## Assessing 21st-Century Skills With Performance Tasks: The Five-Year Journey of a Large School Division

*This article describes the development of the Integrated Performance Task (IPT), a series of assessments designed to measure critical thinking, problem solving, and written communication. The IPT has been administered to students in grades four and seven in Virginia Beach City Public Schools since 2010. The evolution of the IPT is explained in four areas: rubric development, task creation, reviewing and editing, and scoring responses. The ideas presented in the article may be replicated or modified for developing performance tasks at the classroom, school, or division levels.*

**Douglas G. Wren, Ed.D**

*Doug Wren works with the Virginia Beach City Public Schools Department of Planning, Innovation, and Accountability as an Assessment Specialist. He is also an Assistant Adjunct Professor in the Department of Educational Leadership and Foundations at Old Dominion University. Prior to moving to Virginia Beach, Doug taught elementary school and worked in the Department of Research and Evaluation with the DeKalb County School District in Georgia. He can be reached at* **dgwren@vbschools.com**

## Introduction

After years of relying primarily on multiple-choice test results to make important decisions about American students, teachers, and schools, the general attitude towards this type of assessment appears to be shifting. There is widespread disillusionment with the prolific use of state-mandated, multiple-choice tests brought about by *No Child Left Behind* (Darling-Hammond, 2014). Only 26% of over 10,000 teachers surveyed for *Primary Sources 2012: America's Teacher on the Teaching Profession* (Scholastic & Bill & Melinda Gates Foundation, 2012) maintained that the results of standardized tests accurately reflect student

achievement. A more recent survey, the *2013 Phi Delta Kappan/Gallup Poll of the Public's Attitudes Toward the Public Schools* (Bushaw & Lopez, 2013), found that three-fourths of Americans believed increased standardized testing in their local schools either hurt instruction or had no effect.

Performance assessments are making a strong comeback after nearly disappearing from the state and national scene around the turn of this century (Stecher, 2010). In 2014, the Virginia General Assembly passed legislation to eliminate several Standards of Learning (SOL) tests and gave local school boards the option to use authentic performance assessments and integrated assessments to ensure students are learning the content in these areas (Virginia Board of Education, 2014). The terms "performance assessment" and "performance task" are sometimes used interchangeably. However, Stecher (2010) saw performance assessment as a compilation of performance tasks and defined performance task as "a structured situation in which stimulus materials and a request for information or action are presented to an individual, who generates a response that can be rated for quality using explicit standards" (p. 3). This definition is suitable for the performance tasks that are the topic of the present article.

The purpose of this article is to describe the evolution of the *Integrated Performance Task* (IPT), a series of locally-developed performance tasks administered to Virginia Beach City Public Schools (VBCPS) students. This account explains the procedures that were used and the lessons we learned during this five-year excursion through largely uncharted waters.

## Starting With a Strategic Plan

In 2008, the Virginia Beach School Board adopted a new strategic plan, *Compass to 2015.* According to the plan, the primary focus for VBCPS would be on

"teaching and assessing those skills our students need to thrive as 21st century learners, workers, and citizens. All VBCPS students will be academically proficient; effective communicators and collaborators; globally aware, independent, responsible learners and citizens; and critical and creative thinkers, innovators and problem solvers" (Virginia Beach City Public Schools, 2008). (To learn more about Compass to 2015, visit **www.vbschools.com/compass/index.asp.**)

> **"All VBCPS students will be academically proficient; effective communicators and collaborators; globally aware, independent, responsible learners and citizens; and critical and creative thinkers, innovators and problem solvers"**

While we had plenty of existing tests to determine our students' academic proficiency (e.g., benchmark tests, final exams, SOL tests), measuring outcomes such as critical and creative thinking was somewhat of a tall order. Tony Wagner, a prominent author and founder of Harvard's Change Leadership Group, served as a consultant to VBCPS as we were making plans to implement *Compass to 2015*. Dr. Wagner told us about an innovative performance task, the *College and Work Readiness Assessment* (CWRA; Council for Aid to Education, 2007). The CWRA was designed to assess analytic reasoning and evaluation, problem solving, writing effectiveness, and writing mechanics, skills that paralleled some of our *Compass to 2015* student success outcomes. Following a promising field test, we began administering the CWRA to seniors at every VBCPS high school during the 2009-2010 school year. In order for students to see their CWRA results well in advance of graduation, we changed from a senior to a junior administration in 2011-2012. We are now in our fourth consecutive year of administering the CWRA to each student in every English 11

course across the division. (For more information about the CWRA, go to **www.vbschools.com/schools/testing/cwraParents.asp.**)

Once a test to assess students' 21st-century skills at our high schools had been put in place, we turned our attention to the elementary and middle school levels. Several VBCPS administrators attended a two-day Performance Task Academy sponsored by the Council for Aid to Education (CAE), the nonprofit organization that developed the CWRA. At the academy we learned more about the CWRA method of assessing higher-order skills by replicating the use of these skills in the real world. However, our real work began after we returned to Virginia Beach.

## Finding the Right People

With the high school CWRA serving as the model, our next steps were to (a) create performance tasks to measure students' critical-thinking, problem-solving, and written communication skills in grades four and seven, (b) generate rubrics to score the responses, and (c) develop and implement a viable scoring process. My colleague—a coordinator in the VBCPS Department of Teaching and Learning (T & L)—and I shared the responsibility for accomplishing these goals. One of our first and best decisions was to seek out a number of accomplished teachers in our division to assist with the work ahead. Our final list included over 40 teachers from nine elementary schools, nine middle schools, and two high schools. All of these teachers were recommended by principals and central office administrators. The group comprised representatives from the four core subject areas as well as other teachers whose specialties included family and consumer science, gifted education, reading, special education, and technology.

We invited these teachers to an informational meeting of the Compass to 2015 Assessment Development Team in November 2009. Staff who attended the meeting viewed a presentation that explained the role of assessments in the

new strategic plan. My colleague and I then provided an overview of the schema we believed would result in functional performance tasks for elementary and middle school students by the beginning of the next school year. Before leaving, attendees were asked to indicate the areas of the development process in which they felt most qualified to contribute. We envisioned the work to proceed as follows: rubric development, performance task creation, reviewing and editing, and scoring responses. These four areas of work became our focus not just over the next several months, but for the next five years.

## Developing the Rubrics

Because we already had a guiding philosophy and framework for developing the performance tasks (i.e., the CWRA method), our initial efforts were on developing rubrics for the three skills our performance tasks would measure—critical thinking, problem solving, and written communication. We decided ahead of time that the same general rubrics would be used for both fourth and seventh grades. Scoring guides specific to each performance task would supplement the general rubrics to help teachers score students' responses.

My colleague and I chose to label our rubric levels Novice (level 1), Emerging (level 2), Proficient (level 3), and Advanced (level 4). These were the same labels used in an earlier *Compass to 2015* project, the "VBCPS Continuum of 21st Century Skills." (The continuum is at **www.vbschools.com/compass/pdfs/ VBCPSContinuum.pdf**.) Novice level responses would suggest serious student deficits in the skills scored at level 1. Responses at the Emerging level would indicate that students required less attention to get them to the Proficient level, which was where we wanted them to be. The bar would be set much higher for the Advanced level; only responses that went well above the Proficient benchmark would be scored at level 4.

Our first official working group consisted of several teachers from our list of 40 as well as three T & L coordinators with experience in rubric development. First, we created operational definitions of critical thinking, problem solving, and written communication. Next, the group divided into subgroups and drafted three rubrics. The first rubric further operationalized critical thinking (CT) into three elements and provided descriptions of what student responses would look like at each level for each element (i.e., CT1, CT2, and CT3). Likewise, the second and third rubrics provided descriptions at each level for problem solving (PS) and written communication (WC), respectively, but there were only two elements for each of these skills (i.e., PS1, PS2, WC1, WC2). We soon realized the oppressiveness of having to score seven different elements on three rubrics for each response, so we scaled back to four elements—two for CT and one each for PS and WC—and combined them into a single rubric. Eventually, CT2 was eliminated. Our current rubric contains the three elements of CT, PS, and WC. Table 1.1 below illustrates the final rubric elements.

**TABLE 1.1** *IPT Rubric Elements*

| Element | Operational Definition |
|---|---|
| Critical Thinking | Decides if the information is correct and believable. |
| Problem Solving | Makes a choice and gives reasons for the choice. |
| Written Communication | Presents information and ideas that are clear, organized, detailed, and written for the intended audience. |

## Creating the Performance Tasks

In order to engage as many students as possible, performance task situations should be realistic, meaningful, and age-appropriate. Savvy educators realize

students tend to be more interested in instruction with direct connections to their world, rather than learning about something they perceive as not being relevant. Along these lines, we utilized GRASPS (Wiggins & McTighe, 2005), a performance task design framework that stands for Goal, Role, Audience, Situation, Product or Performance, and Standards. (See the appendix for the GRASPS frameworks that describe our original tasks.)

Although we considered possible performance task situations for all elementary and middle school grades, we needed a starting point. Grades four and seven were selected for the same reason we moved CWRA testing from grade 12 down to grade 11—to give students, parents, and school staff the opportunity to view students' responses and scores while the students were still at the same school.

In addition to a description of the situation, the fourth- and seventh-grade performance tasks include several documents to help students arrive at a decision. The fourth-grade situation requires students to choose between two ways to improve their health; the accompanying documents in the booklet are

**Virginia Beach IPT**

a government fact sheet, a news story with a bar graph, and an advertisement. The situation for seventh-grade students involves a controversial mall chaperone policy, and the booklet contains a news story with a line graph, an advertisement, a research brief, and a social media site complete with comments. Because analyzing and interpreting all of this information requires an amalgamation of knowledge and skills, we named our new assessment the Integrated Performance Task (IPT).

## Reviewing, Validating, Administering, and Revising

Every new and existing IPT undergoes changes before and after being administered to all fourth or seventh graders across the division. Systematic revisions are based on feedback from various sources through several processes, which are summarized below.

- **Expert Reviews** – Before field-testing an IPT to a sample of students, we ask a number of experts to carefully read and scrutinize the situation and its accompanying documents. Expert judgment is vital to obtaining evidence of content validity. The process involves appraising "the relevance and representativeness of the test content in relation to the content of the behavioral or performance domain about which inferences are to be made" (Messick, 1990, p. 5). Our experts have included veteran VBCPS teachers and staff from the Office of Programs for Exceptional Children, the Office of Student Assessment, and T & L.

- **Field Tests** – Some of the most insightful commentary we receive comes from students who have taken the IPT. Preceding a field test, the students understand (a) they have been selected to try out a new assessment, (b) they will not receive a grade, and (c) the purpose of the field test is to make the assessment better for other students who will take it later. In short, the

students are evaluating the test instead of the test evaluating the students. These students are encouraged to write their thoughts directly on their IPT booklets so we can review their comments later. After field-testing an IPT with an entire class, we conduct focus groups to draw out more information from the students.

- **Think-Alouds** – Besides field-testing classes, we administer new IPTs to students in one-on-one, "think-aloud" sessions. These sessions involve teacher-recommended students willing to verbalize every thought as they take the test. Each student is paired up with a VBCPS staff member with good note-taking skills and a thorough knowledge of the IPT. Think-alouds can be painstakingly long processes that occasionally yield little information, but we believe they are worthwhile because of the insights we can obtain through no other means.

- **Full Administrations** – Although expert judgments, field tests, focus groups, and think-alouds can generate plenty of information to improve the IPT, there is no substitute for a large-scale administration to obtain comprehensive feedback. When concerns about the IPT come from more than one school, we recognize the legitimacy of the concerns and act accordingly. The best example of this occurred after the first divisionwide IPT administration in fall 2010. While my colleague and I were training teachers to score IPT responses, we learned that having students respond to five open-ended questions (i.e., prompts) was too much; it was equally demanding for a teacher to read five responses and score four different elements for each student. For the spring IPT, we reduced the number of prompts to three and aligned the prompts with specific rubric elements. Prompt 1 responses were scored exclusively for CT1, Prompt 2 responses were scored for CT2 only, and responses to Prompt 3 were scored independently for PS and for WC. (Writing ability is not evaluated
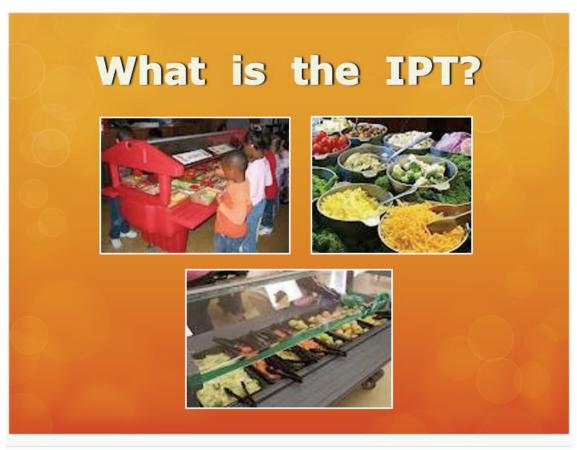
when scoring a student's problem-solving ability and vice versa.)

# Further Improvements

We have made other improvements to various aspects of the IPT to make the tasks more student friendly. Four areas of IPT enhancement are summarized as follows.

- **Readability** – To minimize the effect of reading ability, efforts have been made to ensure that each IPT is written at an appropriate reading level. Furthermore, the fourth-grade situation and documents are read aloud to students as they follow along in their booklets. Seventh-grade examiners read the situation to students before they begin the test. Because the IPT is not a reading test, students at both grade levels are told they may have words or sentences read aloud to them by an examiner or proctor at any time.

- **Glossaries** – The IPT is not a vocabulary test. A comprehensive glossary appears at the back of each seventh-grade IPT booklet. A similar glossary is included in the fourth-grade teacher directions. Students are informed before they begin the IPT that they may ask an examiner or proctor to read the definition of a word or term if they get stuck. Seventh graders can look up definitions on their own.

- **Rubric Explanation** – Students should know in advance how they will be scored on any performance task. The IPT rubric is included in every IPT booklet. Besides simplifying the language in the rubric, we developed short, kid-friendly PowerPoint presentations to explain the purpose of the IPT, its elements, and the rubric levels. Teachers are required to review the rubric and they are encouraged to answer questions to further clarify to students how responses will be scored.

- **Fairness** – According to the Standards for Educational and Psychological

Testing (AERA, APA, & NCTM, 2014), "fairness is a fundamental validity issue and requires attention throughout all stages of test development and use" (p. 49). A good example of a fairness concern was brought to our attention during the first division-wide administrations of the IPT. An assistant principal of a school with many low-income students pointed out that some of her fourth graders had never seen a salad bar. Before this IPT was administered again, we added two slides to the IPT PowerPoint presentation picturing and explaining outdoor fitness courses and salad bars. We later changed "salad bar" to "fruit and salad bar" after several teachers told us some students disliked salad but liked fruit.



What is the IPT?

There is information about a fruit and salad bar in the fall IPT.

A fruit and salad bar is a table or counter where people can choose the food they want.

Fruit and salad bars usually have fresh lettuce, fruits and vegetables, cheese, and other food.

You can see fruit and salad bars in some restaurants, cafeterias, and school lunchrooms.

## Scoring Responses

Prior to the first full-scale administrations during the 2010-2011 school year, we decided that fourth- and seventh-grade students would take two different IPTs annually. Although the spring IPT situation would not be identical to the situation that students had seen in the fall, the types of documents and prompts would be similar for both IPTs at each grade level. The fall IPT allows students to experience a low-stakes performance task while generating data for formative use in the classroom. Many VBCPS teachers give students the opportunity to review and reflect on their IPT responses at the end of the testing window. Black and Wiliam (1998) stated that "self-assessment by pupils, far from being a luxury, is in fact an essential component of formative assessment" (p. 6). The spring, summative IPT is also a low-stakes assessment in that the results are not used to evaluate teachers or make important decisions about students (Popham, 2001). Individual students' scores are available for their parents to view online; aggregate results are used to gauge progress on the *Compass to 2015* student success outcomes at the school and division levels.

For the fall IPT, principals are advised to involve all instructional staff—not just fourth- and seventh-grade teachers—in the scoring process. The Office of Student Assessment provides schools with detailed scoring guides and PowerPoint presentations to train teachers to score IPT responses. If requested, a staff member from the central office will conduct scorer training sessions at schools. After scoring fall IPT responses, each elementary and middle school must submit a form to the VBCPS Department of School Leadership explaining how teachers will use their IPT results to inform instruction.

While responses to the fall IPT are scored locally, spring IPT responses are scored centrally by a trained cadre of teachers. Every summer the cadre comprises over 100 vetted, ten-month VBCPS employees, including a number of teachers from

our original list of 40. Each IPT scorer must attend a full day of training before being allowed to score responses independently. The scoring process for the spring IPT involves independent ratings by two trained scorers. An expert scorer provides a third and deciding score if the first two scores do not match.

Our first experience with the scoring cadre in 2011 taught us invaluable lessons, which are described in the remainder of this section.

- **Training** – For the first summer we scheduled four weeks to score IPT responses. Each week began with a one-day training session. We quickly learned that the revolving door method of training and scoring (i.e., a new group starts each week) was a very bad idea. A major issue in scoring performance tasks is consistency, and achieving an acceptable level of interrater agreement can only be realized if (a) scorers are given the same extensive training on interpreting the rubric, and (b) they consistently apply this interpretation when scoring each and every response. To alleviate this problem, we began conducting training on the first day of a single three-week scoring session in subsequent years.

- **Scoring Assignments** – Another mistake we made during our first scoring adventure was having teachers score multiple elements. This forced scorers to change mindsets when moving from Prompt 1 responses to Prompt 2 responses to Prompt 3 responses. From the second year on, we have trained teachers to score only one element for one grade level. By allowing teachers to become scoring specialists in a single area, we have attained a higher degree of interrater and intrarater agreement. We have calculated the overall percentage of perfect agreement between first and second scorers (i.e., both scorers independently gave the same response identical scores) for each group of scorers since 2012. From 2012 to 2014, the percentages for

these groups have ranged from a low of 68% agreement to a high of 84% agreement. As a common rule of thumb, values of 75% or greater indicate acceptable levels of agreement (Graham, Milanowski, & Miller, 2012). Rooms with agreement levels under 75% require additional training and calibration.

- **Personnel Management** – Managing scorers is another aspect that we shored up after the first summer of scoring. Not only were scorers trained separately by IPT element and grade level, each group was housed in a separate room. To supervise each room, we selected two teachers with extensive scoring experience and strong leadership skills. These room leaders have conducted training and daily calibration sessions, kept up with scored and unscored responses, provided third scores, retrained scorers as needed, and made hiring recommendations for returning scorers every year.

- **Data Management** – Data entry is handled initially by room leaders but by the end of the second week, a designated data team takes on full responsibility for compiling, verifying, and entering tens of thousands of scores for students who took the spring IPT. Effective management of the data room is key to the success of the entire summer scoring operation.

## Further Validation

Our expert reviews of the IPT provided content validity evidence for the IPT, but we wanted to acquire further evidence of validity. By correlating IPT scores with the scores of an established criterion (e.g., a reputable critical-thinking test), we could determine whether the IPT was measuring the same general construct that the criterion test was measuring. This is known as a criterion validation study. The results of this study would indicate if IPT results were a valid measure of critical thinking among students in grades four and seven.

During spring 2013 and spring 2014, we administered age-appropriate

versions of the *California Critical Thinking Skills Test* (CCTST; Facione & Gittens, 2012) to fourth- and seventh-grade students during the IPT testing window. As shown in Table 1.2 below, all of the correlations between the IPT element scores and CCTST overall scores were statistically significant. Correlations of .30 or greater suggested that the IPT and the CCTST are measuring, to a moderate degree, the same general construct (Cohen, 1988). In other words, a student's scores on the IPT can be used to make valid inferences about the student's critical-thinking abilities. (As a reminder, a student's performance on a single assessment does not provide enough information to make consequential decisions about the student.)

**TABLE 1.2** *Correlations Between IPT Element Scores and CCTST Overall Scores*

| Level | n | CT1 | CT2 | PS | WC |
|-------|-----|------|------|------|------|
| Grade 4 | 207 | .36* | .29* | .33* | .35* |
| Grade 7 | 395 | .37* | .41* | .31* | .42* |

*Note. *p < .01, 2-tailed.*

## Conclusion

Long before last year's legislation to amend SOL testing, educators across Virginia have recognized that assessments with open-ended questions are more effective than multiple-choice tests for evaluating their students' depth of understanding. Leaders at the national, state, local, and school levels should be applauded for decisions that encourage teachers to depart from rote instructional methods driven by a multiple-choice mentality. One such decision was the implementation of the IPT in Virginia Beach. With commitment and effort, performance tasks similar to the IPT and CWRA can be developed in Virginia schools to measure SOL objectives as well as 20th-century skills. Performance task developers are urged to heed the lessons of the Virginia Beach

IPT and consider alignment, fairness, readability, scoring consistency, and other vital issues related to this type of assessment as they journey ahead. (For more information about the IPT, visit **www.vbschools.com/schools/ testing/IptFaq.asp.** Literature reviews by this author on the topics of performance assessment and formative assessment are at **www.vbschools. com/accountability/research_briefs/ResearchBriefPerfAssmtFinal. pdf** and **www.vbschools.com/accountability/research_briefs/ researchbriefformassmtfinal.pdf**.)

## References

American Educational Research Association, American Psychological Association, & National Council on Measurement in Education. (2014). *Standards for educational and psychological testing.* Washington, DC: American Educational Research Association.

Black, P., & Wiliam, D. (1998). Inside the black box: Raising standards through classroom assessment. *Phi Delta Kappan, 80(2),* 1-10. Retrieved from **csi.idso.eportalnow.net/uploads/1/1/3/2/11323738/inside_the_ black_box_1998.pdf**

Bushaw, W., & Lopez, S. (2013). The 45th annual PDK/Gallup poll of the public's attitudes toward the public schools. *Phi Delta Kappan (95),* 9–25. Retrieved from **pdk.sagepub.com/content/95/1/8.full.pdf+html**

Cohen, J. (1988). *Statistical power analysis for the behavioral sciences* (2nd ed.). Hillsdale, NJ: Lawrence Erlbaum Associates.

Council for Aid to Education. (2007). *College and Work Readiness Assessment* [Measurement instrument]. New York, NY: Council for Aid to Education.

Darling-Hammond, L. (2014, January/February) Testing to, and beyond the Common Core. *Principal, 93(3)*, 8-12. Retrieved from **www.naesp.org/ principal-januaryfebruary-2014-assessments-evaluations-and-data/ testing-and-beyond-common-core**

Facione P. A., & Gittens, C. A. (2012). *California Critical Thinking Skills Test: M series test manual.* Millbrae, CA: Insight Assessment/California Academic Press.

Graham, M., Milanowski, A., & Miller, J. (2012). *Measuring and promoting inter-rater agreement of teacher and principal performance ratings.* Washington, DC: Center for Educator Compensation Reform. Retrieved from **cecr.ed.gov/pdfs/ Inter_Rater.pdf**

Messick, S. (1990). *Validity of test interpretation and use.* Princeton, NJ: Educational Testing Service, (ERIC Document Reproduction Service No. ED395031). Retrieved from **http://files.eric.ed.gov/fulltext/ED395031.pdf**

Popham, W. J. (2001, March). Teaching to the test. *Educational Leadership, 58(6),* 16-20. Retrieved from **www.ascd.org/publications/educational-leadership/ mar01/vol58/num06/Teaching-to-the-Test%C2%A2.aspx**

Scholastic, & Bill & Melinda Gates Foundation. (2012). *Primary sources 2012: America's teacher on the teaching profession.* New York: Scholastic. Retrieved from **www.scholastic.com/primarysources/pdfs/Gates2012_full.pdf**

Stecher, B. (2010). *Performance assessment in an era of standards-based educational accountability.* Stanford, CA: Stanford University, Stanford Center for Opportunity Policy in Education. Retrieved from **scale.stanford.edu/ system/files/performance-assessment-era-standards-based-educational-accountability.pdf**

Virginia Beach City Public Schools. (2008). *Compass to 2015: A strategic plan for student success.* Retrieved from **www.vbschools.com/compass/index.asp**

Virginia Board of Education. (2014, September 18). *Guidelines for local alternative assessments: Developed in response to 2014 acts of assembly.* Retrieved from **doe.virginia.gov/testing/local_assessments/guidelines_for_local_ alternative_assessments.pdf**

Wiggins, G. P., & McTighe, J. (2005). *Understanding by design* (2nd ed.). Alexandria, VA: Association for Supervision and Curriculum Development.

# Appendix A

## Grade 4

**G**oal: Decide which project would be better for improving students' health at a school

**R**ole: A fourth-grade student at the fictional Smith Elementary School

**A**udience: Mr. Beach, the principal of Smith Elementary School

**S**ituation: A local business is donating money to the school to pay for only one of two projects—an outdoor fitness course or a salad bar

**P**roduct: A persuasive letter to Mr. Beach recommending one of the projects

**S**tandards for success: Described in the rubric

## Grade 7

**G**oal: Decide whether to continue a policy restricting minors' access to the fictional Beach Mall

**R**ole: A teenager serving on a committee formed by Beach Mall officials

**A**udience: Beach Mall officials

**S**ituation: Beach Mall recently implemented a chaperone policy to improve safety, security, and profits for the mall's businesses.

**P**roduct: A written recommendation to the mall officials with rationale for continuing or discontinuing the chaperone policy

**S**tandards for success: Described in the rubric