

Designing Accessible Formative Assessment Tasks to Measure Argumentation Skills for English
Learners

Danielle Guzman-Orth¹, Yi Song², Jesse R. Sparks²

Educational Testing Service

Sacramento, California¹

Princeton, New Jersey²

Abstract

This study investigated the challenges and opportunities in developing a computer-delivered English language arts (ELA) task intended to improve the accessibility of the task for middle school English Learners (ELs). Cognitive labs with eight ELs with varying language proficiency levels provided rich insight to student–task interaction and how the accessibility of the task could be improved to enhance student understanding and to support valid integration of the task as a part of students’ formative assessment process. In this presentation, we will share the results from our research and discuss our iterative approach to improving ELA task accessibility for ELs’ formative assessment process.

Keywords: English learners, formative assessment, accessibility, argumentation

Introduction

Formative assessment is gaining popularity as a method to provide accurate, timely, and actionable information that can be used by teachers and students to improve learning (Heritage, 2010; Heritage, Kim, Vendlinski, & Herman, 2009). For English learners (ELs), this process may be confounded by their varying levels of English language proficiency (ELP), and teachers may be challenged by the need to provide appropriate supports to meet the diverse needs of their students (Shore, Wolf, & Blood, 2013). Informed by learning progressions and the Common Core State Standards, this study investigated how ELs interact with a scenario-based assessment of English Language Arts (ELA) argumentation skills as part of their formative assessment process. Research questions guiding this study were:

- 1) How do middle school ELs interact with an ELA task designed to measure argumentation skills in English?
- 2) How can we improve the task design so that it will elicit valid information about EL students' argumentation knowledge, skills, and abilities?

Theoretical Framework

Argumentation skills are essential for success in college, career, and life; therefore, these skills play a prominent role in recent educational reform efforts such as the Common Core State Standards (Council of Chief State School Officers & National Governors Association, 2010). Despite the importance of argumentation skills, many students cannot write sound arguments or critically evaluate arguments, as evidenced by a variety of large-scale assessments and empirical studies (e.g., Ferretti, Lewis, & Andrews-Weckerly, 2009; National Center for Educational Statistics, 2012; Perkins, Farady, & Bushey, 1991; Song, Deane, & Fowles, 2017). One of the most complex academic skills, argumentation has not been well-supported in instructional

practice, which frequently emphasizes basic written composition and use of specific templates while doing little to develop arguments and critical thinking (Hillocks, 2002). Furthermore, traditional assessments of argumentation, which typically require students to write an essay on a single prompt, offer little information about why students may have failed to accomplish this task.

To gather relevant evidence about students' argumentation skills, we designed a scenario-based assessment aligned to a set of hypothesized learning progressions (LPs; Deane & Song, 2015). Informed by cognitive and learning sciences research (e.g., Bereiter & Scardamalia, 1987; Graham & Perin, 2007; Hayes & Flower, 1980; Kuhn, 1991), argumentation LPs describe how argumentation skills develop into sophistication, characterizing the qualitative shifts that occur as students reach higher levels in four strands of skills: (1) *Appeal building*: understanding an audience's values and beliefs; (2) *Taking a position*: developing a position and understanding other perspectives; (3) *Reasons and evidence*: using reasons and evidence to support an argument and to evaluate others' arguments; and (4) *Framing a case*: organizing and presenting an argument logically. We used LPs as a general framework to determine the targeted skills (i.e., position, reasons, and evidence) and levels and sequences of the activities within the scenario-based assessment.

Furthermore, we recognize that designing opportunities for students to demonstrate their content knowledge, skills, and abilities (KSAs) while minimizing construct-irrelevant variance is a complex process for EL students taking ELA tasks. Design procedures must incorporate attention to diverse student needs (Guzman-Orth, Laitusis, Thurlow, & Christensen, 2016; Pitoniak et al., 2009) as well as a deep understanding of how students learn and express their KSAs (Ketterlin-Geller, 2017) and how tasks can elicit relevant evidence so that the

interpretations of the data can support the assessment claims (Mislevy, Steinberg, & Almond, 2003). To support this process, a multi-disciplinary design team should also be a central feature of any task design, especially when designing tasks for ELs who are a diverse population with a variety of needs, especially linguistic needs (Solano-Flores, Shade, & Chrzanowski, 2014).

Our research occurred in a multi-phase sequential process. First, we conducted cognitive labs with middle school EL students to understand their challenges in an argumentation assessment called *Seaball – Semester at Sea*. In the context of a fictional study abroad program, students need to demonstrate argumentation skills across five increasingly difficult activities (aligned to argumentation LPs). Second, we conducted a literature review to investigate how ELA content and assessment practices are made access for ELs. Finally, we synthesized information from the above two phases of work and made revisions to the *Seaball* task to improve access for middle school ELs.

Phase 1: Cognitive Labs

Methods

Participants. Participants in cognitive labs were eight 7th grade Spanish-English bilingual ELs (three females) from an urban middle school in the mid-Atlantic region. Their overall ELP score ranged from Beginning (n=3), Early Intermediate (n=1), and Intermediate (n=4) (see Table 1 for demographics).

Instruments. The instruments of the cognitive lab study included a background information questionnaire, the *Seaball* task, and a post-task survey.

Background information questionnaire (BIQ). The students' teacher filled out the BIQ, including student demographics such as gender, EL status, and standardized test proficiency levels for English language proficiency, ELA, and Math.

Seaball task. In the *Junk Food* scenario, students engage in argumentation about whether junk food should be sold to students, by interacting with game characters, evaluating claims and evidence, and making policy recommendations. *Seaball Junk Food* includes five activities of varying difficulty. First, students interview characters, and classify their opinions into ban or allow categories (*Interview*). Next, students evaluate four candidates' profiles and then select an expert to address the students about junk food (*Select a Speaker*). In the third activity, students identify the main claim, reasons, and evidence in the speech (*Identify Arguments*). Students then make a recommendation to the student council (i.e., ban or allow junk food) with supporting reasons (*Make Recommendation*). Finally, students sort food items into junk food and healthy food categories, which involves using relevant evidence and evaluating people's arguments (*Establish Criterion*). The task is automatically scored, with 100 points possible.

Post-study survey. A 17-item survey (12 likert-type, five open-ended) elicited information about students' experience and perceptions of the *Seaball* task.

Procedure. Students participated in one-on-one cognitive labs with a trained researcher. Each cognitive lab lasted 60-90 minutes and was audio-recorded. The students individually played through the *Seaball* task on the computer while researchers observed students' interactions (i.e., documenting issues with usability, language, and engagement). Researchers followed these observations with interview questions to learn more about students' perceptions and experience completing the task (e.g., difficult or unknown vocabulary, navigation issues). After completing *Seaball*, students completed the survey.

Phase 1: Cognitive Lab Results

Observation notes, interview transcripts, survey responses, and task log files (i.e., student actions and scores) were analyzed. Overall, these data sources indicated that although the students reported enjoying the experience of *Seaball*, the task was quite challenging for ELs. Notably, task scores were overall very low (range: 27 to 78) due to students being unable to complete each activity before researchers helped them proceed to the next section due to time constraints. Specifically, five major themes emerged from qualitative analysis, related to issues with difficulty, usability, engagement, language, and timing.

Two primary sources of difficulty were related to usability and linguistic complexity, including complexity of the directions, as well as the content. Additionally, cultural accessibility posed an issue, with EL students being relatively unfamiliar with the specific context of studying abroad on a ship. Students were also unfamiliar with colloquialisms (i.e., junk food) and idiomatic expressions (i.e., jokes and humorous dialogue, a game-like design element). The researchers had to help the students as they progressed through each activity in both usability and linguistic aspects (as a result, we recommend interpreting the student performance scores with some caution). These sources of difficulty also greatly affected students' time on task; students generally took extended time to decode and comprehend task directions, or they would try to "figure it out" through trial and error. Students were fatigued by the time they reached the end of the task. In general, EL students experienced difficulty in all aspects of *Seaball*. Specific issues for each *Seaball* activity are described below.

Activity 1: Interview. All students were engaged in this activity, but performance varied widely ($M = 8/16$, or 50% correct; see Table 2). Most students needed assistance with the language, including help understanding the directions and the T-chart graphic organizer (i.e., to

classify opinions as ban/allow). Specific vocabulary was also troublesome, including construct-relevant (e.g., *opinion, ban, allow*), context-relevant (e.g., *junk food, concerned*), and general academic vocabulary (*at least*). Several students also showed some usability issues (i.e., using drag and drop). One student needed support during the entire activity.

Activity 2: Select a Speaker. Similar to Activity 1, students were engaged in the second activity, but their performance indicates that they had difficulty with the task (range 35-75%; $M = 9.89$ or 50% see Table 3). Almost all students needed help with usability; most reported confusion about how to proceed. This confusion could also be related to students' overall language difficulties, in addition to specific vocabulary (i.e., *junk, previous voyage, focus on this side of the debate, issue, buying, and include*). Three students also needed assistance navigating among the speakers' profiles (presented in tabbed format), and one had difficulty submitting a response. One student in particular (ArgCL7) had some difficulty with the typing required in the activity and was unable to produce an interpretable response.

Activity 3: Identify Arguments. Students were engaged, and performed slightly better on this activity compared to the previous two activities (range 13-100%; $M = 10.38$ or 69%; see Table 4). Some students needed assistance with directions and did not understand what to do. Some students did not know where to click on the screen to select the speaker's reasons. One student commented that there was a lot of information to synthesize and process, and another had spelling problems during typing. Two students in particular had difficulty with the English language in their typed responses. One student (ArgCL6) typed: "the soda and candys is bag for the salud of student" [*sic*] (*salud* is the Spanish equivalent of "health"). Another student (ArgCL7) showed reliance on phonetic spelling: "The yunk food is bad to de people and hte yunk food can be bad on the featurer." [*sic*]

Activity 4: Make Recommendation. Overall, students performed in the mid-range again (range 33-100%, $M = 10.75$, or 60%; see Table 5). At this point, researchers reported that students were experiencing some fatigue and needed assistance to move forward.

Activity 5: Establish Criterion. Overall, the data still indicate a wide range in performance (range of 0-88%, $M = 16.50$, or 52%; see Table 6). All aspects of this section were difficult; interviewers reported that most students needed assistance throughout the activity, to either rephrase the directions, overcome usability issues (clicking in the correct areas), or to rephrase the arguments characters provided.

Overall, results of the cognitive lab study identified multiple challenges faced by EL students taking a scenario-based argumentation skills assessment. To inform strategies to improve task accessibility, we next conducted a literature review.

Phase 2: Review of the Literature Using Empirical Evidence to Inform Accessible Task Design

A literature review regarding EL English language arts instruction, assessment, and accommodations was conducted to further validate the information obtained from the cognitive labs. Search terms included combinations of “English learner, English language learner, English language learning, EL, and ELL” and “English language arts, ELA, language arts, argumentation” and “accessible, accessibility, accommodation, support”. Academic journals as well as teacher handbooks and practitioner-focused websites identified through major search engines (e.g., Google, Google Scholar, EBSCOhost). Sources reviewed indicate there is no one-size-fits-all approach to accessibility for ELs. Many means of input and output (e.g., oral and print language in English and students’ home language, visuals, hands-on experiences) are recommended to allow ELs opportunity to access content and to demonstrate their KSAs, and

these elements should be combined with clearly articulated learning goals, pre-teaching activities, multiple opportunities for practice, immediate feedback, collaborative learning with peers, frequent re-teaching, and scaffolding (e.g., Gersten & Baker, 2000; Goldenberg, 2013; Robertson, 2017).

The EL accommodation literature indicates that EL accommodations are typically designed to minimize linguistic construct-irrelevant variance that may impact the students' performance (although evidence points to the mixed efficacy of these supports; see Kieffer, Lesaux, Rivera, & Francis, 2009; Pennock-Roman & Rivera, 2011). Supports range from glossaries, including both in English and in students' home language (Cohen, Tracey, & Cohen, 2017; Graves, August, Mancilla-Martinez, 2012; Solano-Flores, et al., 2014; Wolf et al., 2012a; Wolf, Kim, & Kao, 2012b); read aloud supports (Buzick & Stone, 2014; Higgins & Katz, 2013; Higgins, Russell, & Hoffman, 2005; Laitusis, 2010; Wolf et al., 2012b), linguistic modification or use of plain English (Abedi & Lord, 2001; Abedi et al., 2000; Abedi et al., 2005; Abedi, & Sato, 2007), and home language translations (Solano-Flores, 2006; Solano-Flores, et al., 2014; Solano-Flores, Trumbull, & Nelson-Barber, 2002).

Best practices for assessment design also offer guidelines that can improve student access. Elements such as Universal Design (UD) for assessment (Thompson, Johnstone, & Thurlow, 2002) or for learning (CAST, 2014) emphasize the need to ensure that the language and any related visuals in the assessment are designed to minimize construct-irrelevant variance. Information should be clear, simple, and intuitive, and presented in multiple modalities if it does not compromise measurement of the construct. Best practices for ELs emphasize the importance of attending to student needs throughout the test development process, from conceptualization (e.g., clearly articulating the target population and use cases; Guzman-Orth et al., 2016) to score

reporting (Pitoniak et al., 2009). Placing the compilation of these practices within cognitive models for learning (Ketterlin-Geller, 2017), and evidence-centered design (ECD; Mislevy et al., 2003; Hansen & Mislevy, 2005, 2008) ensures that each step of the development process is conceptually and theoretically tied to the end goal to minimize the need for retrofitting the test for special populations. Taken together, our literature review suggests that these approaches are critical building blocks to implement at the initial stages of development for test design process.

Phase 3: Using Empirical Evidence to Inform Accessible Task Design

In the final phase, we synthesized the findings from the data collection and the literature review. This synthesis indicates that the *Seaball* task could be heavily revised to better promote accessibility for ELs. Consistent with existing recommendations (e.g., Solano-Flores et al., 2014), we assembled a multi-disciplinary review team consisting of experts in the construct, accessibility for ELs, and assessment design for ELs. The multi-disciplinary review team designed and implemented revisions to the *Seaball* task, allowing for multiple rounds of review and consensus building from the team. Several revisions were considered for *Seaball*, but ultimately, the following changes were applied: reducing length, modifying the language, and adding supports (i.e., a glossary, read aloud instructions, a vocabulary-building activity, and visual cues). These revisions are detailed below (see Figure 1 for examples).

Task Length. Most students had difficulty completing the *Seaball* task. To address this issue, four of the five activities were removed, to focus on the first activity, *Interview*. The *Interview* activity content (i.e., identifying people's position on a controversial issue) is considered to be foundational in the argumentation LPs. The vocabulary introduced in this section is also pivotal for the subsequent activities. Thus, we now treat the first activity as a

discrete task that would allow students and teachers to start and stop the task as needed, allowing for re-teaching before proceeding to more complex activities.

Linguistic Modification. We modified the language in the *Seaball* task according to guidelines that were developed and adapted from our synthesis of the literature and best practices for ELs (see Figure 1). Assessment specialists first consulted word lists for both the construct-relevant and context-relevant vocabulary to gain familiarity with the language that should not be modified (see Appendix A). Modification guidelines included the following: (a) modify construct-irrelevant language; (b) do not modify construct-relevant language; (c) do not modify the context-specific language (e.g., do not change the setting of the task from a semester-at-sea program to something more culturally accessible); (d) use familiar/frequently used vocabulary for middle school ELs; (e) reduce sentence length; (f) use active voice and minimize passive voice; (g) shorten nominals/noun phrases; (h) reduce complex question phrases; (i) reduce comparative structures; (j) reduce prepositional phrases; (k) simplify sentence and discourse structure; (l) reduce subordinate clauses; (m) reduce conditional clauses; (n) reduce relative clauses; (o) use specific language rather than abstractions; (p) reduce negation; and (q) increase cultural accessibility (to support ELs who may not have prior knowledge of the semester abroad context). Taken together, these guidelines were applied to simplify the linguistic complexity of the task to an intermediate level of proficiency without changing measurement of the argumentation construct.

Glossary. An English glossary was added to assist in understanding terminology that was considered to be context-relevant (i.e., to the semester-at-sea setting). The embedded, pop up glossary in *Seaball* was designed to function similarly to the pop up glossary supports in multi-state consortia assessments (e.g., PARCC, Smarter Balanced), since students may have more

familiarity with and exposure to this type of design. Similar to the linguistic modification guidelines, assessment specialists first consulted Appendix A to gain familiarity with what context-relevant vocabulary could be glossed (since it was not removed in the linguistic modification step). Our approach to glossary guidelines comes from our synthesis and adaptation of existing guidelines in the field, such as: (a) gloss any remaining words/phrases that are construct-irrelevant; (b) gloss difficult vocabulary, including polysemous words (words with multiple meanings), false cognates, culturally-relevant terms, idiomatic expressions, regional variations, etc.; (c) gloss words the first time they appear (i.e., subsequent words and plural forms do not need a gloss); (d) limit the number of words/phrases that are glossed on each page/screen (i.e., if there are multiple glossed words, consider modifying the language further to reduce the need for multiple glosses); and (e) glossary entries should be relatively short words or phrases, and use concise, clear language and active voice.

Read Aloud. A read aloud (i.e., voice-over audio) component was added to the task directions. This modification provides multi-modal input for students, especially those who may be struggling readers.

Vocabulary-Building Activity. We incorporated a vocabulary-building activity to ensure that students would have an opportunity learn key argument vocabulary prior to interacting with the task (i.e., *opinion, ban, allow, reasons, and evidence*). This activity was designed as a short quiz with immediate feedback to determine the students' prior knowledge of the key vocabulary and to correct misunderstandings. The quiz includes five three-option multiple choice questions. After each response, students receive immediate feedback and the correct definition. At the end of the quiz, students can review all five definitions.

Visual Cues. Because usability and navigation posed difficulties for the students, visual cues were added to help attract the students' attention to places on screen where they need to look or click next, including a glowing, pulsating yellow highlight appearing around buttons (see Figure 1). In sum, these principled revisions may mitigate the difficulties that ELs experienced.

Discussion

This study investigated how ELs interacted with a scenario-based task measuring argumentation skills, and how that task can be improved in a principled fashion to promote greater access for ELs to elicit evidence about their KSAs. Overall, data sources from a cognitive lab study (including observations, surveys, interviews, and performance data) conducted in Phase 1 overwhelmingly indicate that students enjoyed interacting with the task, but also experienced difficulties; these primarily involved usability and linguistic considerations, which affected ELs' ability to independently access the content and demonstrate their knowledge and skills. Phase 2 included a review of relevant literature to investigate best practices for ELA instruction and assessment for ELs, including accommodations. Findings suggest that promoting access is not a one-size-fits-all approach, and that there is no single instructional activity or accommodation that works for all students, suggesting that a combination of approaches should be used.

Taken together, empirical findings from the cognitive labs and the literature review informed our approaches for Phase 3, task revision. Our principled approach to designing EL supports for the *Seaball* task included removing all but one activity, adding the vocabulary-building activity to increase familiarity, adding visual cues to aid navigation, modifying the language and adding a glossary to support comprehension, and a read aloud component to help alleviate reading load and provide an additional means of access to task directions. The revised activity resulting from this principled approach to accommodations will be tested in Fall 2017.

Directions for future research include conducting this pilot test of the revised task to gather validity evidence of the principled approach to designing EL supports for tasks. Classroom teachers should also implement the revised task as part of their formative assessment process, gathering necessary validity evidence to determine *Seaball*'s efficacy as a tool for teachers to use when teaching argumentation skills to ELs.

Further investigation of how ELs understand argumentation and demonstrate their reasoning is also needed. For example, ELs are learning argumentation content in English, while they are simultaneously learning the English language. Some students in the pilot study used their home language to help demonstrate their knowledge, but students may not receive credit from teachers if responses are not in English or are misspelled. Using a different approach, such as a multilingual framework like *translanguaging* (the process of bilingual speakers using all of their linguistic resources in English and their home language to demonstrate meaning; see Canagarajah, 2006) and conceptual scoring (the process of giving credit to correct responses, regardless of what language is used; see Barreuco, Lopez, Ong, & Lozano, 2012) is recommended for bilingual speakers taking assessments measuring their functional KSAs, rather than their limited English proficiency (Guzman-Orth, Lopez, & Tolentino, 2017; Guzman-Orth, Lopez, Tolentino, Sova, & Stolow, 2016; Lopez, Guzman-Orth, & Turkan, 2015; Lopez, Turkan, & Guzman-Orth, 2017). Students' responses to *Seaball* indicated that they partially understood the activity and could produce some relevant rationales (e.g., "the soda and candys is bag for the salud of student" [sic]). Despite typos and grammatical errors, the student states that junk food (soda and candy) are bad for students' health (*salud* is the Spanish equivalent of "health"). The phonetic spelling in the response "The yunk food is bad to de people and hte yunk food can be bad on the featuer" [sic] indicates that, typos aside, the student has internally contextualized the

content of the task and is demonstrating knowledge that junk food can be harmful. This student's understanding includes the use of language at the morphemic level, the smallest, meaningful units of language (e.g., phonemes): the students' use of the "y" in *yunk [sic]* mirrors their accent and pronunciation of the /j/ English phoneme, which sounds like the "ll" letter in Spanish that is pronounced like "y" ("j/yah"). In these cases, although students were having some difficulties throughout the various components of the activities, they were actively seeking ways to meaningfully participate and demonstrate their knowledge using all of their linguistic resources.

This study provides an example of how multiple sources of data can be synthesized to provide evidence to inform assessment design decisions. Despite widespread attention to evidence-based practices for EL accessibility on large scale assessments, less attention has been focused on EL accessibility for classroom-based assessments or activities that can be used as part of the formative assessment process to support teachers and students in monitoring learning. Additionally, this study is unique in that it focuses on accommodating EL students' language needs on ELA tasks – an effort that is highly complex due to the integrated language demands of the content, the need to maintain the measurement of the argumentation construct, students' ongoing acquisition of the English language, and the need to provide language supports for their ongoing language acquisition. We hope that this study can increase attention to the diverse needs for EL students to promote accessibility on ELA tasks so that those ELA tasks can yield more valid, actionable results for teachers and students to monitor and support their learning.

References

- Abedi, J., Courtney, M., Mirocha, J., Leon, S., & Goldberg, J. (2005). Language Accommodations for English Language Learners in Large-Scale Assessments: Bilingual Dictionaries and Linguistic Modification. CSE Report 666. *National Center for Research on Evaluation, Standards, and Student Testing (CRESST)*.
- Abedi, J., Lord, C. (2001). The language factor in mathematics tests. *Applied Measurement in Education, 14*(3), 219–234.
- Abedi, J., Lord, C., Hoffsetter, C., Baker, E. (2000). Impact of accommodation strategies on English language learners' test performance. *Educational Measurement: Issues and Practice, 19*(3), 16–26.
- Abedi, J., & Sato, E. (2007). Linguistic modification. *A report prepared for the US Department of Education LEP Partnership. Washington, DC: US Department of Education*.
- Barrueco, S., Lopez, M., Ong, C., & Lozano, P. (2012). *Assessing Spanish-English bilingual preschoolers: A guide to best approaches and measures*. Baltimore, MD: Brookes.
- Bereiter, C., & Scardamalia, M. (1987). *The psychology of written composition*. Hillsdale, NJ: Lawrence Erlbaum.
- Buzick, H. M., & Stone, E. A. (2014). A meta-analysis of research on the read aloud accommodation. *Educational Measurement: Issues and Practice, 33*(3), 17–30.
- Canagarajah, S. (2006). Changing communicative needs, revised assessment objectives: Testing English as an international language. *Language Assessment Quarterly, 3*(3), 229–242.
- Center for Applied Special Technology (CAST). (2014). UDL guidelines version 2.0. Retrieved from: <http://www.udlcenter.org/aboutudl/udlguidelines>

- Cohen, D., Ryan, T., & Cohen, J. (2017). On the Effectiveness of Pop-up English Language Glossary Accommodations for EL Students in Large-scale Assessments. *Applied Measurement in Education*, (just-accepted).
- Council of Chief State School Officers, & National Governors Association (2010). *Common Core State Standards for English language arts and literacy in history/social studies, science, and technical subjects*. Washington, DC: Author.
<http://www.corestandards.org/the-standards/ELA-Literacy>
- Deane, P., & Song, Y. (2015). The key practice, ‘Discuss and Debate Ideas’: Conceptual framework, literature review, and provisional learning progressions for argumentation. (Research Report No. RR-15-33). Princeton, NJ: Educational Testing Service.
- Ferretti, R. P., Lewis, W. E., & Andrews-Weckerly, S. (2009). Do goals affect the structure of students’ argumentative writing strategies? *Journal of Educational Psychology*, *101*, 577–589.
- Gersten, R., & Baker, S. (2000). What we know about effective instruction practices for English-language learners. *Exceptional Children*, *66*(4), p. 454–470.
- Goldenberg, C. (2013). Unlocking the research on English learners. What we know –and don’t yet know about effective instruction. *American Educator*, *37*, 2, p. 4–11.
- Graham, S., & Perin, D. (2007). *Writing next: Effective strategies to improve writing of adolescents in middle and high schools – A report to the Carnegie Corporation of New York*. Washington, DC: Alliance for Educational Progress.
- Graves, M. F., August, D., & Mancilla-Martinez, J. (2012). *Teaching vocabulary to English language learners*. Teachers College Press.

- Guzman-Orth, D., Laitusis, C., Thurlow, M., & Christensen, L. (2016). Conceptualizing accessibility for English language proficiency assessments (Research Report No. RR-16-07). Princeton, NJ: Educational Testing Service. <http://dx.doi.org/10.1002/ets2.12093>
- Guzman-Orth, D., Lopez, A. A., & Tolentino, F. (2017). A Framework for the Dual Language Assessment of Young Dual Language Learners in the United States. *ETS Research Report Series*. doi: 10.1002/ets2.12165
- Guzman-Orth, D., Lopez, A. A., Tolentino, F., Sova, L., & Stolow, A. (2016, December). *Considerations in Assessing Dual Language Learners*. Presentation at the California Education Research Association, Sacramento, CA.
- Hansen, E. G., & Mislevy, R. J. (2005). Accessibility of computer-based testing for individuals with disabilities and English language learners within a validity framework. In M. Hricko & S. Howell (Eds.), *Online assessment and measurement: Foundation, challenges, and issues*. pp. 214–262, Hershey, PA: Idea Group Publishing, Inc.
- Hansen, E. G., & Mislevy, R. J. (2008). Design patterns for improving accessibility for test takers with disabilities (RR-08-49). Princeton, New Jersey: ETS Research Report Series.
- Hayes, J. R., & Flower, L. S. (1980). Identifying the organization of writing processes. In L. Gregg & E. R. Steinberg (Eds.), *Cognitive process in writing* (pp. 3–30). Hillsdale, NJ: Lawrence Erlbaum Associates.
- Heritage, M. (2010). Formative Assessment and Next-Generation Assessment Systems: Are We Losing an Opportunity? *Council of Chief State School Officers*.
- Heritage, M., Kim, J., Vendlinski, T., & Herman, J. (2009). From evidence to action: A seamless process in formative assessment? *Educational Measurement: Issues and Practice*, 28(3), 24–31.

- Higgins, J., & Katz, M. (2013). A comparison of audio representations of mathematics content. *Journal of Special Education Technology*, 28(3), 59–66.
- Higgins, J., Russell, M., & Hoffman, T. (2005). Examining the effect of computer-based passage presentation on reading test performance. *The Journal of Technology, Learning, and Assessment*, 3(4), 1–35.
- Hillocks, G., Jr. (2002). *The testing trap: How state writing assessments control learning*. New York: Teachers College Press.
- Ketterlin-Geller, L. M. (2017). Understanding and improving accessibility for special populations. In A. A. Rupp and J. P. Leighton (Eds.). *The Handbook of Cognition and Assessment: Frameworks, Methodologies, and Application*. John Wiley & Sons, Inc.
- Kieffer, M. J., Lesaux, N. K., Rivera, M., & Francis, D. J. (2009). Accommodations for English language learners taking large-scale assessments: A meta-analysis on effectiveness and validity. *Review of Educational Research*, 79(3), 1168–1201.
- Kuhn, D. (1991). *The skills of argument*. Cambridge, England: Cambridge University Press.
- Laitusis, C. C. (2010). Examining the impact of audio presentation on tests of reading comprehension. *Applied Measurement in Education*, 23, 153–167.
- Lopez, A. A., Guzman-Orth, D. A., & Turkan, S. (2015). How might a translanguaging approach in assessment make tests more valid and fair for emergent bilinguals? In G. Valdés, K. Menken, & M. Castro (Eds.), *The Common Core and English language learners/emergent bilinguals: A guide for all educators* (pp. 266–267). Philadelphia, PA: Caslon.

- Lopez, A. A., Turkan, S., & Guzman-Orth, D. (2017). Conceptualizing the Use of Translanguaging in Initial Content Assessments for Newly Arrived Emergent Bilingual Students. *ETS Research Report Series*.
- Mislevy, R. J., Steinberg, L. S., & Almond, R. G. (2003). Focus article: On the structure of educational assessments. *Measurement: Interdisciplinary research and perspectives, 1*(1), 3–62.
- National Center for Education Statistics. (2012). *The Nation's Report Card: Writing 2011* (NCES 2012-470). Institute for Education Sciences, U.S. Department of Education, Washington, D.C.
- Pennock-Roman, M., & Rivera, C. (2011). Mean effects of test accommodations for ELLs and non-ELLs: A meta-analysis of experimental studies. *Educational Measurement: Issues and Practice, 30*(3), 10–28.
- Perkins, D. N., Farady, M., & Bushey, B. (1991). Everyday reasoning and the roots of intelligence. In J. F. Voss, D. N. Perkins, & J. W. Segal (Eds.), *Informal reasoning and education* (pp. 83–106). Hillsdale, NJ: Erlbaum.
- Pitoniak, M. J., Young, J. W., Martiniello, M., King, T. C., Buteux, A., & Ginsburgh, M. (2009). *Guidelines for the assessment of English language learners*. Princeton, NJ: ETS.
- Robertson, K. (2017). Supporting ELLs in the mainstream classroom: Language tips. Color in Colorado WETA Public Broadcasting. Retrieved from <http://www.colorincolorado.org/article/supporting-ells-mainstream-classroom-language-tips>
- Shore, J. R., Wolf, M. K., & Blood, I. (2013). ELFA Teacher's Guide. Retrieved from https://www.ets.org/s/research/pdf/elfa_teachers_guide.pdf

- Solano-Flores, G. (2006). Language, dialect, and register: Sociolinguistics and the estimation of measurement error in the testing of English language learners. *Teachers College Record*, 108(11), 2354–2379.
- Solano-Flores, G., Shade, C., & Chrzanowski, A. (2014). Smarter Balanced Assessment Consortium: Item accessibility and language variation conceptual framework. Retrieved from <https://portal.smarterbalanced.org/library/en/item-accessibility-and-language-variation-conceptual-framework.pdf>
- Solano-Flores, G., Trumbull, E., & Nelson-Barber, S. (2002). Concurrent development of dual language assessments: An alternative to translating tests for linguistic minorities. *International Journal of Testing*, 2(2), 107–129.
- Song, Y., Deane, P., & Fowles, M. E. (2017). Examining students' ability to critique arguments and exploring assessment and instructional implications. (Research Report No. RR-17-16). Princeton, NJ: Educational Testing Service.
- Thompson, S. J., Johnstone, C. J., & Thurlow, M. L. (2002). *Universal design applied to large scale assessments* (NCEO Synthesis Report No. 44). Minneapolis, MN: University of Minnesota, National Center on Educational Outcomes.
- Wolf, M. K., Kao, J. C., Rivera, N. M., & Chang, S. M. (2012a). Accommodation Practices for English Language Learners in States' Mathematics Assessments. *Teachers College Record*, 114(3).
- Wolf, M. K., Kim, J., & Kao, J. (2012b). The effects of glossary and read-aloud accommodations on English language learners' performance on a mathematics assessment. *Applied Measurement in Education*, 25(4), 347–374.

Table 1. Student Demographic Characteristics and *Seaball* Performance Score

ID	Gender	English Language Proficiency Test Score Levels					Content Test Score Levels		
		Listening	Speaking	Reading	Writing	Comprehension	Overall	ELA	Math
Arg1CL	M	A	EI	I	I	I	I	Partial	Partial
Arg2CL	M	EA	I	EI	I	I	I	Not Yet Met	Partial
Arg3CL	F	EA	EA	EI	I	EI	I	Partial	Not Yet Met
Arg4CL	F	EA	I	EI	I	I	I	NA	Partial
Arg5CL	M	EI	B	B	EI	B	B	NA	Not Yet Met
Arg6CL	M	EI	B	EI	EI	EI	B	NA	Partial
Arg7CL	M	EI	B	EI	EI	EI	B	Not Yet Met	Not Yet Met
Arg8CL	F	B	B	EI	EI	B	EI	NA	Not Yet Met

Table note: M = Male; F = Female. A = Advanced, EA = Early Advanced, I = Intermediate, EI = Early Intermediate, B = Beginner. ELA = English language arts test scores. Math = Mathematics test scores.

Table 2. Summary of Observations for Activity 1: Interview

ID	Overall ELP	Issues Observed				Activity 1 Classify an Opinion (out of 15)*
		Difficulty	Usability	Engagement	Language	
Arg1CL	I	2	0	1	0	0 (0%)
Arg2CL	I	1	2	2	1	9 (60%)
Arg3CL	I	2	0	2	1	15 (100%)
Arg4CL	I	1	1	2	0	6 (40%)
Arg5CL	B	1	1	2	0	6 (40%)
Arg6CL	B	1	1	2	0	9 (60%)
Arg7CL	B	2	0	2	1	6 (40%)
Arg8CL	EI	0	2	2	2	12 (80%)
<i>Mean Score</i>	--	<i>1.25</i>	<i>0.88</i>	<i>1.88</i>	<i>0.63</i>	<i>8 (53%)</i>

Table note: 0 = no, 1 = partial, 2 = yes. * = Students did receive support from the interviewer in responding to the questions in the task.

Table 3. Summary of Observations for Activity 2: Select a Speaker

ID	Overall ELP	Issues Observed				Activity 2 Select a Speaker (out of 20)*
		Difficulty	Usability	Engagement	Language	
Arg1CL	I	1	1	1	1	7 (35%)
Arg2CL	I	0	1	2	0	10 (50%)
Arg3CL	I	2	0	2	1	9 (45%)
Arg4CL	I	0	2	2	1	15 (75%)
Arg5CL	B	1	1	1	0	9 (45%)
Arg6CL	B	1	1	2	0	7 (35%)
Arg7CL	B	0	2	1	2	9 (45%)
Arg8CL	EI	0	2	2	0	13 (65%)
<i>Mean Score</i>	--	<i>0.63</i>	<i>1.25</i>	<i>1.63</i>	<i>0.63</i>	<i>9.89 (50%)</i>

Table note: 0 = no, 1 = partial, 2 = yes. * = Students did receive support from the interviewer in responding to the questions in the task.

Table 4. Summary of Observations for Activity 3: Identify Arguments

ID	Overall ELP	Issues Observed				Activity 3 Identify Arguments (out of 15)*
		Difficulty	Usability	Engagement	Language	
Arg1CL	I	2	0	1	0	2 (13%)
Arg2CL	I	1	1	2	0	12 (80%)
Arg3CL	I	2	0	2	0	15 (100%)
Arg4CL	I	1	1	2	0	15 (100%)
Arg5CL	B	2	1	1	0	12 (80%)
Arg6CL	B	2	0	2	0	5 (33%)
Arg7CL	B	1	0	2	0	10 (66%)
Arg8CL	EI	2	0	2	0	12 (80%)
<i>Mean Score</i>	--	<i>1.63</i>	<i>0.38</i>	<i>1.75</i>	<i>0</i>	<i>10.38 (69%)</i>

Table note: 0 = no, 1 = partial, 2 = yes. * = Students did receive support from the interviewer in responding to the questions in the task.

Table 5. Summary of Observations for Activity 4: Make Recommendation

ID	Overall ELP	Issues Observed				Activity 4 Make a Recommendation (out of 18)*
		Difficulty	Usability	Engagement	Language	
Arg1CL	I	2	0	1	0	14 (78%)
Arg2CL	I	2	0	2	0	14 (78%)
Arg3CL	I	2	0	1	1	6 (33%)
Arg4CL	I	0	2	0	0	18 (100%)
Arg5CL	B	1	1	2	0	6 (33%)
Arg6CL	B	2	0	1	0	12 (67%)
Arg7CL	B	0	0	0	0	8 (44%)
Arg8CL	EI	2	0	2	0	8 (44%)
<i>Mean Score</i>	--	<i>1.38</i>	<i>0.38</i>	<i>1.13</i>	<i>1.13</i>	<i>10.75 (60%)</i>

Table note: 0 = no, 1 = partial, 2 = yes. * = Students did receive support from the interviewer in responding to the questions in the task.

Table 6. Summary of Observations for Activity 5: Establish Criterion

ID	Overall ELP	Issues Observed				Activity 5 Establish a Criterion (score out of 32)*
		Difficulty	Usability	Engagement	Language	
Arg1CL	I	--	--	--	--	4 (13%)
Arg2CL	I	2	0	2	0	0 (0%)
Arg3CL	I	2	0	1	0	28 (88%)
Arg4CL	I	0	2	0	2	24 (75%)
Arg5CL	B	2	0	2	0	12 (38%)
Arg6CL	B	2	0	2	0	18 (56%)
Arg7CL	B	2	1	1	0	18 (56%)
Arg8CL	EI	2	0	2	0	28 (88%)
<i>Mean Score</i>	--	<i>1.50</i>	<i>0.38</i>	<i>1.25</i>	<i>0.25</i>	<i>16.50 (52%)</i>

Table note: Arg1CL ran out of time and did not complete the activity. 0 = no, 1 = partial, 2 = yes.

* = Students did receive support from the interviewer in responding to the questions in the task.

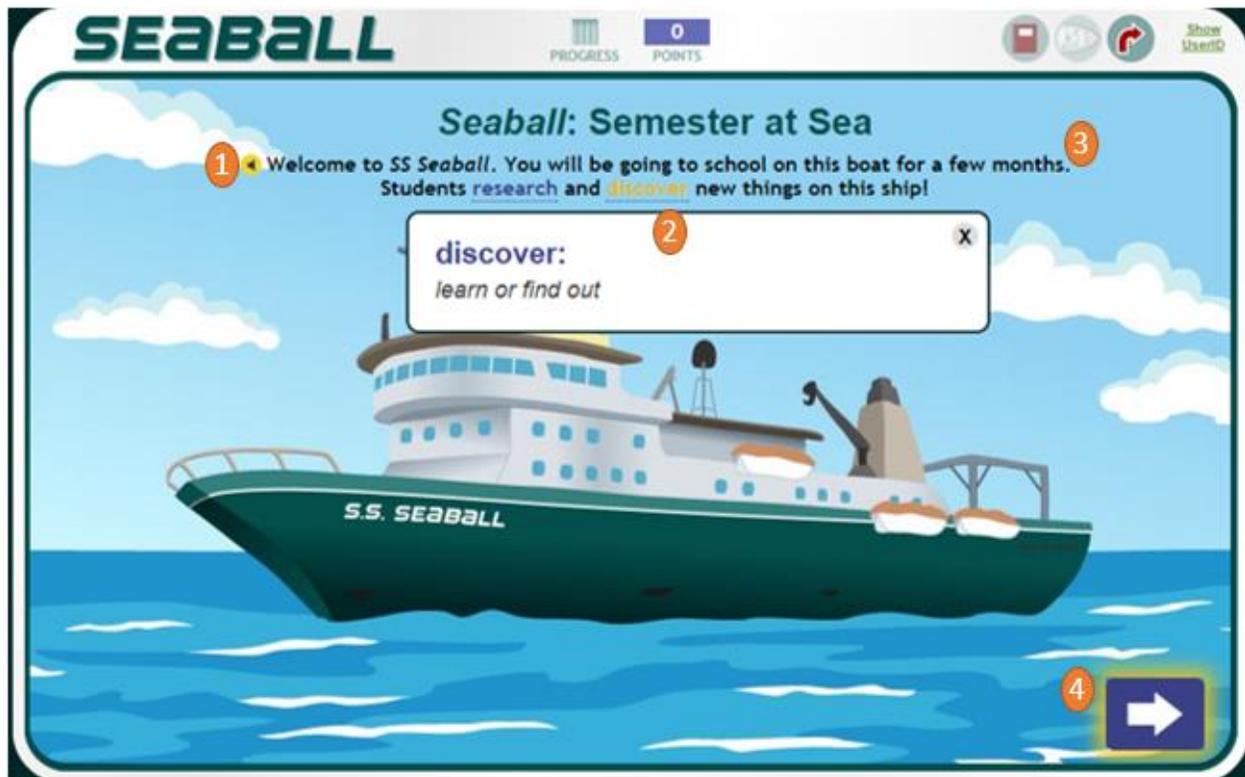


Figure 1. Revised Introduction screen for the *Seaball* task. The redesigned supports for EL students are as follows: 1 = read aloud component for directions. 2 = pop up glossary (activated). 3 = linguistically modified text. 4 = enhanced visual cue to support task navigation. Not shown are the shortened task and vocabulary-building activity.

Appendix A

EL Argumentation: Construct-Relevant Vocabulary List

Construct Relevant Vocabulary (in order of appearance in the task)

Argumentation/argument/argue
evidence
reason/reasoning
opinion/point of view/viewpoint
agree/disagree
concern/concerned
In my opinion/I think/I believe
controversy

Context Relevant Vocabulary (in order of appearance in the task)

ship	cooler
Semester at Sea	notebook
voyage	map
bridge	library
student council	cafeteria
president	student lounge
C-Store	fitness center
junk food	hallway
interview	snacks
ban/allow	captain
meeting room	characters
calorie	