000	1053,000	2227	444,000	558,000	1058,000	2228	446,000	560,000	1062,000	2229	448,000	562,000	1066,000
2230	450,000	565,000	1075,000	2231	452,000	568,000	1080,000	2232	454,000	570,000	1084,000	2233	456,000	572,000	1088,000
2234	458,000	575,000	1093,000	2235	460,000	578,000	1098,000	2236	462,000	580,000	1102,000	2237	464,000	582,000	1106,000
2238	466,000	585,000	1111,000	2239	468,000	588,000	1116,000	2240	470,000	590,000	1120,000	2241	472,000	592,000	1124,000
2242	474,000	595,000	1131,000	2243	476,000	598,000	1136,000	2244	478,000	600,000	1140,000	2245	480,000	602,000	1144,000
2246	482,000	605,000	1153,000	2247	484,000	608,000	1158,000	2248	486,000	610,000	1162,000	2249	488,000	612,000	1166,000
2250	490,000	615,000	1175,000	2251	492,000	618,000	1180,000	2252	494,000	620,000	1184,000	2253	496,000	622,000	1188,000
2254	498,000	625,000	1193,000	2255	500,000	628,000	1198,000	2256	502,000	630,000	1202,000	2257	504,000	632,000	1206,000
2258	506,000	635,000	1211,000	2259	508,000	638,000	1216,000	2260	510,000	640,000	1220,000	2261	512,000	642,000	1224,000
2262	514,000	645,000	1231,000	2263	516,000	648,000	1236,000	2264	518,000	650,000	1240,000	2265	520,000	652,000	1244,000
2266	522,000	655,000	1253,000	2267	524,000	658,000	1258,000	2268	526,000	660,000	1262,000	2269	528,000	662,000	1266,000
2270	530,000	665,000	1275,000	2271	532,000	668,000	1280,000	2272	534,000	670,000	1284,000	2273	536,000	672,000	1288,000
2274	538,000	675,000	1293,000	2275	540,000	678,000	12								

