



The Impacts of a Standards-Based Grading System Emphasizing Formative Assessment, Feedback, and Re-Assessment: A Mixed Methods, Cluster Randomized Control Trial in Ninth Grade Mathematics Classrooms

Steven L. Kramer^a , Michael A. Posner^b, Alexander S. Browman^c ,
Nancy R. Lawrence^d, Jennifer Roem^e and Kathleen Krier^a

^aResearch, The 21st Century Partnership for STEM Education, Wayne, PA, USA; ^bDepartment of Mathematics and Statistics, Villanova University, Villanova, PA, USA; ^cDepartment of Psychology, College of the Holy Cross, Worcester, MA, USA; ^dHuman Development and Family Studies, Colorado State University, Fort Collins, CO, USA; ^eDepartment of Epidemiology, Johns Hopkins Bloomberg School of Public Health, Baltimore, MD, USA

ABSTRACT

We investigated the impact in ninth-grade mathematics classrooms of Proficiency-based Assessment and Reassessment of Learning Outcomes (PARLO), a standards-based grading system. Key components of PARLO are basing student final grades on the number of learning outcomes on which the student is high-performance or proficient, providing students with formative feedback, and encouraging students to reassess for full credit after further study. Our mixed-methods study employed a cluster randomized control trial with 2,736 participating ninth graders at 14 Treatment and 15 Control schools. Data included student achievement tests, interviews with 35 teachers, and student surveys. The program improved student performance on end-of-course algebra and geometry tests by a statistically significant 0.33 SD but did not impact students' value of or expectancies for success in mathematics. However, treatment effects on mathematics performance were moderated by these psychological antecedents of motivation, such that students with higher math expectancies and value benefited more from the treatment. Furthermore, teacher interviews suggested that PARLO may have also had positive effects on growth mindsets, mastery goals, autonomy, and relatedness.

ARTICLE HISTORY

Received 22 February 2022


Revised 25 October 2023

Accepted 27 October 2023

KEYWORDS

Standards-based grading;
high school mathematics;
growth mindsets;
expectancy value theory;
self-determination theory

CONTACT Steven L. Kramer  skramer@21pstem.org  Research, The 21st Century Partnership for STEM Education, Wayne, PA, USA.

 Supplemental data for this article can be accessed online at <https://doi.org/10.1080/19345747.2023.2287594>.

© 2023 The Author(s). Published with license by Taylor & Francis Group, LLC

This is an Open Access article distributed under the terms of the Creative Commons Attribution-NonCommercial-NoDerivatives License (<http://creativecommons.org/licenses/by-nc-nd/4.0/>), which permits non-commercial re-use, distribution, and reproduction in any medium, provided the original work is properly cited, and is not altered, transformed, or built upon in any way. The terms on which this article has been published allow the posting of the Accepted Manuscript in a repository by the author(s) or with their consent.

Introduction

Among the most central and deeply rooted features of most American classrooms is a grading system that is designed to rank, sort, and evaluate students as *learners*, but is not optimally designed to *help students learn*. In the late 19th and early 20th centuries, many educators sought to build a “meritocratic” system that would give students access to a common elementary education, with students encouraged to persist according to their interests and abilities. Specifically, the system was designed so that more naturally talented and quicker learning students would receive higher grades than those with less natural talent or who learned more slowly, and only those with high grades would be encouraged or even permitted to progress to the next level of the education system (Farrington & Small, 2008; Schneider & Hutt, 2014).

Thus, the grading system that is still typical in U.S. classrooms today was adopted at a time when a primary goal of public education was to sort students based on their supposed inherent intellectual abilities. As a result, the grading system does not optimally support the theoretical goal of the modern education system: learning. Specifically, classroom grading is built around *summative assessments* that typically only provide students with a limited amount of time to learn a topic and a limited number of opportunities to demonstrate how much they have learned. After assessment, students are often not provided with additional opportunities to improve their skills or to demonstrate that they have increased their learning. Further, early assessments are typically averaged with later assessments. This can discourage persistence in the face of initial difficulty because no matter how well students eventually learn the material they will be evaluated partly based on their initial difficulties.

A Proposed Improvement: Adding Formative Assessment

To better support student learning, educators in the 1960s began recommending formative assessment as a supplement to summative assessment (Bloom, 1968). In contrast to the traditional system, where assessments are used solely to provide a final, summative judgment of student performance, formative assessment uses assessments intermittently to provide teachers and students with feedback about each student’s progress. This enables teachers and students to adapt their teaching and learning strategies to help each student progress based on their individual needs. Numerous meta-analyses have shown that formative assessment has significant positive implications for student engagement, learning, and performance, especially for previously low-performing students (e.g., Black & Wiliam, 2010; Kingston & Nash, 2011).

To date, most formative assessment programs that have been implemented, especially in mathematics, have been teacher-facing. This approach to formative assessment provides teachers with feedback about their students’ current level of proficiency, allowing the teachers to adapt their teaching strategies to meet the individual needs of their students (e.g., Supovitz et al., 2018). However, formative assessment can also be student-facing, wherein assessments are used to provide students with feedback about their own current level of proficiency and the nature of their mistakes so they can adapt their learning and become more proficient.

Numerous educators have called for more widespread implementation of student-facing formative assessment (e.g., Leahy et al., 2005; Wiliam, 2011). Mills and Silver (2018, p. 180) articulated the argument for mathematics: “Teachers on their own cannot make students learn mathematics; students must become partners in their own learning in order for effective teaching and learning to occur.” Supporting student agency by incorporating student-facing formative assessment may be especially efficacious for middle and high school students, an age when students are exploring their independence and have a particular need for status and respect (Yeager, 2017; Yeager et al., 2017).

Nonetheless, prior research on student-facing formative assessment in math class has been limited, mostly focusing on improving the feedback that students receive about their work. For example, Murphy et al. (2020) reported results of a cluster randomized control trial (RCT) that evaluated “Assistments,” an online program that provides students immediate feedback on the correctness of homework problems. “Assistments” increased the mathematics achievement of seventh graders by a statistically significant 0.2 standard deviations (SD). Programs like “Assistments,” are philosophically compatible with and potentially complementary to, but differ from, the program we investigated, which focused on reengineering the grading system to better support student learning.

Standards-Based Grading: The PARLO Assessment System

The present work examines the effects of the *Proficiency-based Assessment and Re-assessment of Learning Outcomes (PARLO)* system. The system was developed by a partnership between Dylan Wiliam, a pioneer researcher in formative assessment, and the Math Science Partnership of Greater Philadelphia (MSP-GP), an NSF-funded project that brought together 13 institutions of higher education and 46 Pennsylvania and New Jersey school districts to improve grades 6-12 mathematics and science education. It was partly inspired by a program at the Young Women’s Leadership Charter School (YWLCS), a public charter high school in Chicago that changed its grading system to allow students to reassess for full credit whenever they had mastered classroom content. The goal was to promote students’ agency by fostering a sense of partnership with their teachers in ensuring that they mastered the big ideas from each course. As a result of these changes, YWLCS regularly achieved the highest graduation rate of any nonselective public school in Chicago, despite serving a similar student body as neighboring schools (mostly low-income Black and Latina students; Farrington & Small, 2008).

The PARLO assessment system is designed to be implementable in a single teacher’s classroom, embedded in a school that may have a traditional American grading system. Consequently, PARLO teachers do assign final summative grades. However, PARLO reengineers the classroom assessment system by better integrating summative assessment and student-facing formative assessment. Specifically, while students receive summative scores on quizzes and assignments, such scores are not recorded indelibly to be used in a weighted average that determines final grades. Rather, the teacher uses such assessment evidence both to evaluate the extent to which a student is proficient in each of the course’s learning outcomes *at that particular moment in time*, and then to provide students with personalized feedback designed to guide further learning. Students are then

given opportunities to do further work, at home or in school, and to be reassessed for full credit. In other words, summative assessments become formative tools designed to promote further learning, instead of merely yardsticks to measure the student's efficiency as a learner.

The first step in implementing PARLO is clarifying and sharing learning intentions and success criteria (Wiliam, 2011). To do this, the teacher organizes their course instruction around 10-15 learning outcomes per semester, which together define the material to be mastered. These learning outcomes are then shared with students and their parents so that they can be partners in the student's educational progress. An important characteristic of PARLO that distinguishes it from similar programs such as Mastery Learning (Bloom, 1968) is the way it defines success criteria. In addition to establishing criteria for "proficient" performance, teachers also share criteria for "high performance." Achieving high performance on a learning outcome is intended to challenge talented students, while potentially being achievable through hard work by nearly all students. In practice, teachers in the current study used Webb's Depth of Knowledge (Webb, 2002) to define success criteria, with high-performance requiring students either to solve problems at a greater depth of knowledge (application or strategic thinking), or else to tutor a fellow student and successfully bring them from not-yet-proficient up to proficient.

The second step in implementing PARLO employs the central idea of student-facing formative assessment: eliciting evidence of learner's achievement and providing feedback that moves learning forward (Wiliam, 2011). PARLO teachers use short quizzes, end-of-lesson written "exit tickets," notes from observing students working in groups, and other formative assessment techniques to provide students feedback about how well they are progressing toward proficient or high performance on each learning outcome. Note that other student-facing formative assessment programs like Assistments (Murphy et al., 2020) implement this step. In future iterations of PARLO, teachers might use programs like Assistments to support PARLO implementation.

The remaining two steps in implementing PARLO focus on changing the way teachers assign summative grades. As will be explained in more detail below, the changes to summative grading were the key features differentiating the Treatment from the Control schools in the current study. These two steps are:

1. Reassessment for full credit after further learning. Student grades are not averaged over the semester; instead, students are rated on each learning outcome as not-yet-proficient, proficient, or high-performance based on the best work they can show by the end of the semester. While students can reassess for full credit, before reassessment they must first engage in further learning, completing activities such as:
 - a. Error logs, wherein students explain what they did wrong on a problem, rework it, and describe what they would need to remember so as not to make the same mistake again.
 - b. Remediation plans, or contracts with students about activities to be undertaken before reassessment.

- c. Flashback days, or in-class opportunities for students to work individually or together to revisit learning outcomes and learn material at proficient or high-performance level.
2. Final Grades Based on Learning Outcomes. With project help, each school developed its own algorithm to determine a letter grade based on the number of Learning Outcomes scored Proficient and the number scored High Performance. Other factors sometimes used to assign grades, such as attendance, attitude, and homework completion, are not averaged as part of the final grade. Instead, they are viewed as means to the end of understanding course content.

At the time that we designed the intervention, a grading system with components like the four we described above was commonly known as “proficiency-based”—hence the name “PARLO.” However, terminology has since shifted, and PARLO-type grading systems have become known under the name *standards-based grading* (see, e.g., Marzano, 2010).

Theoretical Framework

We hypothesized that the PARLO system would positively impact student learning in two ways: by providing *additional time and opportunity for students to learn*; and by *encouraging students’ motivation/desire to learn*.

Opportunity to Learn

As noted above, traditional summative assessment systems encourage teachers to provide a limited amount of time for students to learn a topic and a limited number of opportunities to demonstrate how much they have learned. In contrast, PARLO teachers provide students with more opportunities to learn from their mistakes, improve, and be reassessed to demonstrate their increased proficiency. Further, student persistence is encouraged by giving them full credit for material learned, regardless of initial difficulties they may have had.

Motivation

According to expectancy-value theory (Wigfield & Eccles, 2000), students’ motivation and classroom engagement are largely determined by the degrees to which they both *value* the material being taught and *expect to succeed* in class. Yet, studies have found that students exposed to an assessment system designed to compare and rank students report weaker feelings of intrinsic value (enjoyment of learning for its own sake), utility value (the belief that learning is useful and will have purpose and relevance in one’s life), and expectancies for success than those exposed to a system that provided students with opportunities to improve their skills or to demonstrate that they have increased their learning (Covington & Omelich, 1984; Haley, 2015; Sanchez et al., 2017), and that

these effects may be especially pronounced among students from less advantaged backgrounds (Jury et al., 2015; Smeding et al., 2013).

Additional Motivational Antecedents Addressed

This study was designed using concepts from the expectancy-value theory of motivation (Wigfield & Eccles, 2000). Additionally, after we had collected our data and begun to analyze it, we found concepts from three additional theories of motivation to be helpful in understanding what the teachers told us in interviews: growth mindset theory (Dweck, 2007), self-determination theory (Ryan & Deci, 2020), and achievement goal theory (Senko, 2016). Growth mindset theory suggests that students show greater motivation, engagement, learning, and performance in school, especially in the face of difficulty, if they believe that they can grow and improve their intelligence and academic skills through hard work (a growth mindset), versus if they instead believe that their ability levels are innate and unchangeable (a fixed mindset). Self-determination theory suggests that students will be more internally motivated to engage with their schoolwork when three needs are met: *competence* (roughly analogous to expectancies for success, described above), *autonomy* (the feeling that they have opportunities to make meaningful choices about their learning), and *relatedness* (the feeling that they are valued and belong in class). Finally, achievement goal theory suggests that students experience greater academic outcomes when they adopt *mastery goals*, i.e., engaging in schoolwork either to learn as much as possible, or else to meet self-determined standards of success in learning. As will be discussed below, each of these themes emerged in interviews with teachers and open-ended survey responses from students who were exposed to PARLO—a program that encouraged the capacity of students to improve their mathematics abilities (i.e., growth mindsets and mastery goals), and provided students with flexibility in how they could go about acquiring and demonstrating this improvement (i.e., self-determination). Below, we refer to feelings, attitudes, beliefs, and goals that contribute to student motivation and engagement as *motivational antecedents*.

Conceptual Model

Figure 1 displays our current conceptual model regarding how the PARLO system may influence students' academic experiences and learning. Because it was developed using the still-to-be described results of the present work, we return to the complete model in the Discussion section. For now, we wish to note the most important difference between the PARLO Treatment and the Control group in the present work: the two white boxes on the bottom left of the figure. Unlike the Control teachers, PARLO teachers were trained and expected to reassess students for full credit after further learning, and to base a student's summative grade on the number of proficient and the number of high-performance learning outcomes the student had mastered. The two gray boxes on the top left—organizing instruction by learning outcomes and providing formative assessment and feedback—are also necessary parts of the PARLO system and were implemented by teachers in our study. However, these two PARLO components were popular innovations already in place in the school districts where we conducted the study. Given this environment, the project delivered professional development to teachers in the

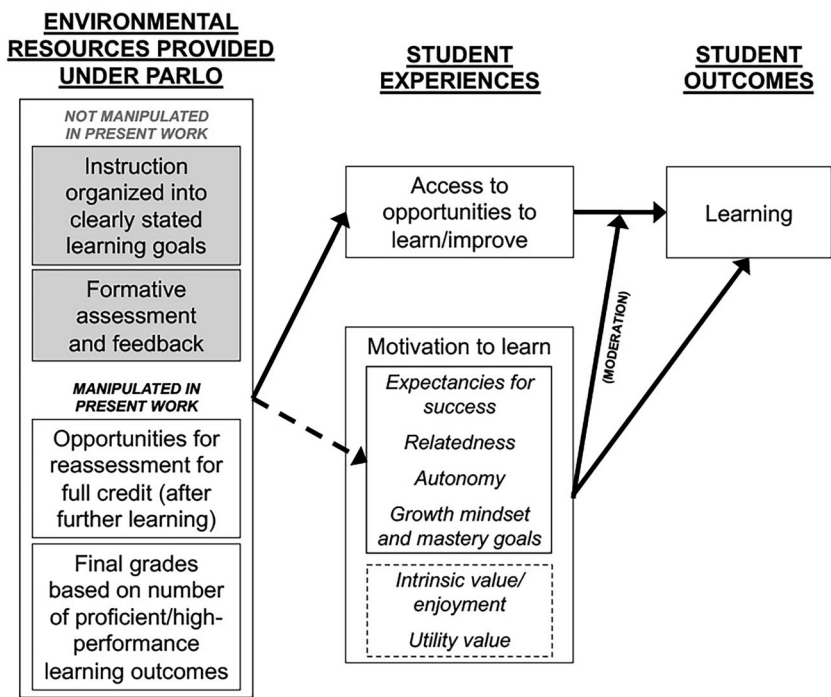


Figure 1. Conceptual model.

Control group to further assist them in implementing these two formative assessment-related practices. We anticipated that PARLO teachers would be more likely than Control teachers to implement the PARLO components displayed in the grayed-out boxes, but only because implementing the bottom two components encouraged them to do so. Thus, the current study tests the impact of experimentally adding the lower two components of the PARLO standards-based grading system—reassessment for full credit after further learning, and summative grades based on the number of proficient and the number of high-performance learning outcomes—to classrooms that were already endeavoring to implement the top two components.

We also note that the focus of the conceptual model, and of the current article, is on *changing teacher’s assessment practices*. Specific inputs, such as software tools made available to teachers and specific professional development provided, evolved over the course of the project, and will likely be implemented differently in the future. Consequently these “program inputs” are not displayed in the model. Our intent is to generalize results to future implementations of standards-based grading that are similar to PARLO. See our Conclusion section for ideas about possible future implementations.

Prior Research on Standards-Based Grading and PARLO

While numerous researchers have written about standard-based grading, published work has focused on case studies of small-scale implementation in single classrooms or schools (e.g., Farrington & Small, 2008). Consequently, researchers have decried the lack

of studies evaluating the practice (Marzano, 2010; Scarlett, 2018). The current study should help to fill that gap.

There have been two pilot studies of the PARLO system. First, Clymer and Wiliam (2007) implemented PARLO for a full school year in a single eighth grade science classroom. Average achievement on an end-of-class exam increased by 0.4 SD over the previous year's scores. In addition, student interviews indicated that these positive effects on performance may have emerged because students developed stronger growth mindsets. Second, Posner (2011) taught two sections of an introductory undergraduate statistics class for non-majors, nonrandomly employing PARLO in one section and traditional grading in the other. He found that the PARLO group scored significantly higher than the traditional grading group on measures of perceived intrinsic value, utility value, expectations for success, and motivation in statistics.

The Current Study

While these initial results were promising, the studies described were only small pilots that did not provide a randomized control-style investigation of the effects of the PARLO assessment system. The present work describes the results of such an investigation. Specifically, we employed a mixed-methods, cluster randomized control approach with school as the unit of analysis, to address the following two research questions: *What is the impact of the PARLO system on 9th grade students' expectancies for success, utility value, intrinsic value, long-term motivation, and academic performance in mathematics class? What are the mechanisms by which this impact occurs?*

In addition, while not part of our *a priori* hypotheses, our examinations of the qualitative data unexpectedly revealed that the relationship between the PARLO system and student motivation might be bi-directional: not only might PARLO improve student motivation, but students with higher initial motivation might also benefit more from the PARLO system. Consequently, we used our quantitative data to investigate an additional research question: *Is the impact of the PARLO system on student academic performance moderated by established antecedents of student motivation—specifically, students' value of and expectancy for success in mathematics?*

Method

This was a mixed-methods study, collecting qualitative and quantitative data throughout program implementation. Our data analysis followed a sequential explanatory design, performed in three phases. In the first phase, we conducted quantitative analyses to see whether the PARLO system had its predicted effects on motivation and on student achievement. In the second phase, we used qualitative data to look for explanations of the results we found in phase 1. Our qualitative analysis began inductively, thus allowing for the emergence of trends we did not hypothesize *a priori*. We summarized preliminary results as memos describing major themes observed and evidence for those themes. We then reviewed the memos in light of the quantitative findings, looking for explanations of those findings. The second phase yielded an unexpected result, indicating that student motivation might be moderating the PARLO system's effects on student

achievement. Consequently, we conducted a third phase, investigating whether the quantitative data we had collected could confirm the moderating effect.

Participants

We recruited urban, suburban, and rural schools; public, charter, and religious schools; schools from high- and low-performing districts; and schools with a wide variety of racial make-ups. In order to ensure reasonably high fidelity of implementation, we only assigned a school to the project if both the administration and the ninth-grade mathematics teachers agreed in advance that they would participate as assigned either in the treatment or in the control condition, and would maintain participation regardless of subsequent random assignment. All ninth-grade algebra and geometry teachers at each school were asked to participate. Teachers received a stipend in exchange for their participation.

We ultimately recruited two cohorts of schools: a cohort of 20 schools (14 public schools, 3 charter schools, and 3 Catholic all-girl schools) that participated during the 2010–11 and 2011–12 school years; and a second cohort of 15 schools (14 public schools from one large urban school district and 1 additional public school) that participated during the 2011–12 and 2012–13 school years. Cohorts participated for two years (i.e., two separate ninth grade classes) because implementing PARLO required significant changes to teachers' instructional practices that we anticipated might take more than a single year to take root.

Several schools dropped out of the study after randomization but before data collection was completed. The final sample of schools in the quantitative study included 14 PARLO and 15 control schools across both cohorts. One additional PARLO school provided qualitative data but dropped out before quantitative data collection could be completed. In all cases, school-level attrition was caused by a change in administration, which is unlikely to be related to the treatment. The overall attrition rate was 17%, and the differential attrition between treatment and control schools was 10.4%. At the student level, the overall attrition rate was 16.9% and the differential attrition was 3.9%. These rates are within acceptable standards when the intervention is unlikely to affect attrition (USDOE, 2022). See [Tables A1, A3, and A4 in the online appendix](#) for details about attrition.

Our Institutional Review Board (IRB) required active consent (i.e., parental opt-in) before administering our quantitative measures during the 2010–11 school year—the first implementation year for Cohort 1. However, many students failed to return a consent form, creating a risk of self-selection and bias in the data. As a result, our external reviewer and advisory board for this project determined that the quantitative data collected in 2010–11 could not be used for valid analyses. With the approval of our IRB, this issue was resolved beginning in the 2011–12 school year through the use of passive consent (i.e., parental opt-out). Thus, our quantitative analyses focus on data collected during academic years 2011–12 and 2012–13.

Across two years of data collection, 3,273 students completed algebra or geometry pretests and had available demographic data for the analysis—1,936 in treatment schools and 1,337 in control schools. Of these, 2,736 students provided complete quantitative

data, including baseline motivation and motivational antecedent scores, race, gender, and post-test scores, and thus were included in the analytic sample for our primary quantitative analysis. Of the total 2,736 students, 1,649 were from the classes of 38 teachers at the 14 PARLO treatment schools, and 1,087 were from the classes of 27 teachers at the 15 control schools. [Tables A5 and A6 in the online appendix](#) compare the demographic characteristics of treatment and control group students and teachers in the analytic sample.

Random Assignment and Conditions

In June 2010, all ninth-grade algebra and geometry teachers at all Cohort 1 schools who had agreed to participate attended three 6-hour days of professional development (PD) in the summer. After teachers completed this PD, project staff then randomly assigned participating schools either to the treatment or the control condition. Randomization was done in three blocks: charter vs. Catholic vs. public. The following year, we randomly assigned Cohort 2 schools to treatment or control conditions following a process that was identical but for two exceptions: the PD before random assignment lasted only two days, and blocking by school type was not necessary as all Cohort 2 schools were public. Participating teachers were aware of their condition assignment. See the online appendix for details about random assignment and for a depiction of the study timeline.

Control Condition

When this study was conducted, two formative assessment practices were coming into widespread use, i.e., becoming business-as-usual for many teachers in local school districts: sharing learning intentions and success criteria by organizing instruction around learning outcomes tied to state standards; and using formative assessment strategies such as providing frequent feedback. Both practices were necessary but not sufficient to implement PARLO. In order to further focus the study on the unique aspects of PARLO that were not already coming into popular use, before random assignment we provided all participating teachers professional development supporting these two widely used formative assessment practices. We had each teacher work with state standards for their course to translate them into 10 to 15 learning outcomes each semester that would be the focus of teaching in their classroom. We introduced Webb's Depth of Knowledge (Webb, 2002) to help them think about addressing each learning outcome at application and strategic-thinking levels. We used ideas from Wiliam (2011) to teach techniques for eliciting evidence from students and providing feedback to move their learning forward. The initial professional development lasted three days for Cohort 1 and two days for Cohort 2. Because the Control teachers were encouraged to implement the two aspects of PARLO reflected in the gray boxes of [Figure 1](#), the project often described the Control schools as a "limited treatment condition," to be contrasted with the "full treatment condition" that implemented all four aspects of PARLO.

Control teachers did not participate in any project-related PD beyond the two or three days that they participated in before randomization, and they were not obligated to modify their instructional practices. In addition to administering the student content exam and Attitudes Toward Mathematics Inventory described below, control teachers

completed questionnaire measures of the extent to which they used various practices in their teaching.

PARLO Treatment Condition

Teachers in treatment schools participated in three training opportunities not offered to control teachers. First, in August prior to first implementing the PARLO system in their classrooms, they attended three (for Cohort 1) or four (for Cohort 2) additional six-hour days of PD. Second, during the two years of project participation, they were given the opportunity to participate in monthly Professional Learning Community (PLC) meetings. The PLC meetings were held at five different locations and facilitated by teacher leaders who were mathematics content and PARLO experts. Third, they attended two six-hour days of follow-up PD during the summer between their two years of PARLO participation. Total participation time for treatment teachers in PD and PLC meetings amounted to about 88 hours over the two years. Seventy-nine percent of the treatment teachers attended 70% or more of the PD available to them.

The PARLO-specific PD focused on the two distinguishing characteristics of the PARLO standards-based grading system: 1) Reassessment for full credit after further learning; and 2) Basing a student's semester grade on their number of proficient and number of high-performance learning outcomes. Treatment teachers were provided with a project-developed software tool, *PARLO Tracker*, and trained in how to use it. Teachers could enter into Tracker the learning outcomes they had developed and enter evidence collected, including assessments and reassessments, about each student's progress on each learning outcome. Students and parents could access Tracker online to track student progress. (Tracker was developed with user input and became available at the beginning of the second year of the project, which is the first year that quantitative data used in the current study was collected. Use was optional, and when surveyed during Year 3 of the project, 16 of 25 PARLO teachers (64%) reported using Tracker regularly.) Additional topics covered in the professional development included: creating scoring rubrics for learning outcomes and sharing them with students; helping students track their proficiency; teaching students what to do when they are not yet proficient; requiring proof of new learning before reassessment, including use of tools like error logs and remediation plans; scoring proficiency in learning outcomes based on the best evidence to date instead of traditional methods like averaging scores; and converting learning outcome scores into semester grades. PARLO teachers were responsible for developing their own learning outcomes based on state standards, developing their own methods for informing students about the PARLO system and grading scheme, establishing their own classroom procedures, and deciding on their own means of collecting assessment and/or reassessment evidence. However, through PLCs they were able to share ideas and support each other in developing routines, assessments, and so on.

Quantitative Measures

Implementation Measures

In the spring of 2012, we administered a 7-item measure to Treatment and Control teachers, investigating whether they implemented key aspects of the PARLO system. In

the spring of 2013 we administered an additional survey to PARLO teachers only, investigating in more detail to what extent they had implemented the PARLO system. See the online appendix section *Implementation: Additional Details* for a detailed description of the measures.

Algebra and Geometry Content Exams

The state where we conducted our research had recently adopted content standards for both geometry and algebra based on the Common Core State Standards for Mathematics (National Governors Association Center for Best Practices & Council of Chief State School Officers, 2010), but had not yet implemented an end-of-course test of the material. Consequently, we created a test to mimic as closely as possible a state-designed test of its published content standards. To do so, we administered (with permission) an adapted version of the Virginia Standards of Learning multiple choice Algebra and Geometry Tests. For each course (algebra or geometry), we selected items from the Virginia test that addressed our state's standards. Then, for any standards not covered, we added items adapted from textbook sources or from released items from older assessments in our participating state. In each participating classroom, we administered the appropriate content exam as both a pretest during the first five days of the school year (Algebra: $M = 34.2\%$, $SD = 11.0\%$, Cronbach's $\alpha = 0.72$; Geometry: $M = 36.0\%$, $SD = 11.3\%$, $\alpha = 0.75$) and, with items in a different order, during the last month of the school year (Algebra: $M = 46.9\%$, $SD = 17.5\%$, $\alpha = 0.90$; Geometry: $M = 49.7\%$, $SD = 14.9\%$, $\alpha = 0.94$). Following suggestions by May et al. (2009) for combining differing tests into a single measure, we used a linear transformation to standardize the pretests within each subject area (algebra or geometry) to have a mean of 0 and a standard deviation of 1, and then combined the two sets of scores to create a single measure. We did the same for post-test scores. See the online appendix for the complete tests used and technical details about computing student scores.

Motivation and Motivational Antecedents Measures

We assessed concepts from the expectancy-value theory of motivation using the Attitudes Toward Mathematics Inventory (ATMI; Tapia & Marsh, 2004). The four ATMI subscales are described below. In each participating classroom, we collected baseline scores on the four ATMI subscales during the first five days of the school year, and we collected outcome data on the same subscales during the last month of the school year. Students responded to all of the subscales using a 5-point scale (1 = "strongly disagree"; 5 = "strongly agree"). See the online appendix for the full set of items from the ATMI and technical details about computing student scores.

Perceived Intrinsic Value of Mathematics. The extent to which students saw mathematics as having intrinsic value was assessed using the ATMI's 10-item "enjoyment" subscale. Items included "I have usually enjoyed studying mathematics in school" (baseline: $M = 3.27$; $SD = 0.88$, $\alpha = 0.89$; outcome measure: $M = 3.08$; $SD = 0.86$, $\alpha = 0.89$).

Perceived Utility Value of Mathematics. The extent to which students saw mathematics as having utility value was assessed using the ATMI's 10-item "value" subscale. Items included "I can think of many ways that I use math outside of school" (baseline: $M = 3.73$; $SD = 0.67$, $\alpha = 0.86$; outcome measure: $M = 3.54$; $SD = 0.74$, $\alpha = 0.87$).

Expectancies for Success in Mathematics. The strength of students' expectancies for success in mathematics was assessed using the ATMI's 15-item "self-confidence" subscale. Items included "I believe I am good at solving math problems" (baseline: $M = 3.45$; $SD = 0.77$, $\alpha = 0.79$; outcome measure: $M = 3.31$; $SD = 0.83$, $\alpha = 0.81$).

Long-Term Motivation to Engage with Mathematics. The ATMI's 5-item "motivation" subscale assessed the degree to which students enjoyed the challenge of mathematics and planned to pursue the subject over the long term. Items included "I plan to take as much mathematics as I can during my education" (baseline: $M = 3.21$; $SD = 0.79$, $\alpha = 0.79$; outcome measure: $M = 3.10$; $SD = 0.85$, $\alpha = 0.81$).

Student-Level Covariates

We coded each student's race and gender based on information provided by the school or, when school-provided data were not available, based on a self-report survey administered simultaneously with the ATMI. A preliminary analysis using Akaike's Information Criterion (AIC; Akaike, 2011) indicated that the most informative way to operationalize "race" was to create a dichotomous variable, with "1" indicating the student was identified as White (51%) or Asian (4%), and 0 indicating that the student was identified as Black (30%), Hispanic (8%) or Multiracial/other (7%).

School-Level Covariates

We obtained from state databases the schoolwide percent low-income students enrolled in the academic year 2010–11, the first year of PARLO implementation for the study. We also aggregated student-level covariates (race, gender, and content pretest and baseline motivation-related scores) to the school level as potential additional covariates to control for in our analysis.

Qualitative Data

Qualitative data were collected throughout the three operational years of the project. The qualitative data sources used for this paper included 84 PARLO teacher interviews and an open-ended survey of 678 students from a nonrandom sample of 6 of the 14 PARLO treatment schools.

Teacher Interviews

During the first year of participation for Cohort 1 and Cohort 2, PARLO teachers were randomly selected to be interviewed; however, if a teacher was the only participant at their school, they were automatically selected to be interviewed. During the 2011–12

school year, the 13 Cohort 1 teachers who were interviewed were the same teachers who were interviewed the previous year. During the 2012–13 school year, all 25 remaining participating teachers were interviewed. Interviews were conducted in the fall and spring of each year, and some teachers were interviewed twice during the same year. Over the three years of the study, 11 teachers were interviewed once, 8 were interviewed twice, 7 were interviewed three times, and 9 were interviewed four times. Overall, 35 teachers at 15 schools participated in 84 interviews.

Interview data for the current article was drawn primarily from questions that addressed how PARLO affected students—for example, “Does student engagement ‘look’ different? By that, I mean are students interacting with math, with you, with each other differently under PARLO?” See the online appendix for detailed interview questions. Interviews were semi-structured, permitting teachers to expound upon questions that warranted further elaboration. With few exceptions, interviews were conducted following a lesson observation, either during teachers’ preparatory or lunch periods, or, rarely, after school. Interviews lasted about 25 to 35 minutes and were audio recorded with the teacher’s permission. Every interviewee agreed to be recorded. Each interview was transcribed, then entered into NVivo software for analysis.

Student Survey

During the 2011–12 school year, six of the treatment schools agreed to have their students complete an open-ended online 8-question survey (e.g., “How would you describe PARLO to next year’s ninth grade students?”) A total of 678 students completed the survey. See the online appendix for the full list of questions asked on the student survey.

Data Analysis

Quantitative Analyses

We conducted all statistical analyses using the *lmer* command from the *lmerTest* package of the R programming language (Kuznetsova et al., 2017). We used Restricted Maximum Likelihood (REML) to test all research questions, but as recommended by Zuur et al. (2009), we used maximum likelihood in preliminary analyses that employed AIC (Akaike, 2011) to choose variance components and covariates, or that used log likelihood to evaluate heteroscedasticity. We calculated degrees of freedom for statistical tests using Satterthwaite’s approximation (Satterthwaite, 1946).

Our data set identified each student by course (geometry vs. algebra), teacher, school, and study year. Consequently, the data were structured as a five-level hierarchical linear model: students within course within study-year within teachers within schools. We used AIC to choose the most appropriate variance structure for our analysis, using students’ mathematics content post-test as our dependent variable. This analysis indicated that after entering variance components for school and for course-within-year-within-teacher, adding additional variance components for year-within-teacher or teacher-within-school reduced the efficiency of the model. We therefore used the following three-level variance structure in our analyses: “student,” within “course-within-study-year-within-teacher,” within “school.”

Our first quantitative analysis investigated the main effects of PARLO using mathematics content post-test scores as the dependent variable. We used AIC to select covariates for our model, choosing the following set: school level percent low-income students in 2011; a dichotomous indicator identifying students enrolled in geometry; a dichotomous variable indicating female sex; a dichotomous variable indicating White or Asian race (vs. Black, Hispanic, or multiracial/other); students' baseline expectancy score; linear and quadratic terms for the pretest; and interactions between the geometry indicator and the linear and quadratic pretest variables. We also included indicators to ensure that students were compared to other students in the same study-year and in the same blocking group used for random assignments. The reference group was Cohort 1 public school students. Other groups were: Catholic school students; Charter School students; Cohort 2 Public school students tested in 2010–11; and Cohort 2 public school students tested in 2011–12. All non-dichotomous student level variables were centered around the student mean, and percent low-income was centered around the mean of the school means.

Our second quantitative analysis investigated the main effects of PARLO on the four motivation-related subscales obtained by the ATMI. We used the same covariates that we used the first analysis, with the following exceptions: we used scores for all four motivation and antecedent subscales, whereas previously only expectancy scores were used. We also removed the quadratic term for the pretest and its interaction with the geometry indicator. See the online appendix for further details about how we selected covariates to use in each of our models.

Our two main-effects analyses used two approaches to account for missing data. First, we used Full Information Maximum Likelihood (FIML; Allison, 2012), an approach that, like multiple imputation, minimizes bias due to missingness. However, because it is difficult to analyze interactions between treatment and student-level covariates using FIML, and because these interactions were important in addressing *for whom* the PARLO system might be effective, we also ran the analyses using listwise deletion. The main effects estimated by the two analyses were very similar; consequently, we report detailed results of the listwise deletion analysis here. See [Tables A10](#) and [A12](#) of the online appendix for details of the FIML analysis.

Our third quantitative analysis, unlike the first two, was not pre-planned. We conducted the analysis to see whether we could corroborate an unanticipated finding from our qualitative data: According to teacher interviews, student motivation moderated PARLO's impact, such that more highly motivated students benefited more from the PARLO program (see Results for details). We therefore used interaction terms to investigate how our four motivation-related quantitative variables moderated PARLO's impact on mathematics content post-test scores.

Note that for each of the four measures of motivation, we had two observations per student: a score from the first week of the school year, and a score from the last month. Because motivation can change throughout the school year, our preferred approach to measuring motivation across the year would be to average the baseline and end-of-year scores together. We note that end-of-year motivation scores might be endogenous, which can make interpretations of some models confusing or bias results under certain circumstances. However, there is preliminary evidence (based on Monte Carlo results)

that if x is an endogenous variable and w is an exogenous variable, the coefficient of the interaction term xw (i.e. the moderator effect) will be consistently estimated as long as w is binary and x is homoscedastic conditional on the other variables in the model (Bun & Harrison, 2018). PARLO treatment was a binary variable, and all four motivation variables were homoscedastic, so we could safely analyze the interactions using the averaged baseline-and-end-of-year scores. These moderation analyses used the same covariates we had used when testing the main effects of PARLO on mathematics post-test scores, except that instead of using the baseline confidence score as a covariate, we included the main effect for each moderator being tested. As a sensitivity check, we also re-ran the moderation analysis using baseline scores instead of year-average scores for all four motivation-related moderators. See the online appendix for details.

As recommended by the What Works Clearinghouse Procedures Handbook (USDOE, 2020), whenever we used multiple statistical tests to address a single issue we used the Benjamini-Hochberg (BH) correction to account for multiple comparisons. The BH procedure controls the false discovery rate for multiple comparisons, ensuring that the expected proportion of falsely identified statistically significant results equals the intended alpha, in this case 5%.

Qualitative Analyses

At the conclusion of teachers' participation in the study (June 2013), two qualitative researchers began coding three years' worth of interview transcripts. To enhance reproducibility, the researchers established a coding scheme using a procedure similar to the three-stage process outlined by Campbell et al. (2013). In the first stage, each researcher read the same three interview transcripts, jotting notes in the margins to identify possible broad themes. Once themes were identified, the researchers met to compare notes and look for agreement or disagreement. If there was a lack of consensus, they discussed why they had coded the text under a particular theme, until agreement was reached that either only one of their themes was accurate, or that both of their themes were accurate. In the latter instances, the text was double coded. At the end of this process, 17 broad themes were identified.

In the second stage, the remaining transcripts were divided, and each researcher coded half of the interview transcripts. In the third stage, each researcher assumed responsibility for the analyses of a defined set of themes. The researchers then prepared memos by theme and study year, with each memo summarizing major evidence for the existence of that theme in the interviews conducted that year. Initial analysis of the student survey followed a similar process, with an additional analysis that identified key words or phrases that appeared frequently in students' responses to survey prompts.

To produce the findings described in this article, a third researcher read all three years' worth of memos to both identify any claims that addressed how or why the PARLO system affected students' motivational antecedents, motivation, or academic performance, and to summarize evidence supporting those claims. Finally, for claims that were supported by numerous teachers' interviews, the research team made a final pass through all 84 transcribed interviews to quantify how many teachers supported that particular claim.

Results

PARLO Implementation

All PARLO teachers interviewed indicated they had implemented the new system, allowing reassessment for full credit and basing final grades on the number of proficient learning outcomes and the number of high-performance learning outcomes. In practice, most teachers rated individual assignments using a “traffic light” system: Red (not yet proficient), Yellow (approaching proficient), Green (proficient), or Blue (high performance; available only for high cognitive demand tasks). In interviews, teachers confirmed that learning outcome grades were developed by teacher judgments based on cumulative evidence, with more recent evidence weighted most heavily. Most PARLO teachers also required proof of learning before reassessment. In an online questionnaire administered to the 25 PARLO teachers participating during the 2012–13 school year, 88% reported using resubmission/correction of work, error logs, remediation plans, and/or flashback days at least “fairly often.”

When surveyed in the spring of 2012, PARLO teachers were significantly more likely than Control teachers to agree that they planned to utilize “a ‘high performance’, ‘proficient’, or ‘not yet proficient’ assessment system” (76% vs. 22% agreement) and use “the traffic light system to monitor student progress” (50% vs. 15% agreement). Both these contrasts were statistically significant after using the BH adjustment to control for a false discovery rate of 5%.

Two additional differences between groups were significant after controlling for a false discovery rate of 10%, but not 5%. First, PARLO teachers were somewhat more likely to implement “holding students accountable for learning outcomes” (94% vs. 74%) and somewhat more likely to implement “collecting evidence of student learning” (88% vs. 67%). Both Treatment and Control groups indicated they were highly likely to implement “using formative assessment strategies in your classrooms”: 91% for PARLO teachers and 93% for Control teachers. See the online appendix for additional details about the implementation measures and their results.

Overall, the above data indicate that the contrast between the PARLO teachers and Control teachers appears to have been as intended. A large majority of both groups implemented the PARLO characteristics depicted in the gray boxes of [Figure 1](#), although there is evidence that the PARLO group may have implemented learning outcomes more extensively. The PARLO teachers were much more likely than Control teachers to base a student’s grades on the number of proficient and number of advanced learning outcomes. While we do not have statistics on the percentage of Control teachers who gave students the opportunity to reassess for full credit after additional learning, all PARLO teachers interviewed indicated that they had adopted this practice.

PARLO’s Impact on Quantitative Measures of Mathematics Performance

Main Effects: Did PARLO Increase Mathematics Performance?

[Table 1](#) summarizes the PARLO treatment’s main effects on student mathematics learning. Recall that students’ mathematics content test scores were standardized, so that within each subject area (algebra or geometry), the post-tests had a mean of 0 and a

Table 1. PARLO treatment main effect on mathematics achievement.

	Estimate	Standard Error	Degrees of Freedom	p-value
<i>School Level Fixed Effects</i>				
Intercept	−0.66	0.12	20	<.0001
Catholic Girls' School	0.30	0.22	21	.18
Charter School	−0.05	0.29	33	.86
Cohort 2 Public School, 2011-2012 school year	0.23	0.15	14	.14
Cohort 2 Public School, 2012-2013 school year	0.34	0.15	14	.035
Proportion Disadvantaged	−0.53	0.23	12	.039
PARLO Treatment	0.33	0.12	14	.014
<i>Course Level Fixed Effects</i>				
Geometry Student	0.20	0.16	50	.23
<i>Student Level Fixed Effects</i>				
Baseline Expectancy of Success	0.20	0.02	2,687	<.0001
White or Asian	0.15	0.04	2,654	<.0001
Female	0.10	0.03	2,672	.0005
Pretest	0.26	0.02	2,710	<.0001
Pretest-squared	0.04	0.01	2,673	.0028
Geometry X Pretest	0.20	0.05	2,720	.0002
Geometry X Pretest-squared	−0.02	0.03	2,700	0.54
<i>Random Effects</i>				
	<i>N</i>	<i>Variance</i>		
School	29	0.03		
Course x Teacher x Year	85	0.14		
Residual	2,736	0.53		

Notes: This analysis utilized data from 2,736 students, 65 teachers, and 29 schools. "Mathematics Achievement" was defined as the Algebra post-test z-score for algebra students and the Geometry post-test z-score for geometry students. All non-dichotomous variables are grand mean centered. The reference group for the blocking variables was Cohort 1 Public School students who participated during the 2011-12 school year. The intercept represents the post-test score for students at average values for pretest, expectancy of success, and school-level proportion disadvantaged and in the reference group for all dichotomous variables (i.e. in the Cohort 1 Public School randomization block; in a control school, taking algebra, not White/Asian, and not Female).

standard deviation of 1. Consequently, the results reported in Table 1 can be interpreted as effect sizes in standard deviation units. After controlling for random-assignment grouping, students' pretest scores, schoolwide percent low-income, geometry vs. algebra class, race, sex, and baseline expectancy for success, students in the PARLO system scored 0.33 *SD* higher than students in the control group on the project-administered end-of-course tests ($p = 0.014$).

Using post-test standard deviation as the denominator and combining Treatment with Control groups, Geometry students gained approximately 0.91 standard deviations and Algebra students gained approximately 0.73 standard deviations from our pretest to our post-test. Thus, an impact of 0.33 *SD* might be thought of as roughly equivalent to about 36% (for Geometry students) to 45% (for Algebra students) of a year's learning.

Interaction Effects: For Whom Did Academic Performance Increase?

We investigated interactions between the PARLO program and each of the student level covariates in our model to determine whether baseline student characteristics moderated the effects of the PARLO system on mathematics content learning. There were no statistically significant interactions. In other words, students benefited from the PARLO system regardless of their course enrollment (geometry or algebra), their content pretest scores, their baseline expectancies for success in mathematics, their race, or their gender. See the online appendix for complete details regarding these analyses.

PARLO's Impact on Quantitative Measures of Students' Mathematics Motivation

Main Effects: Did PARLO Impact Motivation?

Table 2 summarizes the PARLO treatment's main effects on students' self-reported intrinsic value, utility value, expectancies for success, and long-term motivation to engage with mathematics. (See Table A11 of the online appendix for more detailed information including covariate coefficients and variance components.) As can be seen in the table, the PARLO system had substantively small and statistically non-significant effects on post-test measures of the four constructs, with effects ranging from -0.08 to -0.02 on a 5-point scale.

Interaction Effects: Were There Any Subgroups for Whom Average Motivation Changed?

While there were no main effects of treatment on any of the motivation related measures, it seemed possible that one or more subgroups of students might have been impacted. We investigated this possibility by testing whether any of our four measures were impacted by interactions between treatment condition and each of our student level covariates: gender, race, geometry vs. algebra class, and baseline expectancies, intrinsic value, utility value, long-term motivation, and mathematics content test scores. After using the Benjamini-Hochberg (BH) procedure to control for multiple comparisons, none of the interactions were significant. (See the online appendix for additional details.) Thus, we found no quantitative evidence that the PARLO system affected the intrinsic value, utility value, expectancies for success, or long-term motivation of any student subgroup.

Mechanisms Leading to PARLO's Positive Impact on Mathematics Performance

Opportunities to Learn

We had hypothesized that PARLO would improve student achievement partly by improving opportunities to learn. Teacher interviews provided detail about how increased opportunities may have operated.

Table 2. Effects of PARLO treatment on students' motivational antecedents and long-term motivation in math class.

Dependent Variable (Subscale)	PARLO effect Estimate	Standard Error	df	95% conf. interval	p-value	Effect size in SD units
Intrinsic Value	-0.06	0.06	15	$(-0.18, +0.06)$.31	-0.07
Utility Value	-0.02	0.05	14	$(-0.12, +0.07)$.60	-0.03
Expectancy	-0.06	0.04	51	$(-0.15, +0.02)$.14	-0.08
Long-term Motivation	-0.08	0.06	19	$(-0.21, +0.05)$.23	-0.09

Notes: This analysis utilized data from 2,698 students, 65 teachers, and 29 schools. Dependent variables are measured on a 1-5 Likert scale. df = degrees of freedom. Standard deviations for the dependent variables were as follows: Intrinsic Value: 0.86; Utility value: 0.74; Expectancies for success: 0.83; Long-term motivation: 0.85. Reported results controlled for the following covariates: Assignment block for randomization; School-level proportion disadvantaged; course assignment (geometry or algebra), student-level gender, race, pretest, and baseline score on all four subscales of the ATMI.

Reassessment Opportunities. The central feature of PARLO is encouraging students to re-study and reassess to achieve proficiency or high performance on each learning outcome, and a recurring theme—emerging among 89% of teachers interviewed (31/35)—was that teachers felt this central feature of PARLO led directly to better student learning of mathematics content. Specifically, teachers reported that reassessment opportunities were useful to students of all ability levels. Regarding struggling students, teachers made comments like the following:

...some of the students that have typically struggled in math are able to see how their hard work pays off... I think it's something that they have not had before, because if they struggled on a test, even if they came in and got help, that test grade still remained, and they weren't able to demonstrate their understanding. (Teacher #31)

Teachers noted that high achievers also benefited, specifically from opportunities to reassess for high performance. For example:

The kids that really, really want to learn it will do anything that they can to achieve that high performance. In the traditional class, they've got that one shot to master it, and they can't revisit it. So, they can't ever get to that achievement. They can't push themselves to reach that next level. It's either one or none. [With] the PARLO, they can go back and push themselves to get it. (Teacher # 32)

Clearer Focus on Learning. Sixty percent of teachers interviewed (21/35) also noted that providing feedback and reassessment opportunities for each learning outcome helped students focus their learning efforts fruitfully, asking better questions and showing better understanding of what they needed to work on—for example:

They're coming to me with better questions. They're not just coming to me and saying, "What can I do?" They come in knowing where they need to focus, knowing what they need to work on. (Teacher #15)

Peer Interactions. Although the interviews did not explicitly ask about student peer support, 49% of teachers interviewed (17/35) brought up peer collaboration under PARLO as one contributor to better student learning—for example:

They tend to actually be getting involved more with other students, because... they also recognize quickly that in teaching others, they're getting a better understanding themselves because now they can present it and defend. (Teacher #6)

Revisiting Topics. Thirty-one percent of teachers interviewed (11/35) noted the pedagogical value of allowing students to revisit topics over time—for example:

One big thing is that PARLO allows me time to come back and go over things that we already covered... And by going over the learning outcomes... I see students that have gaps in their learning. Like, one student, she was truant in February and March, but she is a very high-level student, so now I see her learning all of that stuff and connecting the dots between our most recent learning outcomes and stuff that we learned three months ago. And also, just kids filling in the gaps in their own learning as we go through it. "Oh, I remember this" or "I forget how we do this" and then re-remembering it... I noticed that they're getting it very easily now the second time around. (Teacher #23)

Increased Engagement

While our quantitative analysis did not detect any program impact on motivation, teacher interviews indicated that student engagement, which is often associated with motivation, did increase. Eighty percent of teachers (28/35) interviewed said that under PARLO, their students were more likely to participate in discussions about content with both the teacher and their fellow students. The same percent of teachers said that students were more likely to ask questions when concepts were not well understood.

When asked to describe “successes” under the PARLO program, 63% of teachers (22/35) reported increased student persistence in the face of initial difficulty, increased ownership of the learning process, and increased responsibility—for example:

Students are taking ownership in what they are learning... They are learning how to organize and keep up with it, each marking period they are getting better at it. (Teacher #15)

Earlier this year, the K-8 math coordinator visited my classroom, to follow up on some of her students from last year. After class she said to me, “Was that Rachel [pseudonym] In the front row, with her head up?” I said, “Yes.” And the coordinator said she could not believe it was the same girl because last year, Rachel, when she came to class, sat in the back of the classroom with her head on her desk. This year, Rachel told me during one-on-one conference that she wants high performance on every Outcome. And she’s always asking for more and more tests. She’s also begging me to teach her next year as she wants to stay in a PARLO classroom. (Teacher #4)

Examining Potential Reasons behind PARLO’s Non-Significant Effects on Quantitative Measures of Student Motivation

Posner’s (2011) pilot study of the PARLO system found positive effects on intrinsic value, utility value, expectancy for success, and motivation in college statistics. In contrast, the current study found no effects of PARLO on the motivation-related constructs we measured. Our qualitative analysis investigated why this might have occurred.

Decreased Engagement for Some Students

Despite teachers’ reports that under PARLO overall engagement increased, a minority of teachers expressed concern that the PARLO system might interfere with some students’ engagement. They identified two potential problems, which we labeled “contentment” and “procrastination.”

Contentment. Twenty-nine percent of teachers interviewed (10/35), indicated that PARLO might discourage some students from doing their best. Specifically, they noted that some students seemed content with proficient performance and reluctant to try high cognitive demand problems in order to achieve high performance—for example:

What I stopped doing was marking the problems with [an] asterisk if it were blue [high performance], because students that are not that motivated would just choose not to do those... (Now) they don’t know if it’s green [proficient] or blue [high performance] until we get the assessments back. (Teacher #28)

Procrastination. Twenty-six percent of teachers interviewed (9/35) indicated that the PARLO system might encourage some students to put off studying under the assumption that if they did poorly at first, they could always reassess later—for example:

As far as kids that don't take advantage, we have so many. They're the kids that aren't motivated... John [pseudonym] has gotten to the point where he pretty much hands in tests with nothing on them. And I think that he does subconsciously think, "Oh, I can reassess," but he doesn't follow through. (Teacher #36)

Potential Net Zero Impact on Expectancies for Success

Re-assessment opportunities under PARLO appeared to increase students' expectancies for success, but this positive effect was perhaps balanced by a negative effect on expectancies caused by student uncertainties about their final grades.

Positive Effects of Re-Assessment Opportunities. In 2011–12, 678 PARLO students completed an open-ended survey about their experiences with the new assessment system. When asked "What do you like about PARLO?" the most common word or phrase students mentioned was "retake" ($N=90$), followed by "make up" ($N=52$) and "another chance" ($N=48$). Students made comments like, "It gives you the chance to retake tests, quizzes, or homework. You can work at your own pace without having to rush to learn something new," and "I think PARLO is very good for math students struggling."

Negative Effects of Uncertainties about Final Grades. When asked, "What do you not like about PARLO?" 41% of students (277/678) indicated that they felt they never knew their exact grade until report card time or the end of the year, which they found confusing. When asked, "How would you describe PARLO to next year's ninth grade students?," the single most commonly mentioned word was "confusing," mentioned by 73 of the 678 students interviewed. Even students who found the system helpful complained about uncertainty in grading, making comments like, "PARLO is very helpful. You can always bring your grade up. However, one downer is that you never know what your grade is until report cards arrive."

No Detected Impact on Intrinsic Value and Utility Value

As noted in our methods section, the first step in our qualitative analysis was inductive, identifying major themes that emerged from the data. Changes in students' intrinsic value for mathematics or changes in students' beliefs about mathematics utility value did not emerge as themes. Unlike student expectancies for success, which seemed to have a balance of positive and negative impacts, we found no evidence that the PARLO system had either a positive or a negative impact on students' valuing of mathematics.

Unmeasured Motivational Constructs

A majority of the teachers interviewed (24 out of 35, or 69%) indicated that students' experiences with PARLO enhanced their motivation. The discrepancy between this finding and our quantitative results (i.e., no program effect on motivational antecedents or

long-term motivation) may have occurred because, by “motivation,” teachers were referring to constructs we did not measure quantitatively.

Mastery Goals. Of the 35 teachers interviewed, 24 (69%) indicated that students’ experiences with PARLO promoted the adoption of mastery goals. For example:

...last year all they were concerned about were points and grades. This year they are talking about math more than my students last year. They know their content. (Teacher #26)

Growth Mindsets. Of the 35 teachers interviewed, 20 (57%) also reported that, under PARLO, students were more likely to adopt a growth mindset—for example:

They are finally, really understanding that knowledge is gained and built over time, you don’t just know something, you don’t just get something and that’s it. You have to work at it, and if you want to keep it you have to continue to work at it. (Teacher #24)

I think their work ethic is better. They realize they need to keep on top of stuff, keep working. If you fail one thing, don’t just say, ‘Oh, it’s over.’ It’s not over! You can still learn stuff. (Teacher #1)

Autonomy. A smaller but still substantial number of teachers (34%; 12/35) described a different mechanism by which PARLO increased motivation: that under PARLO, some students felt that they had autonomy and could control their own mathematics destinies (Ryan & Deci, 2020)—for example:

They have the ownership of their grade now. It’s no longer “What can he give me?” So, the ball is in her court now. It’s no longer, “You failed me, or you gave me this grade.” They talk about getting their grade up to where they want it to be. (Teacher #16)

It should be noted, however, that some teachers highlighted one characteristic of the PARLO assessment system that caused some students to perceive *decreased* autonomy: Under PARLO, a student’s grades are based solely on evidence of content understanding. Compliance measures like homework completion and attendance no longer count toward grades. Some students expressed concern that they could control their compliance, but not necessarily ensure their learning.

Relatedness. Finally, PARLO may also have helped meet students’ needs for relatedness (Ryan & Deci, 2020). As noted above, teachers reported that increased peer support was one of the reasons that mathematics performance improved under PARLO. This increased peer support may have been effective because it improved students’ sense of relatedness. For example, a teacher described classmates’ response to a PARLO student coming up to the chalkboard to demonstrate proficiency on a learning outcome:

They were still supportive ... And they’re like, “Come on, Jenny (pseudonym), you can do it! You just have a few more steps to go! You’re almost there!” Never did they yell out, “Hey stupid, wrong step ...” They’re using the language [of PARLO]. (Teacher # 4)

The Emergence of Student Motivation as a Potential Moderator of PARLO's Impact

We turn now to an unexpected finding that emerged from our qualitative analysis. We had theorized that the PARLO system might improve students' math motivation. In contrast, teacher interviews suggested that the converse might also be true: higher student motivation might improve the effectiveness of the PARLO system.

Why did some students respond to the PARLO system with active and enthusiastic engagement, whereas other students responded with apathy and work avoidance? As illustrated in the quote from Teacher #36 in our discussion (above) of procrastination, teachers attributed the different student behaviors to differences in student motivation. In fact, 71% of teachers (25/35) noted that "motivation," "caring," "willingness" or a similar concept was important in determining the effectiveness of PARLO, making comments like the following:

Motivation is huge. And if the kids don't have the motivation, then who cares if they can take a test? They don't care. (Teacher #2)

The thing I really like about PARLO is the kids that are willing to work and seem to actually care, are the ones who really seem to benefit from it. (Teacher #9)

Quantitative Verification of Motivation as a Moderator of PARLO Impact

As noted in our Methods section, when we designed this study we did not plan to conduct a quantitative analysis to see whether student motivation might moderate the PARLO treatment's effect on student mathematics learning. However, after the qualitative data indicated that this might be the case, we conducted a quantitative analysis to see whether we could corroborate the unanticipated qualitative finding. If student motivation did moderate treatment effects on student achievement, we would expect to see a positive interaction between PARLO treatment and each of our three measured psychological antecedents of motivation, as well as between PARLO and our measure of long-term motivation. Table 3 summarizes the results of the analysis.

As can be seen in the table, all four interactions were positive and three of the four were statistically significant. To interpret the table substantively, the positive impact of PARLO on the end of the year mathematics content score increased by 0.09, 0.11, 0.13, and 0.07 *SD* for each 1-point increase on the 5-point measure of intrinsic value, utility value, expectancies, and long-term motivation, respectively. Thus, PARLO's impact on individual students' mathematics achievement was positively affected by their levels of motivation in mathematics.

We used the BH adjustment for multiple comparisons to account for the fact that Table 3 employs four statistical tests. The BH procedure confirmed the statistical significance of intrinsic value, utility value, and expectancies for success as moderators. Note that the analysis reported in Table 3 averaged the baseline and end-of-year scores to measure each of the four motivation-related constructs. As a sensitivity check, we also ran the moderation analysis that instead used only baseline scores for all four motivation-related moderators. The results pointed in the same direction as those reported in

Table 3. Interactions with treatment condition: do students' perception of intrinsic value, utility value, expectancy of success, and long-term motivation moderate PARLO's impact on mathematics achievement?

Interaction	n	Effect size Estimate	Standard Error	Degrees of Freedom	95% conf. interval	p-value
PARLO \times Average Intrinsic Value	2,529	0.09	0.04	2,482	(+0.01, +0.17)	0.026
PARLO \times Average Utility Value	2,532	0.11	0.05	2,483	(+0.01, +0.21)	0.028
PARLO \times Average Expectancy	2,529	0.13	0.04	2,481	(+0.05, +0.22)	0.002
PARLO \times Average Longterm Motivation	2,533	0.07	0.04	2,485	(-0.02, +0.15)	0.13

Notes: This analysis utilized data from 65 teachers, and 29 schools. "Mathematics Achievement" was defined as the Algebra Post-test z-score for algebra students and the Geometry post-test z-score for geometry students. n = number of students with data available for the model testing each specific moderator. Because the analysis utilized average of baseline and end-of-year scores, students who were missing end-of-year scores were not included in the analysis. Reported results controlled for the following covariates: Assignment block for randomization; School-level proportion disadvantaged; course assignment (geometry or algebra), student-level gender, race, pretest, pretest-squared, main effects of the motivation subscale being studied and main effects of the PARLO treatment.

Table 3 but were not statistically significant. See the online appendix for details about the BH procedure and the sensitivity check.

Discussion

The present work investigated the impact of the PARLO standards-based grading system on ninth graders' learning of algebra and geometry, hypothesizing that the reengineered grading system would positively impact student learning both by directly providing students additional opportunities to learn, and in a mediated manner by enhancing student motivation. Partly supporting these predictions, our quantitative analysis indicated that PARLO had a positive impact on ninth graders' learning of mathematics content, with an estimated effect size of 0.33 SD . In normally distributed data, an effect size of 0.33 SD would move a student from the 50th percentile on a test up to the 63rd percentile. Comparing the PARLO impact to typical student growth from pre- to post-test in our data set, we estimate 0.33 SD to be the equivalent to about 36% to 45% of a year's learning. The program was effective regardless of students' race, gender, or prior achievement.

As discussed, we had also hypothesized that one mechanism through which PARLO might improve achievement would be the mediating effect of motivation: the PARLO program would improve motivation, which in turn would improve student engagement and learning. Our quantitative analysis did not support this hypothesis. We found no significant effect of PARLO on any of our four motivation-related measures.

Our qualitative analysis provided insights into the reasons behind our quantitative findings. As predicted, teachers reported a number of ways that increased opportunities to learn under PARLO improved student achievement. These opportunities included the direct effects of encouraging students to re-study and reassess, helping students to focus on the content they needed to learn, increased support from students' peers, and providing opportunities to revisit topics over time.

Qualitative data also provided insight into why students did not report increased expectancies for success under PARLO. There appear to have been two countervailing tendencies, potentially leading to a net zero impact. As predicted, students' appreciation

of the opportunities to reassess appears to have strengthened their expectancies. Balancing this, however, one aspect of our PARLO implementation may have had the unintended consequence of decreasing students' expectancies. The assessment system at YWLCS, which inspired PARLO, did not assign students letter grades, instead using long report cards that described proficiency by learning outcome. In our professional development sessions, we showed these report cards to teachers and suggested, "When all outcomes have been rated, students' overall proficiency can be converted to a letter grade if necessary." As a result, many PARLO teachers provided students with little information about their overall grade until all outcomes had been rated at the end of the semester or year. This appears to have created anxiety and confusion among some students about their grades, which may have reduced their expectancies for success. We anticipate that future implementations of the PARLO system will correct this problem, instead encouraging teachers to help students understand where they stand relative to a final mathematics grade and what they need to do if they want to improve it.

Our quantitative analysis did not indicate any program effects on students' valuing of math class, and the qualitative data seldom mentioned the topic. While we had hypothesized that changing the assessment system to focus less on quick learning and ranking and more on mastery of course content would positively impact students' value of mathematics, this expectation may have been unreasonable. Making mathematics more enjoyable (i.e., increasing intrinsic value) or increasing the perceived utility value of mathematics may be most directly impacted by designing engaging instructional activities or by explicitly connecting curriculum to potential applications, neither of which was a focus of PARLO.

Teachers interviewed did, however, report that the PARLO system positively impacted student engagement and motivation in ways that our quantitative analysis did not address. Motivational antecedents that may have been positively affected include growth mindsets and mastery goals, autonomy, and relatedness.

An additional unanticipated and noteworthy finding emerged from our qualitative analysis. Teacher interviews indicated that student motivation moderated PARLO's impact—that is, that more motivated students benefited more strongly from PARLO. Once this finding emerged from our qualitative analysis, our quantitative data provided some corroboration. Students who scored higher on year-average scores for intrinsic value, utility value, or expectancies for success in mathematics benefited more from the PARLO system than did students who scored lower on those same measures.

To unify these distinct but interconnected findings, we return to [Figure 1](#), which depicts our current conceptual model regarding how the PARLO system may influence students' academic experiences and learning. Solid lines in the figure indicate relationships that have been supported either by prior research or that were supported both by our present quantitative *and* qualitative analyses. Dashed lines indicate relationships that we hypothesize based on patterns that emerged in our qualitative data, but have not received robust quantitative support in the present or past research, and the lack of a connection indicates relationships that did not emerge in the present or prior research.

We note that this current conceptual model includes two refinements to our initial hypotheses. First, we now theorize that the PARLO system impacts motivation indirectly through improving some, but not all, of the motivational antecedents we have

considered. Specifically, the current study provides evidence that the PARLO standards based grading system impacts expectancies, autonomy, relatedness, growth mindset, and mastery goals—but not perceived intrinsic value or utility value. Second, motivation is believed to moderate the relationship between PARLO-provided opportunities to learn and actual learning. Thus, supplementing the PARLO program with additional support for motivational antecedents would in theory magnify the impact of the program on learning.

Study Limitations

We designed our study using an expectancy-value framework (Wigfield & Eccles, 2000), and our quantitative measures reflect that design. Findings that emerged from our qualitative analysis of teacher interviews caused us to expand the framework to include constructs from growth mindset theory (Dweck, 2007), self-determination theory (Ryan & Deci, 2020), and achievement goal theory (Senko, 2016). However, because these connections emerged through inductive qualitative analyses conducted following data collection, quantitative data on growth mindsets, autonomy, relatedness, and mastery goals were not collected. This may have contributed to the discrepancy between teacher interviews, which reported increased motivation under PARLO, and quantitative analyses, which did not detect a program impact on motivation.

The study would have been stronger had we interviewed the control teachers. Doing so might have provided important insight into the ways in which treatment and control students' experiences in the classroom differed.

The PARLO system works by changing the grading system to provide students with more opportunities to learn. To be effective, students must understand the system. We did not provide direct assistance to help teachers explain the PARLO grading system to their students, and our data indicate that some students did not understand it. In a free response survey, 73 out of 678 students described PARLO as “confusing.” Future implementations might be able to increase PARLO's effectiveness by providing materials and techniques teachers can use to explain the new system to their students.

Informal aspects of our intervention may have influenced outcomes in important but untracked ways. For example, some but not all teachers encouraged peer-to-peer support by grading a student as high performance on a learning outcome if they helped another student attain proficiency. “Flashback days” were an innovation developed by some participating teachers and shared with all PARLO teachers through PLCs, but not implemented by all of them. Teachers also developed techniques to reduce workload requirements under PARLO. Some of these innovations may be key ingredients that greatly enhance program success, but we were not able to determine which innovations were the most important. Future research should seek to track how such elements—in addition to factors like the level of transparency with which proficiency ratings are translated into final grades, the way the program is introduced to students, the amount and types of supports provided to parents and guardians, and the number of courses or grade levels in a school that implement the PARLO program—might affect program results.

Finally, we note that treatment and control teachers had all volunteered, if assigned to treatment, to change a central pillar of their classroom instruction: the grading system. Implementing the PARLO system required considerable effort from the teachers involved. Results may not generalize to teachers or schools who are less willing to make such a change.

Conclusion: The Future

Working with teachers who were already implementing formative assessment and mathematics instruction designed around clearly state learning outcomes, the PARLO standards-based grading program reinforced those teaching practices and added two new elements: reassessment for full credit after further learning; and basing a student's final grade on the number of learning outcomes the student learned at a proficient level and the number of learning outcomes the student learned at a high performance level. This treatment improved the amount of mathematics students learned by 0.33 standard deviations, which is roughly 36% to 45% of the effect an entire year's learning had for our sample.

Moving forward, we anticipate that future work will not replicate the current project precisely, but rather build on lessons learned. For example, we recommend that future implementations correct elements of the program that may have unintentionally reduced program effectiveness. Specifically, we recommend that procedures for computing final grades be made transparent and explained clearly to students; and that the program provide materials and procedures to assist teachers in introducing the PARLO program to students. If possible, future implementations should provide optional-use assessment, reassessment, and re-teaching materials to reduce the amount of work required of PARLO teachers. Future work might also consider implementing standards-based grading in more courses or more grade levels (vs. ninth grade math only), so that the PARLO grading system does not seem so unusual (and thus potentially confusing) compared to grading in other courses.

Our data also provides evidence that the PARLO standards-based grading system, while effective regardless of race, gender, or prior knowledge, is especially effective for students who are more highly motivated. This finding makes sense: the PARLO system provides students additional opportunities to learn, and more motivated students are more likely to take advantage of those opportunities. If one views motivation as a fixed attribute of students, then this finding could be seen as troubling. A standards-based grading system would have positive effects, but it would be especially helpful for the "haves" who are already motivated and engaged.

Fortunately, however, research has consistently shown that motivation is malleable: the mindsets that are psychological antecedents to motivation, including beliefs in intrinsic and utility value of subject matter, expectancies for success, growth mindsets, autonomy, and relatedness are greatly impacted by teacher behaviors and classroom learning conditions (Binning & Browman, 2020; Farrington et al., 2012). Furthermore, educational psychologists have developed interventions that are designed to target the specific underlying psychological processes that can promote student motivation and engagement. Importantly, a growing body of research has found that such interventions

can have strong positive impacts on student learning and well-being, but only in contexts that are fertile ground for supporting the target psychological factors (Walton & Yeager, 2020, Yeager & Walton, 2011). To create such “fertile ground” researchers have suggested that it may be especially important to reform grading systems in order to make them less focused on summative assessment and ranking (Farrington et al., 2012). This raises an exciting avenue for future work: implementing the PARLO standards-based grading system, while simultaneously supporting PARLO with interventions that foster student motivation. Such a combined intervention could theoretically magnify the impact of PARLO by supporting student motivation, while simultaneously magnifying the impact of the motivation interventions by creating a grading environment in which they can have a larger impact.

Finally, we note that there is no theoretical reason to expect the PARLO standards-based grading system to be less successful in other grade levels or in other disciplines than it was in ninth grade mathematics. We encourage other researchers to build on the potential identified in our work, finding effective ways to implement standards-based grading. We hope that the results will have large benefits for cultivating student learning.

Open Research Statements

Study and Analysis Plan Registration

There is no study and analysis plan registration associated with this manuscript.

Data, Code, and Materials Transparency

The data, code, and materials underlying the results reported in this manuscript are available on the Open Science Framework: <https://osf.io/htkfs>.

Design and Analysis Reporting Guidelines

This manuscript is accompanied by a completed JREE Randomized Trial Checklist.

Transparency Declaration

The lead author (the manuscript’s guarantor) affirms that the manuscript is an honest, accurate, and transparent account of the study being reported; that no important aspects of the study have been omitted; and that any discrepancies from the study as planned (and, if relevant, registered) have been explained.

Replication Statement

This manuscript reports an original study.

Acknowledgments

The authors thank Helen Kramer, Holly Bozeman, John Baker, Teresa Harrison, Maurice Bun and two anonymous reviewers for their feedback. This research was supported by the National Science Foundation under grant number DRL-0918474. Research involving human subjects was reviewed and approved by the Institutional Review Board (IRB) at Research for Better Schools. Any opinions, findings and conclusions or recommendations expressed in this material are those of the author(s) and do not necessarily reflect the views of the National Science Foundation.

Disclosure Statement

No potential conflict of interest was reported by the author(s).

Open Scholarship



This article has earned the [Center for Open Science](#) badges for Open Data and Open Materials through Open Practices Disclosure. The data and materials are openly accessible at [10.35542/osp.io/pzc3f](https://doi.org/10.35542/osp.io/pzc3f).

ORCID

Steven L. Kramer  <http://orcid.org/0000-0003-2900-6484>

Alexander S. Browman  <http://orcid.org/0000-0002-2957-3262>

References

- Akaike, H. (2011). Akaike's information criterion. In M. Lovric (Ed.), *International encyclopedia of statistical science* (pp. 25–25). Springer. https://doi.org/10.1007/978-3-642-04898-2_110
- Allison, P. D. (2012). Handling missing data by maximum likelihood. In *SAS Global Forum*, Paper 312-212. Retrieved September 8, 2021, from <https://statisticalhorizons.com/wp-content/uploads/MissingDataByML.pdf>
- Binning, K. R., & Browman, A. S. (2020). Theoretical, ethical, and policy considerations for conducting social-psychological interventions to close educational achievement gaps. *Social Issues and Policy Review*, 14(1), 182–216. <https://doi.org/10.1111/sipr.12066>
- Black, P., & Wiliam, D. (2010). Inside the black box: Raising standards through classroom assessment. *Phi Delta Kappan*, 92(1), 81–90. <https://doi.org/10.1177/003172171009200119>
- Bloom, B. S. (1968). Learning for mastery. Instruction and curriculum. Regional education laboratory for the Carolinas and Virginia, topical papers and reprints, number 1. *Evaluation Comment*, 1(2), 1–12. <https://eric.ed.gov/?id=ED053419>
- Bun, M. J. G., & Harrison, T. D. (2018). OLS and IV estimation of regression models including endogenous interaction terms. *Econometric Reviews*, 38(7), 814–827. <https://doi.org/10.1080/07474938.2018.1427486>
- Campbell, J. L., Quincy, C., Osserman, J., & Pedersen, O. K. (2013). Coding in-depth semistructured interviews: Problems of unitization and intercoder reliability and agreement. *Sociological Methods & Research*, 42(3), 294–320. <https://doi.org/10.1177/0049124113500475>
- Clymer, J. B., & Wiliam, D. (2007). Improving the way we grade science. *Educational Leadership*, 64(4), 36–42. <https://eric.ed.gov/?id=EJ766296>
- Covington, M. V., & Omelich, C. L. (1984). Task-oriented versus competitive learning structures: Motivational and performance consequences. *Journal of Educational Psychology*, 76(6), 1038–1050. <https://doi.org/10.1037/0022-0663.76.6.1038>

- Dweck, C. S. (2007). *Mindsets: The new psychology of success*. Ballantine.
- Farrington, C. A., Roderick, M., Allensworth, E., Nagaoka, J., Keyes, T. S., Johnson, D. W., & Beechum, N. O. (2012). *Teaching adolescents to become learners. The role of noncognitive factors in shaping school performance: A critical literature review*. University of Chicago Consortium on Chicago School Research. Retrieved March 18, 2013, from <http://ccsr.uchicago.edu/publications/teaching-adolescents-become-learners-role-noncognitive-factors-shaping-school>.
- Farrington, C. A., & Small, M. (2008). *A new model of student assessment for the 21st century*. American Youth Policy Forum. <https://www.researchgate.net/profile/Camille-Farrington/publication/253449779>
- Haley, J. M. (2015). *To curve or not to curve? The effect of college science grading policies on implicit theories of intelligence, perceived classroom goal structures and self-efficacy* [Unpublished doctoral dissertation]. Boston College. <http://hdl.handle.net/2345/bc-ir:104165>
- Jury, M., Smeding, A., & Darnon, C. (2015). First-generation students' underperformance at university: The impact of the function of selection. *Educational Psychology*, 6, 710. <https://doi.org/10.3389/fpsyg.2015.00710>
- Kingston, N., & Nash, B. (2011). Formative assessment: A meta-analysis and a call for research. *Educational Measurement: Issues and Practice*, 30(4), 28–37. <https://doi.org/10.1111/j.1745-3992.2011.00220.x>
- Kuznetsova, A., Brockhoff, P., & Christensen, R. (2017). lmerTest package: Tests in linear mixed effects models. *Journal of Statistical Software*, 82, 1–26. <https://doi.org/10.18637/jss.v082.i13>
- Leahy, S., Lyon, C., Thompson, M., & William, D. (2005). Classroom assessment: Minute by minute, day by day. *Educational Leadership*, 63(3), 19–24. <https://eric.ed.gov/?id=EJ745452>
- Marzano, R. J. (2010). *Formative assessment & standards-based grading*. Marzano Research Laboratory.
- May, H., Perez-Johnson, I., Haimson, J., Sattar, S., & Gleason, P. (2009). *Using state tests in education experiments: A discussion of the issues* (NCEE 2009-013) National Center for Education Evaluation and Regional Assistance, Institute of Education Sciences, U.S. Department of Education.
- Mills, V. L., & Silver, E. A. (2018). Putting it all together and moving forward: Concluding thoughts. In E. A. Silver & V. L. Mills (Eds.), *A fresh look at assessment in mathematics teaching*. National Council of Teachers of Mathematics.
- Murphy, R., Roschelle, J., Feng, M., & Mason, C. A. (2020). Investigating efficacy, moderators and mediators for an online mathematics homework intervention. *Journal of Research on Educational Effectiveness*, 13(2), 235–270. <https://doi.org/10.1080/19345747.2019.1710885>
- National Governors Association Center for Best Practices, Council of Chief State School Officers. (2010). *Common core state standards*. National Governors Association Center for Best Practices, Council of Chief State School Officers.
- Posner, M. A. (2011). The impact of a proficiency-based assessment and reassessment of learning outcomes system on student achievement and attitudes. *Statistics Education Research Journal*, 10(1), 3–14. <https://eric.ed.gov/?id=EJ925274> <https://doi.org/10.52041/serj.v10i1.352>
- Ryan, R. M., & Deci, E. L. (2020). Intrinsic and extrinsic motivation from a self-determination theory perspective: Definitions, theory, practices, and future directions. *Contemporary Educational Psychology*, 61, 101860. <https://doi.org/10.1016/j.cedpsych.2020.101860>
- Sanchez, C. E., Atkinson, K. M., Koenka, A. C., Moshontz, H., & Cooper, H. (2017). Self-grading and peer-grading for formative and summative assessments in 3rd through 12th grade classrooms: A meta-analysis. *Journal of Educational Psychology*, 109(8), 1049–1066. <https://doi.org/10.1037/edu0000190>
- Satterthwaite, F. E. (1946). An approximate distribution of estimates of variance components. *Biometrics Bulletin*, 2(6), 110–114. <https://doi.org/10.2307/3002019>
- Scarlett, M. H. (2018). “Why did I get a C?”: Communicating student performance using standards-based grading. *InSight: A Journal of Scholarly Teaching*, 13, 59–75. <http://files.eric.ed.gov/fulltext/EJ1184948.pdf> <https://doi.org/10.46504/14201804sc>
- Schneider, J., & Hutt, E. (2014). Making the grade: A history of the A–F marking scheme. *Journal of Curriculum Studies*, 46(2), 201–224. <https://doi.org/10.1080/00220272.2013.790480>

- Senko, C. (2016). Achievement goal theory: A story of early promises, eventual discords, and future possibilities. In K. R. Wentzel & D. B. Miele (Eds.), *Handbook of motivation at school* (2nd ed., pp. 75–95). Routledge. <https://www.researchgate.net/profile/Corwin-Senko/publication/312591307>
- Smeding, A., Darnon, C., Souchal, C., Toczek-Capelle, M.-C., & Butera, F. (2013). Reducing the socio-economic status achievement gap at university by promoting mastery-oriented assessment. *PloS One*, 8(8), e71678. <https://doi.org/10.1371/journal.pone.0071678>
- Supovitz, J. A., Ebby, C. B., Remillard, J., & Nathenson, R. A. (2018). Experimental impacts of the ongoing assessment project on teachers and students. CPRE Research Report # RR 2018-1. Consortium for Policy Research in Education. <https://eric.ed.gov/?id=ED593465>
- Tapia, M., & Marsh, G. E. (2004). An instrument to measure mathematics attitudes. *Academic Exchange Quarterly*, 8, 16–22. Downloaded 8/30/2021 from <http://www.rapidintellect.com/AEQweb/cho25344l.htm>
- U.S. Department of Education. (2020). *What works clearinghouse procedures handbook, v4.1*. Institute of Education Sciences, National Center for Education Evaluation and Regional Assistance, What Works Clearinghouse. Downloaded 8/30/2021 from <https://ies.ed.gov/ncee/wwc/handbooks>
- U.S. Department of Education. (2022). *What works clearinghouse procedures and standards handbook, v5.0*. Institute of Education Sciences, National Center for Education Evaluation and Regional Assistance, What Works Clearinghouse. Downloaded 12/15/2022 from <https://ies.ed.gov/ncee/wwc/handbooks>
- Walton, G. M., & Yeager, D. S. (2020). Seed and soil: Psychological affordances in contexts help to explain where wise interventions succeed or fail. *Current Directions in Psychological Science*, 29(3), 219–226. <https://doi.org/10.1177/0963721420904453>
- Webb, N. L. (2002). Depth-of-knowledge levels for four content areas. *Language Arts*, 28, 1–9. http://facstaff.wcer.wisc.edu/normw/All_content_areas_DOK_levels_2032802.pdf
- Wigfield, A., & Eccles, J. S. (2000). Expectancy-value theory of achievement and motivation. *Contemporary Educational Psychology*, 25(1), 68–81. <https://doi.org/10.1006/ceps.1999.1015>
- Wiliam, D. (2011). *Embedded formative assessment*. Solution Tree.
- Yeager, D. (2017). Social-emotional learning programs for adolescents. *The Future of Children*, 27(1), 73–94. <https://doi.org/10.1353/foc.2017.0004>
- Yeager, D. S., Lee, H. Y., & Dahl, R. E. (2017). Competence and motivation during adolescence. In A. J. Elliot, C. S. Dweck, & D. S. Yeager (Eds.), *Handbook of competence and motivation: Theory and application* (pp. 431–448). The Guilford Press. <https://www.researchgate.net/publication/337856583>
- Yeager, D. S., & Walton, G. M. (2011). Social-psychological interventions in education: They're not magic. *Review of Educational Research*, 81(2), 267–301. <https://doi.org/10.3102/0034654311405999>
- Zuur, Alain, Ieno, E. N., Walker, N., Saveliev, A., & Smith, G. M. (2009). *Mixed effects models and extensions in ecology with R*. Springer. https://doi.org/10.1007/978-0-387-87458-6_1

Online Appendix

Details of Random Assignment and Attrition.

Our study design called for us to recruit 42 schools, 21 Treatment and 21 Control, in order to have an 80% chance of detecting an effect size of 0.25. Recruitment was to begin after the program received funding. Because the funding arrived later than anticipated only 20 schools were recruited to participate by the time the study began in spring of 2010: 14 public schools, 3 charter schools, and 3 Catholic all-girl schools. In June of 2010, all ninth-grade algebra and geometry teachers at all Cohort 1 schools who had agreed to participate attended three days of professional development (PD) focusing on formative assessment techniques and on developing clear Learning Outcomes. The 20 participating schools were then randomly assigned either to the Treatment or the Control condition. Randomization for the original 20 schools was done in three blocks by type of school: Charter vs. Catholic vs. Public, with each school having a 50% chance to be assigned to each condition. Among the Charter schools two were randomly assigned to the Treatment condition and one to the Control condition. Similarly, among the Catholic schools two were randomly assigned to Treatment condition and one to the Control condition. Among the Public schools, 6 were assigned to the Treatment condition and 8 to the Control condition. Note: the uneven assignment of public schools (6 to Treatment vs 8 to Control) occurred because before random assignment we agreed with one participating district to assign that district's two schools randomly to the SAME treatment. Consequently, we assigned 13 "units" randomly to Treatment or Control: 6 to Treatment and 7 to Control. However 1 of the 7 Control units consisted of two schools. See "Two Middle Schools" below for more information.

After randomization, one Public Control school, one Public Treatment school, and one Charter Treatment school dropped out of the project. Additionally, a Charter Treatment school

delayed implementation for a year.

The next year, we recruited 15 additional schools: 14 high schools from one large urban school district and 1 additional public high school. We provided Cohort 2 teachers two days of PD, and then randomly assigned the 15 schools to conditions: 8 to Treatment and 7 to Control. After randomization for Cohort 2, one Control and one Treatment school dropped out of the program.

Table A1 summarizes the number of schools assigned and the number that remained in the project, by treatment condition and cohort.

Table A1

Number of Schools by Assignment Condition and Cohort

<i>Cohort</i>	<i>Assignment Group</i>	<i># Assigned</i>	<i># With Data Available for Analysis</i>	<i>% Attrition</i>
Cohort 1	Treatment	10	7	30.0%
	Control	10	9	10.0%
	Treatment + Control	20	16	20.0%
Cohort 2	Treatment	8	7	12.5%
	Control	7	6	14.3%
	Treatment + Control	15	13	13.3%
Cohorts 1 & 2	Treatment	18	14	22.2%
	Control	17	15	11.8%
	Treatment + Control	35	29	17.1%

Idiosyncrasies of Random Assignment and Sensitivity Analyses to Address Them

Charter School that Delayed Participation. As described above, of the three Charter schools we recruited, one dropped out of the program; a second, randomly assigned to Control status, provided unusable data in 2010-11 and usable data in 2011-12; and the third, randomly assigned to Treatment status, provided usable data in 2011-12 and 2012-13. Our reported analyses did not use the 2012-13 data from the Charter Treatment school, because there was no comparison data available from a Charter Control school for that year. We conducted sensitivity

analyses, reported below, to see whether either including both years of data from the Charter Treatment school or dropping both Charter schools altogether from the analysis made any substantive change in our results.

Two Middle Schools. One Cohort 1 district that agreed to participate in our study contained one high school and two middle schools. That district taught Algebra in eighth grade, not ninth grade, and requested that the two middle schools participate in lieu of the high school and be randomly assigned together either to Treatment or Control condition. After the 3 days of professional development, the two schools were randomly assigned to Control, and they are counted among the 10 Cohort 1 Control schools reported in Table A1. We conducted sensitivity analyses, reported below, to see whether dropping those two schools from the analysis made any substantive change in our results.

Sensitivity Analyses. We conducted five sensitivity analyses to see whether our results were robust to alternative decisions about which school data should be counted in our data set. Our sensitivity check analyzed the following five alternative data sets to the analytic data set used in the main article:

Alternative 1: Use both 2011-12 and 2012-13 data from Charter Schools. (Student n= 2,775 for the analysis of PARLO effects on student achievement, and between 2,568 and 2,572 for models using an ATMI subscale as the dependent variable.)

Alternative 2: Drop both Charter schools from the analysis. (Student n= 2,669 for the analysis of PARLO effects on student achievement, and between 2,463 and 2,467 for models using an ATMI subscale as the dependent variable..)

Alternative 3: Drop the two middle schools from the analysis. (Student $n = 2,675$ for the analysis of PARLO effects on student achievement, and between 2,470 and 2,474 for models using an ATMI subscale as the dependent variable.)

Alternative 4: Use both 2011-12 and 2012-13 data from Charter Schools AND drop the two middle schools from the analysis. (Student $n = 2,714$ for the analysis of PARLO effects on student achievement, and between 2,509 and 2,513 for models using an ATMI subscale as the dependent variable.)

Alternative 5: Drop both Charter schools from the analysis AND drop the two middle schools from the analysis. (Student $n = 2,608$ for the analysis of PARLO effects on student achievement, and between 2,404 and 2,408 for models using an ATMI subscale as the dependent variable.)

Table A2 reports the results of the five sensitivity analyses. Overall, the results confirm those reported in the main article. In all cases, the PARLO effect on academic outcomes is minimally changed from the 0.33 standard deviations estimated in the main article, and remains statistically significant. In all cases, the PARLO effect on all motivation-related outcomes remains small, negative, and statistically insignificant, as in the main article. In all cases, the interaction terms “Average Expectancy of Success” and “Utility Value” remain significant moderators of the PARLO effects, as they are in the analysis reported in the main article—although under Alternative 3 the interaction term for Utility Value does not meet the .025 significance level required by the Benjamini-Hochberg criterion for ensuring the false discovery rate is kept below 0.05. (See the relevant section below, as well as Table A20, for more details about how we performed Benjamini-Hochberg adjustments). The interaction term for “Intrinsic Value” remains a significant moderator in some models, but not when the two middle schools are dropped from the analysis, i.e., alternatives 3 through 5.

Table A2

Sensitivity of Results to Sample Selection Choices

	Alternative 1 Extra Charter School Year		Alternative 2 No Charter Schools		Alternative 3 No Middle Schools		Alternative 4 Extra Charter School Year, No Middle Schools		Alternative 5 No Charter Schools, No Middle Schools	
Main Effects of PARLO treatment										
Dependent variable	Effect Size		Effect Size		Effect size		Effect Size		Effect Size	
	Estimate	p-value	Estimate	p-value	Estimate	p-value	Estimate	p-value	Estimate	p-value
Content Tests	0.33	.014	0.34	.018	0.36	.013	0.36	.012	0.36	.016
Intrinsic Value	-0.06	.31	-0.06	.30	-0.06	.29	-0.07	.27	-0.07	.26
Utility Value	-0.02	.73	-0.02	.64	-0.03	.50	-0.03	.56	-0.04	.48
Expectancy of Success	-0.06	.17	-0.05	.23	-0.07	.11	-0.07	.12	-0.06	.16
Long-term Motivation	-0.05	.45	-0.05	.47	-0.07	.32	-0.07	.36	-0.07	.37
Moderation Effects on Student Content Knowledge										
Moderator	Effect Size		Effect Size		Effect Size		Effect Size		Effect Size	
	Estimate	p-value	Estimate	p-value	Estimate	p-value	Estimate	p-value	Estimate	p-value
Intrinsic Value x PARLO	0.09	.031	0.09	.027	0.07	0.077	0.07	.089	0.08	.063
Utility Value x PARLO	0.11	.024	0.13	.014	0.11	.028	0.12	.024	0.13	.015
Expectancy x PARLO	0.13	.0024	0.14	.0012	0.12	.0066	0.12	.0082	0.13	.004
Long-term Motivation x PARLO	0.07	.10	0.08	.084	0.06	.18	0.07	.15	0.07	.13

Student-level Attrition

During the first week of each school year, Treatment and Control teachers administered a content pretest and a survey collecting demographic information and student responses to motivation-related Likert scales. Only students at non-attributing schools who supplied data on both the pretest and the survey were included in our *reference sample* and used in the analysis. Table A3 summarizes by study-year, cohort, and treatment condition: the number of students in the reference sample; the number of students in the reference sample who provided outcome data and thus were included in the analysis; and the attrition rate.

Table A3
Number of Students by Assignment Condition, Cohort, and Year of Participation

<i>Cohort</i>	<i>Assignment Group</i>	<i># In Reference Sample</i>	<i># With Data Available for Analysis</i>	<i>% Attrition</i>
Cohort 1, 2011-12 school year	Treatment	797	647	18.8%
	Control	483	421	12.8%
	Treatment + Control	1,280	1,068	16.6%
Cohort 2, 2011-12 school year	Treatment	537	465	13.4%
	Control	448	358	20.1%
	Treatment + Control	985	823	16.4%
Cohort 2, 2012-13 school year	Treatment	602	537	10.8%
	Control	406	308	24.1%
	Treatment + Control	1,008	845	16.2%
All students	Treatment	1,936	1,649	14.8%
	Control	1,337	1,087	18.7%
	Treatment + Control	3,273	2,736	16.4%

Teacher-level Counts and Attrition

Cohort 1 provided two years of data, i.e., data collected in 2010-11 and data collected in 2011-12, but only the 2011-12 data was deemed valid data for the study. We did not maintain records on teacher participation in 2010-11. In 2011-12 Cohort 1 included 16 Control teachers

and 20 PARLO teachers. Since we did not maintain records about teacher participation in 2010-11, we did not control for first vs. second year of program participation.

It is common for multi-year studies, especially those conducted in urban districts, to experience significant year-to-year teacher churn, as teachers leave their positions, are transferred to new schools within a district, or simply switch grade levels or courses they teach while staying at the same school (Baker & Baruch, 2015; Boruch, Merlino, & Porter, 2012). Our data from Cohort 2 confirmed this trend. Cohort 2 provided usable data for 2011-12 and 2012-13. Table A4 reports by treatment condition the number of Cohort 2 teachers who participated in 2011-12, the number of those who remained to participate in 2012-13, and the number of new teachers who joined the project (replacing Algebra or Geometry teachers in participating schools) during 2012-13.

Table A4

Cohort 2: Teacher Assignment, Retention, and Joining

Assignment Group	<i># Teachers who participated 2011-12</i>	<i># Teachers who stayed 2012-13</i>	<i># Teachers who joined 2012-13</i>	<i>Total # of teachers who participated 2012-13</i>
Treatment	13	7	5	12
Control	9	6	2	8
Treatment + Control	22	13	7	20

Student Demographics

Table A5 reports by treatment condition the demographic characteristics of our analytic sample of 2,736 students attending 14 Treatment and 15 Control schools.

Table A5

Demographic Characteristics of the Analytic Sample Students, Treatment vs Control

	<i>PARLO (n = 1,649)</i>		<i>Control (n = 1,087)</i>	
<i>Demographics</i>	<i>Percent</i>	<i>n</i>	<i>Percent</i>	<i>n</i>
Female	56%	1,649	54%	1,087
Asian	3%	1,649	7%	1,087
Black	29%	1,649	31%	1,087
Hispanic	6%	1,649	11%	1,087
White	57%	1,649	41%	1,087
Multi/Other	6%	1,649	10%	1,087
<i>Baseline Scores</i>	<i>Mean (SD)</i>	<i>n</i>	<i>Mean (SD)</i>	<i>n</i>
Algebra PreTest	+0.07 (1.04)	1,574	-0.14 (0.91)	886
Geometry PreTest	-0.13 (0.80)	75	+0.05 (1.06)	201
Combined Algebra & Geometry Tests	+0.06 (1.03)	1,649	-0.10 (0.95)	1,087
Intrinsic Value	3.28 (0.88)	1,649	3.28 (0.87)	1,087
Utility Value	3.74 (0.65)	1,649	3.71 (0.71)	1,087
Expectancy	3.46 (0.76)	1,649	3.47 (0.78)	1,087
Long-term Motivation	3.20 (0.77)	1,649	3.25 (0.79)	1,087

Note: In preparing Table A5, we standardized the algebra and geometry pretests to have a mean of 0 and SD of 1 for the full data set. Mindset subscales are measured from 1 to 5. *n* = number of students providing data for this statistic. SD= “standard deviation” and is included in parentheses where appropriate.

Teacher Demographics

Table A6 is based on data from all teachers in the analytic sample who completed a demographic survey.

Table A6

Demographic Characteristics of the Analytic Sample Teachers, Treatment vs. Control		
<i>Demographic Characteristic</i>	<i>PARLO Teachers (n=33)</i>	<i>Control Teachers (n=25)</i>
Median Years Taught	7	6
Percent Certified	88%	100%
Percent Math Major	82%	78%
Percent Female	67%	69%
Percent Asian	0%	9%
Percent Black	10%	3%
Percent White	90%	86%
Percent Multi-racial/Other	0%	3%

Timeline

Figure A1 displays the timeline for program implementation. As the timeline shows, each Treatment and Control school participated in the project for two years: Cohort 1 during 2010-11 and 2011-12 and Cohort 2 during 2012-12 and 2012-13. At each participating school, the sample of students included ninth graders enrolled in Algebra or Geometry. We collected two years of data at each participating school, but as explained in the main article the data collected during 2010-11 could not be used. Consequently, we analyzed one year of data from one set of ninth graders who attended Cohort 1 schools and two years of data (one each from two separate sets of ninth graders) from Cohort 2 schools. Students in our data sets were anonymized, with pretests, surveys, and posttests matched using a barcode assigned to each student. We had available only the barcode to use as student identifiers. This means that if some students in Cohort 2 repeated ninth grade and provided data both in 2011-12 and 2012-13 we could not know it. Of necessity, we treated the two years of data from Cohort 2 as coming from separate sets of ninth graders. This applied to both Treatment and Control schools.

Figure A1 PARLO RCT Timeline
2010-2011

School Cohort 1	<p>Early summer Randomization: 20 participating schools identified.</p> <p>PD: 3 days for all Grade 9 Algebra and Geometry teachers in participating schools.</p>	<p>July Randomization: 10 Treatment and 10 Control schools randomly assigned.</p>	<p>August Attrition: 2 Treatment schools and 1 Control school drop from the program. 1 additional Treatment school delays implementation until next year, yielding 7 Treatment and 9 Control schools participating in Cohort 1, Year 1.</p> <p>PD: 3 days for all Treatment teachers.</p>	<p>September Data Collection: COHORT 1, 1ST GROUP OF 9TH GRADERS: achievement pretest and baseline attitudes survey administered to volunteer students. (These data were later deemed invalid due to the volunteer nature of the sample.)</p>	<p>September-June PD: Approximately monthly PLC sessions at Treatment schools</p> <p>Implementation: COHORT 1, 1ST GROUP OF 9TH GRADERS: Full PARLO implementation at 7 Treatment Schools.</p> <p>Data Collection: Teacher Interviews at all Treatment schools.</p>	<p>June Data Collection: COHORT 1, 1ST GROUP OF 9TH GRADERS: End-of-year achievement test and attitudes survey administered to volunteer students. (These data were later deemed invalid due to the volunteer nature of the sample.)</p>
------------------------	--	--	--	--	--	---

2011-2012

School Cohort 1	Early summer No activity.	July No activity.	August Attrition: 1 Treatment school drops, but an additional treatment school that had delayed active participation rejoins Cohort 1. Thus, the Cohort 1 group remains: 7 Treatment and 9 Control schools. PD: 2 days for all Treatment teachers.	September Data Collection: COHORT 1, 2 ND GROUP OF 9 TH GRADERS: achievement pretests and baseline attitudes survey administered to all students.	September-June PD: Approximately monthly PLC sessions at Treatment schools Implementation: COHORT 1, 2 ND GROUP OF 9 TH GRADERS: Full PARLO implementation at 7 Treatment Schools. Data Collection: Teacher Interviews at all Treatment schools.	June Data Collection: COHORT 1, 2 ND GROUP OF 9 TH GRADERS: End-of-year achievement test and attitudes survey administered to all students.
School Cohort 2	Early summer Randomization: 15 participating schools identified. PD: 2 days for all Grade 9 Algebra and Geometry teachers in participating schools.	July Randomization: 8 Treatment and 7 Control schools randomly assigned.	August Attrition: 2 Treatment and 1 Control school drop from the program, yielding 7 Treatment and 6 Control schools. PD: 4 days for all Treatment teachers.	September Data Collection: COHORT 2, 1 ST GROUP OF 9 TH GRADERS: achievement pretests and baseline attitudes survey administered to all students.	September-June PD: Approximately monthly PLC sessions at Treatment schools Implementation: COHORT 2, 1 ST GROUP OF 9 TH GRADERS: Full PARLO implementation at 7 Treatment Schools. Data Collection: Teacher Interviews at all Treatment schools.	June Data Collection: COHORT 2, 1 ST GROUP OF 9 TH GRADERS: End-of-year achievement test and attitudes survey administered to all students.
Both Cohorts						June Data Collection: Qualitative survey of 678 students.

Timeline: 2012-2013

School Cohort 2	Early summer No activity.	July No activity.	August Attrition: none PD: 2 days for all Treatment teachers.	September Data Collection: COHORT 2, 2 ND GROUP OF 9 TH GRADERS: achievement pretests and baseline attitudes survey administered to all students.	September-June PD: Approximately monthly PLC sessions at Treatment schools Implementation: COHORT 2, 2 ND GROUP OF 9 TH GRADERS: Full PARLO implementation at 7 Treatment Schools. Data Collection: Teacher Interviews at all Treatment schools.	June Data Collection: COHORT 2, 2 ND GROUP OF 9 TH GRADERS: End-of-year achievement test and attitudes survey administered to all students.
------------------------	-------------------------------------	-----------------------------	--	---	--	---

Implementation: Additional Details

Spring 2012 Survey

Treatment and Control teachers completed seven PARLO-related questions in a survey administered in the spring of 2012. At that time, teachers in Cohort 2 were just completing the first of two years of project participation. Meanwhile, teachers in Cohort 1 were just finishing their second year of participation in the project.

The PARLO-related questions were set up as a checklist responding to the following prompt: “Looking ahead to the next school year (2012 -13), what elements of PARLO, if any, are you likely to implement? Please check all that apply.” Note: While data were not collected from Cohort 1 classrooms during the Year 3 of the project (2012-2013), Treatment teachers in Cohort 1 were permitted to continue participating in PARLO PLCs during Year 3, and many of them opted to do so.

Table A7 describes the percent of teachers who checked off each prompt on the list, broken down by Treatment vs. Control.

Table A7

Percent of Teachers Saying They Were Likely to Implement Each Practice During the Coming School Year, by Treatment Condition

	PARLO	Control	Chi-square	p-value
Holding students accountable for learning outcomes	94%	74%	1.723	.043*
Utilizing a “high performance”, “proficient”, or “not yet proficient” assessment system	76%	22%	2.431	<.001**
Allowing students to resubmit	32%	41%	-0.363	.499

assignments for full credit				
Reassessment	85%	78%	0.505	.451
Collecting evidence of student learning	88%	67%	1.322	.049*
Using the traffic light system to monitor student learning	56%	15%	1.986	.002**
Using formative assessment strategies in your classrooms	91%	93%	-0.190	.841

Note: N = 34 Treatment and 27 Control teachers. “Chi-square” = Chi-square value for a logistic regression analysis with 1 degree of freedom. “p-value” = significance of the Chi-square test before Benjamini-Hochberg adjustment.

** Statistically significant difference between groups after using the Benjamini-Hochberg (BH) adjustment to control for a false discovery rate of 5%.

* Statistically significant difference between groups after using the BH adjustment to control for a false discovery rate of 10%.

As noted in the published article, PARLO teachers were significantly more likely than Control teachers to say they planned to use a “high performance”, “proficient” or “not yet proficient” assessment system and that they planned to use a traffic light system to monitor student learning. The fact that 15% of Control teachers (4 out of 27) reported planning to use a traffic light system could indicate some cross-contamination. However, it is likely that instead of referring to a grading system they were referring to a similarly named comprehension-monitoring formative assessment system that was popular in local school districts, in which students displayed red, yellow, or green drinking cups to indicate how well they were understanding a teacher’s presentation.

Two items on the questionnaire were aimed at assessing whether teachers were implementing the PARLO practice of encouraging students to reassess for full credit after further work. As shown in the table, there was little difference between groups on either item,

“Reassessment” or “Allowing students to resubmit assignments for full credit”. In retrospect, we think both items were poorly phrased. Regarding the “reassessment”: Math teachers commonly reassess topics (e.g., on a quiz and then on a final) but unlike PARLO teachers, they average early scores with later scores. Another common practice is to allow students who failed a test to retake that test and average the two grades. Such practices emphasize test-performance (e.g., “Maybe the student had a bad day”), while PARLO emphasizes partnering with each student to ensure they all reach high levels of understanding of each learning outcome. Regarding “allowing students to resubmit assignments for full credit”: PARLO teachers often required students to turn in corrected assignments before reassessment, e.g. using error logs—but the resubmitted assignment did not receive credit by itself. Instead, students had to prove separately that they had mastered the material.

It is also worth noting the two items that approached statistical significance (significant at $p=.10$ but not $p=.05$ after BH adjustment). At the time of our study, it was becoming increasingly common practice for teachers to teach by learning outcome, and we chose to lean into this version of “business as usual” by encouraging both Treatment and Control teachers to do so. However, the PARLO system emphasized and reinforced use of learning outcomes, and the data seem to indicate that PARLO teachers were especially likely to use them. Similarly, while under business-as-usual it was common for all teachers to regularly collect evidence of student learning, participating in PARLO may have made teachers especially likely to collect evidence.

Finally, we note that although using formative assessment is a necessary part of PARLO, as expected both Treatment and Control teachers were familiar with formative assessment and both groups at least tried to use formative assessment techniques.

Spring 2013 Survey

In the spring of 2013 we administered a survey to the participating PARLO teachers only, asking about various classroom practices. Many of the questions overlapped the questions that had been asked of both PARLO and Control groups in spring of 2012. Because the 2012 questionnaire had been administered to both groups, we focus on the 2012 data, reported above, in our analysis of those questions. However, four questions addressed an aspect of the PARLO system not previously measured: whether teachers required evidence of further study before permitting reassessment.

The questions responded to the following prompt: “During the 2012-13 school, year how often, if at all, did you use the following strategies?” Teachers could respond: Never; Rarely (i.e., once or twice); A few times (i.e., 3-5 times); Fairly Often (i.e. 6 times to about once per month); Frequently (i.e., a couple of times a month); Almost every week. Table A8 describes the percent of teachers who selected each response option for each of the three questions.

Table A8

Percent of Year 3 PARLO Teachers Selecting Each Response Option

	Never	Rarely	A Few Times	Fairly Often	Frequently	Almost Every Week
Allowed students to resubmit assignments for full credit	8%	8%	8%	12%	20%	44%
Had students complete error logs for work that was "Not Yet Proficient"	16%	20%	24%	16%	16%	8%
Had students complete remediation plans in order to become proficient	16%	12%	16%	12%	28%	16%
Held "Flashback Fridays" or catch up days for students who were "Not Yet Proficient"	20%	12%	20%	12%	20%	16%

Note: N = 25 teachers

The PARLO system did not envision teachers using all four of the above techniques all the time. Instead, they were expected to use some technique to encourage students to reassess after further learning. Therefore, we asked, “How frequently did teachers use at least one of the

above four techniques?” Consequently, for each teacher we calculated the maximum score across the four techniques. Results are reported in Table A9.

Table A9

How Often Year 3 PARLO Teachers Reported Using At Least One of the Four Re-learning Techniques (Percent of Teachers in Each Category)

Never	Rarely	A Few Times	Fairly Often	Frequently	Almost Every Week
0%	0%	12%	12%	28%	48%

Content Tests Used

Addenda 1 and 2 display the algebra and geometry post tests used in this study. The Algebra PreTest contained 31 items that were identical to those displayed in Addendum 1 but presented in a different order. Similarly, the Geometry PreTest contained 35 items that were identical to those displayed in Addendum 2 but presented in a different order.

After all tests had been administered but before student scores were computed, we dropped two items from the algebra exam displayed in Addendum 1. We dropped item #7 on the post test, and the identical item #30 on the pretest, because either response option “a) $1/5$ ” or response option “b) $1/25$ ” could be correct, depending on the student’s interpretation of the question. We dropped item #22 on the post-test and the identical item #4 on the pretest, because including the item reduced reliability (Cronbach’s alpha) on both the pretest and post-test. Fewer than 25% of students answered the item correctly on either test, indicating that students were guessing the answer. It appeared that neither Treatment nor Control students had the opportunity to learn the content of this item, which tested students’ ability to find the Least Common Multiple of Monomials.

When scoring the tests, we awarded 1 point for each correct answer, 0 points for each incorrect answer, and 0.25 points for items left blank. We scored it in this manner because some students chose to leave items they could not answer blank, whereas other students guessed and thus had a 25% chance of getting the items correct. As compared to scoring blank items 0 points, when we assigned 0.25 points for blank items the pretest and post-test scores correlated more highly with each other, more highly with students' responses on the Attitudes Towards Mathematics Inventory (ATMI) four subscales, and, for the subset of students who had state testing data available, more highly with state-administered 8th grade math test scores.

We also investigated whether a piecewise linear model might be better to use for modeling the covariance of pretest with post-test. For algebra, we found that the most effective model had one slope for all students who correctly answered 7 or fewer of the 29 items on the test (i.e., scored 24% or less) and a different slope for all students who correctly answered 8 or more of the questions. Indeed, the piecewise model showed a near-zero, slightly negative slope for scores below 25%. Consequently, before analysis, we rescaled the algebra pretest score by subtracting 7 from the number correct and setting all negative scores equal to 0. In contrast to algebra, for geometry we found that a piecewise model was not superior to a simple linear model in predicting post-test from pretest. Further, rescaling the geometry test did not improve appropriate correlations. Consequently, we did not rescale the geometry pretest in this way before performing analyses. Neither did we rescale either of the post-tests in this way.

Once we identified our analytic sample of 2,736 students, we used linear transformations to convert the rescaled algebra pretests, as well as the geometry pretests, and the algebra and geometry post-tests, into z-scores with a mean of 0 and a standard deviation of 1 within the analytic sample.

We computed a new “Post” score that we used in the analysis, defined as the student’s algebra post-test z-score for algebra students, and the student’s geometry post-test z-score for geometry students. Similarly, we computed a new “Pretest” score, defined as the student’s algebra pretest test z-score for algebra students and the student’s geometry pretest test z-score for geometry students. We also computed the square of each student’s pretest score.

Attitudes Toward Mathematics Inventory (ATMI)

Addendum 3 displays the four ATMI subscales. When completing the ATMI, some students left one or more items on one or more subscales blank while completing the rest of the subscale. Because items within a subscale differed in their mean and variance, we could not use a simple average of these students’ responses as a valid reflection of their score on the subscale. For example, mean responses on the 5-item Motivation scale ranged from a low of 2.7 for item #33 to a high of 4.0 for item #2. In order to compute a meaningful average score even for students who completed only part of a subscale, we first converted each item to a z-score with mean of 0 and standard deviation of 1. For each subscale, we then averaged all available z-scores for each student subscale. Using only students with complete information on each subscale, we used OLS regression to compute a linear transformation from the average z-scores back to a 1-5 scale computed by averaging raw scores. We then used this linear transformation to translate each student’s subscale score back to an easily interpretable 1-to-5 scale.

Teacher Interviews

Addendum 4A displays the 23 questions in the Year 1 teacher interview. Non-exit interviews in later years asked the same questions, except that the reference to the software program Ease was replaced with a reference to its replacement, PARLO Tracker. Teachers interviewed at the end of the year in the last year of their school’s participation in the project were administered the 17-question Exit Interview, which is displayed in Addendum 4B.

Open-ended Student Survey

Addendum 5 displays the 8 open-ended questions on the student survey. The survey was administered during the second year of program operation, i.e., 2011-12. The 678 respondents came from 6 treatment schools with 18 PARLO teachers and 31 PARLO class sections.

Choosing Covariates for the Baseline Model.

Unlike a quasi-experiment, our experimental study would be expected to produce unbiased results even if we selected an incomplete covariate model. Nonetheless, by selecting appropriate covariates we were able to reduce random variance in our model, thus improving the precision of our model and narrowing confidence intervals around estimated results. It was also important that we follow a systematic procedure to select covariates before addressing our main research questions. By doing so we could avoid the temptation of unconsciously choosing a set of covariates that most closely produced our preferred results.

We chose a variance structure following procedures described in the main article under the heading *Variance Components and Covariates Selected for Use*. After deciding on a variance structure, we investigated all student-level covariates to determine which were statistically significant predictors of the content post-test. Next, we tested school-level covariates. Since intervention assignment was based on school, in order to maximize model precision, we predetermined to include any school-level variable that reduced between-school variance. The only school-level covariate that met that criterion was *proportion low-income students*, which had been measured during the 2010-11 school year. In analyses, the school proportion low-income students was rescaled by subtracting the mean across schools. Finally, at the suggestion of an anonymous reviewer, we added four blocking variables: “Catholic-school,” “Charter-school,” “Cohort Two Public School, students tested in 2011-12,” and

“Cohort Two Public School, students tested in 2012-13” as fixed effects in the model. The Cohort 1 Public School random-assignment block was the reference group.

To investigate PARLO treatment impacts on mindsets and motivation we used the same covariates we had used to investigate the PARLO treatment impact on academic outcomes, with two exceptions made for substantive reasons. First, we dropped the quadratic term for the algebra pretest because it was not significant and made analyses less parsimonious. Second, we used baseline scores from all four ATMI subscales, instead of only from the Self-Confidence (i.e., Expectancy for Success) subscale.

Full Information Maximum Likelihood Analysis of PARLO impact on Academic Performance.

Instead of deleting cases with missing pretests or post-tests, an alternative is to handle missing data by employing a Full Information Maximum Likelihood (FIML) approach that takes advantage of hierarchical linear modeling. In this approach, one treats pretests and post-tests as repeated measures of an underlying construct (Allison, 2012; Raudenbush & Bryk, 2002). Test scores are nested within students, who are nested within higher level clusters.

For this analysis, we modeled test scores (pretest or post-test) as Level 1, Student as Level 2, Course within Teacher within Study-Year as Level 3, and School as Level 4. We converted pretest scores to z-scores within the data set, with a mean of 0 and a standard deviation of 1. In order to maintain comparability to the analysis using listwise deletion that we reported in the main article, we converted each student’s post-test score to the same scale used in Table 1 of the main article, by subtracting the mean for the 2,736 students in the listwise-deleted analytic data set and dividing by the standard deviation of the post-test score for those 2,736 students.

We created four new dichotomous variables as indicators of the four types of content

test: Algebra PreTest, Geometry PreTest, Algebra Post-test, or Geometry Post-test. We allowed each of the four indicators to vary randomly within Level 3 and Level 4. Correlations among the random variables enabled us to control for the covariation between geometry pretest and post-test and between algebra pretest and post-test. Note that unlike the analysis using listwise deletion, this analysis did not model a quadratic term for the effect of the algebra pretest test on the algebra post-test.

We also created a *time* fixed effect variable, coded 0 for the pretest test and 1 for the post-test. In this model, the PARLO treatment effect on post-test after controlling for the pretest would be the *PARLO x time* interaction.

In order to avoid deleting cases (students) who were missing data for other covariates, we aggregated each remaining student level covariate, i.e., Female, White or Asian, and Baseline Expectancy of Success, to the school level and used the school mean as the covariate. This approach was appropriate since school was the unit of assignment and thus between-school variation was more important than between-student variation in accounting for any differences between PARLO and Control schools. We recentered all school level variables so that 0 represented the mean of the school means. We used the REML option of the lmer command in the lmerTest package of the R programming language to run the analysis.

Similar to the PARLO treatment effect, the effect of each covariate on the post-test after controlling for the pretest, would be the *Covariate x time* interaction.

By taking this approach, our FIML model was able to use data from every student who completed a content pretest or a content post-test. Students with missing observations on one or the other variable nonetheless contributed to our estimate of within-cluster means and slopes. Consequently, unlike the analysis reported in the main article, this analysis sample includes

joiners, i.e., students who were not present during the first week of school when the pretest and attitudes survey were administered. The sample also includes students who did not stay for the entire school year, and students who completed the school year but for some reason did not provide post-test scores. Data for this analysis was provided by 4,116 students and 73 teachers from 29 schools: 2,385 students and 43 teachers at 14 PARLO treatment schools and 1,731 students and 30 teachers at 15 control schools. Table A10 summarizes the results.

Table A10

FIML Analysis of PARLO Treatment Main Effect

<i>Fixed Effects</i>	<i>Estimate</i>	<i>Standard Error</i>	<i>Degrees of Freedom</i>	<i>p-value</i>
Intercept	-0.04	0.11	29	.72
<i>Main Effects (impacts on Pretest score)</i>				
PARLO Treatment	0.02	0.09	24	.82
Catholic Girls' School	-0.84	0.44	42	.060
Charter School	0.40	0.25	47	.12
Cohort 2 Public School, 2011-2012 school year	0.05	0.12	24	.69
Cohort 2 Public School, 2012-2013 school year	0.19	0.13	27	.15
Geometry Student	-0.09	0.11	8	.42
School Mean baseline Expectancy	0.82	0.28	10	.016
School Proportion White or Asian	0.99	0.28	31	.0015
School Proportion Female	2.19	0.72	40	.0045
School Proportion Disadvantaged	0.21	0.32	38	.51
Time	-0.11	0.16	27	.50
<i>Interaction Effects (impacts on Posttest after controlling for Pretest score)</i>				
PARLO Treatment x time	0.28	0.12	32	.030
Catholic Girls' School x time	0.16	0.61	27	.80
Charter School x time	-0.52	0.34	27	.14
Cohort 2 Public School, 2011-2012 school year x time	-0.32	0.15	52	.033
Cohort 2 Public School, 2012-2013 school year x time	-0.27	0.15	59	.076
Geometry Student x time	0.16	0.12	8	.22
School Mean baseline Expectancy x time	0.57	0.28	87	.047
School proportion White or Asian x time	-0.25	0.40	24	.53
School proportion Female x time	0.18	1.02	27	.86
School Proportion Disadvantaged x time	-0.10	0.46	24	.63
<i>Number of Observations</i>				
School	29			
Course by Teacher by Year	95			
Student ID	4,116			
Total Observations (pre + post)	6,944			

Notes: The intercept of the model can be interpreted as the expected baseline test z-score for the covariate reference values, i.e., for an algebra student in a control classroom at a Cohort 1 Public school with average scores on the ATMI expectancy of success measure and on proportion White or Asian, proportion Female, and proportion Disadvantaged. The interactions with Time

estimate the impact of each independent variable on the post test, after controlling for pretest scores and school-level aggregate covariates.

As shown in Table A10 the maximum likelihood estimate of the PARLO program's impact on algebra and geometry content test scores was statistically significant ($p = .030$) and the effect size was close to what we estimated using listwise deletion: +0.28 standard deviations.

Other things worth noting from the table: Schools with high means on the Expectancy subscale had higher scores than other schools on the pretest and tended to improve achievement more from pretest to post-test. Schools with high proportion White or Asian had higher scores on the pretest but appeared to gain at roughly the same rate as other schools from pretest to post-test. There may also have been a difference among assignment blocks in achievement growth, with Cohort 2 schools showing less pre- to post-test growth during 2011-12 than Cohort 1 schools. Finally, the reader is cautioned that because the Catholic schools were the only gender-specific schools in our data set, the variables "Catholic Girls' School" and "School Proportion Female" were highly correlated. Those two variables were entered as controls, but their slopes should not necessarily be interpreted as meaningful.

Additional Details for Listwise-Deletion Analysis of PARLO Impact on Student Mindsets and Motivation

Table A11 provides a more detailed breakdown of information contained in Table 2 of the main article, including covariate coefficients and variance components. The most interesting finding from this more detailed breakdown is that females in the sample had somewhat lower scores on all four motivation measures than did males.

Table A11

Effects of PARLO Treatment on Students' Motivational Antecedents and Long-term Motivation in Math Class (analytic data set prepared using listwise deletion)

	<i>Intrinsic Value</i>			<i>Utility Value</i>			<i>Expectancy of Success</i>			<i>Long-term Motivation</i>		
	β	SE	p-value	β	SE	p-value	β	SE	p-value	β	SE	p-value
<i>School Level Fixed Effects</i>												
Intercept	3.20	0.06	<.0001	3.62	0.05	<.0001	3.40	0.05	<.0001	3.19	0.07	<.0001
Catholic Girls' School	-0.13	0.10	.19	0.12	0.08	.88	-0.05	0.08	.55	0.06	0.10	.59
Charter School	-0.14	0.13	.31	-0.23	0.11	.043	-0.22	0.12	.072	-0.24	0.15	.11
Cohort 2 Public School, 2011-12 school year	0.04	0.07	.61	0.05	0.06	.38	0.03	0.05	.64	0.11	0.08	.20
Cohort 2 Public School, 2012-13 school year	-0.04	0.07	.54	-0.04	0.06	.47	0.01	0.05	.88	-0.05	0.08	.56
Proportion Disadvantaged	0.06	0.11	.57	0.22	0.09	.026	0.04	0.08	.61	0.16	0.12	.21
PARLO Treatment	-0.06	0.06	.31	-0.02	0.05	.60	-0.06	0.04	.14	-0.08	0.06	.23
<i>Course Level Fixed Effect</i>												
Geometry Student	0.01	0.07	.85	-0.03	0.03	.63	0.03	0.07	.72	0.13	0.07	.069
<i>Student Level Fixed Effects</i>												
Baseline Expectancy	0.14	0.03	.91	0.003	0.03	.91	0.48	0.03	<.0001	-0.01	0.03	.87
Baseline Intrinsic Value	0.43	0.03	.10	0.05	0.03	.10	0.13	0.03	<.0001	0.19	0.03	<.0001
Baseline Long-term Motivation	0.11	0.03	<.0001	0.10	0.03	<.0001	0.08	0.03	.0027	0.37	0.03	<.0001
Baseline Utility Value	-0.08	0.03	<.0001	0.44	0.03	<.0001	-0.08	0.03	.0032	0.08	0.03	.011
White or Asian	-0.01	0.03	.65	0.01	0.03	.65	0.03	0.03	.33	0.006	0.04	.87
Female	-0.08	0.03	.0021	-0.08	0.03	.0021	-0.12	0.03	<.0001	-0.08	0.03	.0086
Pretest	0.07	0.01	.0002	0.05	0.01	.0002	0.10	0.01	<.0001	0.03	0.02	.038
Geometry X Pretest	0.07	0.05	.82	0.01	0.04	.82	0.03	0.04	.56	0.04	0.05	.39
<i>Random Effects</i>												
	<i>n Variance</i>			<i>n Variance</i>			<i>n Variance</i>			<i>n Variance</i>		
School	29	0.008		29	0.003		29	<0.0001		29	0.014	
Course x Teacher x Year	78	0.012		78	0.011		84	0.018		84	0.012	

Residual	2,694	0.44	2,697	0.36	2,694	0.40	2,698	0.47
----------	-------	------	-------	------	-------	------	-------	------

Note. This analysis utilized data from 2,698 students, 65 teachers, and 29 schools. All non-dichotomous variables are grand mean centered. Dependent variables are measured on a 1-5 Likert scale. β = effect estimate in raw 1-5 units; SE = Standard Error; n = number of observations. The intercept of the model can be interpreted as the end-of-year survey scale score for students at average values for pretest, baseline scores on all four ATMI subscales, and school-level proportion disadvantaged; and in the reference group for all dichotomous variables (i.e. in the Cohort 1 Public School randomization block; in a control school, taking algebra, not White/Asian, and not Female).

Full Information Maximum Likelihood Analysis of PARLO Impact on Intrinsic Value, Utility Value, Expectancy of Success, and Long-term Motivation

As we did with our analysis of PARLO effects on academic performance, we also addressed PARLO effects on mindsets and motivation using FIML instead of listwise deletion of missing data. The approach was the same as the maximum likelihood analysis described above. For each survey subscale, a student's baseline score and post-score were treated as repeated measures. (Note: instead of recentering them as z-scores, we kept the outcome ATMI variables on a 1-5 Likert scale.) The nesting structure was the same as the one described above for the maximum likelihood analysis of program effects on academic performance.

We created two new dichotomous variables as indicators of “baseline” or “post-score,” allowing each of the two variables to vary randomly within Level 3 (Course within Teacher within Study-year) and Level 4 (school). Correlations among the two random variables enabled us to control for the covariation between baseline and post-score. We also created a fixed-effect “time” variable.

The covariates for this analysis were the same ones used in Table 1 of the main article, albeit (with the exception of the baseline score on the target attitude) aggregated to the school level and centered around the mean of the school means. For the algebra and geometry pretests, we used the school mean z-score for students who took the pretest and set the value to zero for students who did not take that particular pretest. As with the maximum likelihood analysis described above, the parameter of interest was the *PARLO* \times *time* interaction effect on the outcome variable. Data for these analyses was provided by 4,189 students and 68 teachers from 29 schools: 2,412 students and 40 teachers at 14 PARLO treatment schools and 1,777 students and 28 teachers at 15 control schools. Table A12 summarizes the estimated PARLO effect on

each ATMI subscale. More detailed output tables are available from the first author upon request.

Table A12
FIML Estimate of PARLO Effects on ATMI Subscales (1-5 Likert Scales)

<i>Dependent Variable (Subscale)</i>	<i>PARLO effect Estimate</i>	<i>Standard Error</i>	<i>Degrees of Freedom</i>	<i>95% conf. interval</i>	<i>p-value</i>	<i>Effect size in SD units</i>
Intrinsic Value	-0.13	0.08	15	(-0.30, +0.04)	.15	-0.15
Utility Value	-0.08	0.05	13	(-0.19, +0.03)	.15	-0.11
Expectancy	-0.07	0.07	57	(-0.21, +0.07)	.29	-0.08
Long-term Motivation	+0.01	0.09	14	(-0.18, +0.20)	.90	+0.01

As shown in Table A12, the maximum likelihood estimate of the PARLO program's impact was similar to what we estimated using listwise deletion, with confidence intervals overlapping zero for all four dependent variables. The most noticeable differences between the results of the maximum likelihood analysis and those of the analysis using listwise deletion are that in the maximum likelihood analysis the point-estimate slope for Long-term Motivation was positive instead of negative, and the point-estimate slope for Intrinsic Value, at -0.13, was of somewhat greater magnitude than the point-estimate calculated when using listwise deletion.

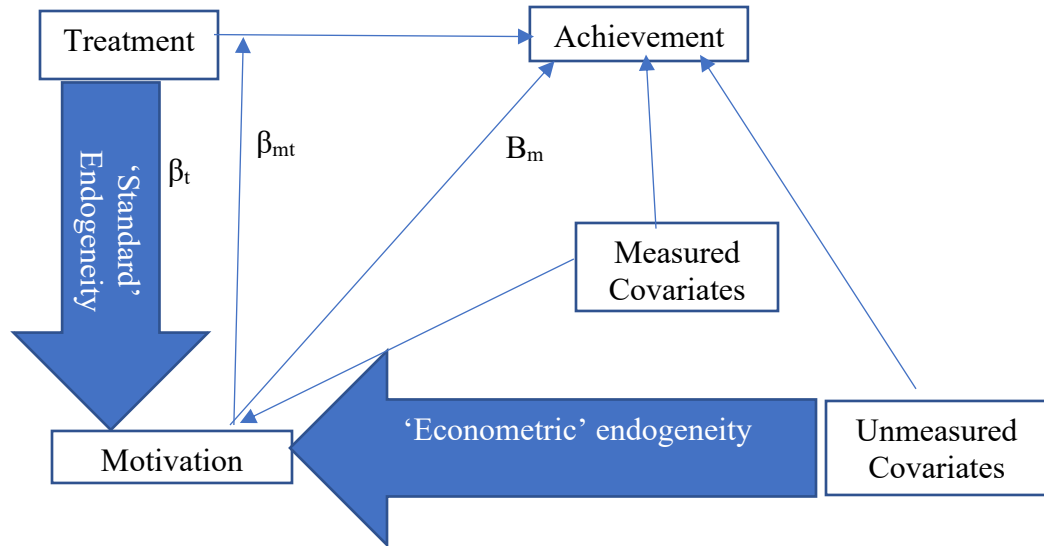
Details About the Moderation Analysis (Table 3 of the Main Article)

Potential problems created by endogeneity. The term “endogenous” has two distinct meanings in the research literature. One form of endogeneity, which we refer to as *standard endogeneity* is explained as follows in the International Journal of Social & Behavioral Sciences: “Exogenous variables are thought of as causes, endogenous as their effects. But there is no necessary connection; one may use the terms without implying causality A simpler terminology for exogenous and endogenous is independent and dependent variables.” (Peterson, 2001). The other form of endogeneity, which we refer to as *econometric endogeneity*, is

commonly used in econometrics and is defined as a predictor variable that is correlated with the error term (Bun and Harrison, 2018). Our measures of average motivation scores might be endogenous in either sense.

Figure A2 displays the situation. The two thicker arrows are used to display potentially endogenous relationships. Our measures of average intrinsic value, utility value, expectancy of success, and long-term motivation are potentially “standard endogenous” because they could be impacted by Treatment. Including any of these standard-endogenous variables in the model would make main effects very hard to interpret. The main effect of Treatment, β_t would change depending on where we centered the motivation scores. Because motivation scores are endogenous, the average score for Treatment and Control would be different, and we could change the reported Treatment effect depending on whether we centered motivation on average for the Treatment group, average for the Control group, average for the entire study population, or somewhere else. Similarly, the main effect of a motivation variable, β_m , would change depending on whether our reference group was Treatment or Control. Fortunately, however, we use the potentially standard-endogenous motivation related measures only in our moderation analysis. That analysis reports only the interaction slope β_{mt} . The interaction slope, and its statistical significance, is invariant no matter where we center other variables in the model. Thus, there is no difficulty in interpreting the results.

Figure A2 Standard Endogeneity and Econometric Endogeneity



Our measures of average intrinsic value, utility value, expectancy of success, and long-term motivation are also potentially “econometric endogenous”. For example, one of the scores might be predicted by an unmeasured variable like grit (Duckworth, et al., 2007), which might also predict student achievement. An econometric-endogenous variable will be correlated with the error term, but calculation procedures for Multilevel and OLS regression calculate regression slopes that force the error term to be uncorrelated with any of the covariates. Doing so usually forces the regression slopes to be biased.

Once again, however, our moderation analysis is valid. This is because, while main effect estimates might be biased by inclusion of an econometric-endogenous variable, as explained below the estimated slope for the *interaction* between an endogenous variable and a binary Treatment will be unbiased, as long as the endogenous variable is homoscedastic.

Confirming the Homoscedasticity of end-of-year scores on Intrinsic Value, Utility Value, Expectancy of Success, and Long-term Motivation in mathematics.

In the context of OLS regression, Bun and Harrison (2018) demonstrated that using an interaction term to test the moderation effects of an endogenous variable on a Treatment would

provide consistent results assuming two conditions: 1) the Treatment was binary; and 2) the endogenous moderator is homoscedastic conditional on the Treatment and on the full set of covariates in the model. Further, they presented Monte Carlo simulations confirming that under these constraints the interaction term was likely to be unbiased as well as consistent. Finally, although additional work has not yet been finalized or published, preliminary Monte Carlo analyses indicate that the results can be generalized to multilevel regression, (Bun, personal communication, January, 2023).

Consequently, we investigated the heteroscedasticity of the potentially endogenous portion of the moderator variables of interest, i.e., the post-test scores for Intrinsic Value, Utility Value, Expectancy of Success, and Long-term Motivation in mathematics.

We employed a modification of the Breusch and Pagan (BP) test, as updated Koenker (1981). The BP test evaluates heteroscedasticity by first regressing the variable of interest on the other variables in the model and saving the residuals. The BP test then uses a Chi-square statistic with degrees of freedom equal to the number of predictors in the model to evaluate whether the same model is a significant predictor of the squared residuals. If the squared residuals are a function of the predictors in the model, then the assumption that the variable is homoscedastic can be rejected. This is the method we used.

First, we used a multilevel model with the same variance components used in our main effects analysis to regress each of the potentially endogenous variables in our model (the post-test scores on intrinsic value, utility value, expectancy of success, and long-term motivation) on the following 13 independent variables: indicator variables for PARLO Treatment, for the four blocking variables (Catholic School, Charter School, Public School 2011-12, and Public School 2012-13), for Geometry (vs. Algebra) and for demographic variables White-or-Asian and

Female. We included student pretest, student pretest squared, Geometry-times-pretest, and Geometry-times-pretest-squared. Finally, we included school-level proportion disadvantaged. We saved the residuals and squared them.

Finally, we followed procedures described by Raudenbush and Bryk (2002) to test whether the squared residuals were related to the thirteen predictor variables in our model. Using Maximum Likelihood estimation and the same variance structure, we regressed the squared residuals on the thirteen predictor variables. We tested the significance of the relationship by subtracting -2 times the Log-likelihood, called the *deviance*, of the 13-variable model from the deviance of a null model using only the variance components. The test statistic was Chi-Square with 13 degrees of freedom. The results are shown in Table A13. As can be seen from the table, we could not reject the assumption of homoscedasticity for any of the four variables.

Table A13
Heteroscedasticity Check for Post-test Value Variables
(Significant p-value Indicates Heteroscedastic)

Variable	-2 Log Likelihood, null model	-2 Log Likelihood, with 13 Predictors	Chi-square value (13 df)	p- value
Intrinsic Value (Post)	6232.5	6215.6	16.9	0.20
Utility Value (Post)	5580.2	5564.3	15.9	0.25
Expectancy of Success (Post)	6025.4	6004.6	20.8	0.08
Long-term Motivation (Post)	6322.8	6308.4	14.4	0.35

It is important to note here that, when using potentially econometric-endogenous variables as moderators, we are generalizing results from OLS regression (Bun and Harrison, 2018) and applying them to multilevel regression. Preliminary work indicates that the results can be generalized to multilevel regression, but this research has not yet been finalized or published (Bun, personal communication, January, 2023). Consequently, it is important to conduct a

sensitivity analysis to see if we achieve similar findings when using only the endogenous baseline scores as moderators. As explained in the main article, the baseline scores are theoretically less important as moderators than are the average scores, and we expected them to have a weaker moderation effect, but nevertheless to have an effect in the same direction.

Sensitivity Analysis using only Baseline Scores as Moderator

Table A14 reports the results of the sensitivity analysis we conducted using only endogenous variables to investigate the moderation effects of student motivation. As can be seen in the table, the point-estimate moderation effects of baseline moderation were somewhat smaller but in the same direction as the moderation effects of average moderation across the year, as reported in Table 3 of the main article and Tables A15 through A18 below. The effects, however, are not statistically significant at the $p=.05$ level.

Table A14
Interactions with Treatment Condition: Do Students' *Baseline* Perception of Intrinsic Value, Utility Value, Expectancy of Success, and Long-term Motivation Moderate PARLO's Impact on Mathematics Achievement?

<i>Interaction</i>	<i>n</i>	<i>Effect size Estimate</i>	<i>Standard Error</i>	<i>Degrees of Freedom</i>	<i>95% conf. interval</i>	<i>p-value</i>
PARLO x Baseline Intrinsic Value	2,736	0.06	0.03	2,685	(-0.01, +0.12)	0.11
PARLO x Baseline Utility Value	2,736	0.08	0.05	2,675	(-0.01, +0.17)	0.088
PARLO x Baseline Expectancy	2,736	0.08	0.04	2,691	(-0.00, +0.15)	0.052
PARLO x Baseline Long-term Motivation	2,736	0.06	0.04	2,678	(-0.02, +0.14)	0.13

Notes: This analysis utilized data from 65 teachers, and 29 schools. "Mathematics Achievement" was defined as the Algebra Post-test z-score for algebra students and the Geometry post-test z-score for geometry students. n = number of students. Reported results controlled for the following covariates: Assignment block for randomization; School-level proportion disadvantaged; course assignment (geometry or algebra), student-level gender, race, pretest, pretest-squared, main effects of baseline scores of the motivation subscale being studied, and main effects of the PARLO treatment.

Detailed Expansions of Table 3 in the Main Article, Including Covariate Coefficients and Variance Components

Tables A15 through A18 provide an expanded view of the analyses used to produce results displayed in Table 3 of the main article, estimating slopes of the covariates, as well as standard errors and p-values.

We caution that covariate slopes in tables A15 through A18 could potentially be biased. Based on Bun and Harrison (2018) the interaction terms we were investigating are likely to be consistent and unbiased. However, if the average values for the four motivation related variables are in fact econometric-endogenous (i.e., correlated with the error term), we cannot guarantee that the other parameter estimates reported in the four tables below will be consistent or unbiased. Indeed, Bun and Harrison (2018) found that, unlike the interaction effect presented in bold, in each table the estimated slopes for the main effects of the motivation variable and of Treatment are likely to be biased.

Table A15

PARLO x Intrinsic Value Interaction Effect on Content Knowledge

<i>Fixed Effects</i>	<i>Estimate</i>	<i>Standard Error</i>	<i>Degrees of Freedom</i>	<i>p-value</i>
Intercept	-0.65	0.13	19	<.0001
PARLO Treatment	0.32	0.12	14	.021
Catholic Girls' School	0.33	0.22	20	.14
Charter School	-0.02	0.28	30	.96
Cohort 2 Public School, 2011-2012 school year	0.26	0.15	14	.11
Cohort 2 Public School, 2012-2013 school year	0.37	0.15	14	.030
Geometry Student	0.05	0.18	77	.77
White or Asian	0.17	0.04	2,444	<.0001
Female	0.10	0.03	2,463	.0018
Pretest	0.27	0.02	2,505	<.0001
Pretest Squared	0.03	0.01	2,467	.035
Geometry x Pretest	0.23	0.06	2,512	.0002
Geometry x Pretest-squared	-0.02	0.04	2,490	0.48
School Proportion Disadvantaged	-0.58	0.24	12	.03
Intrinsic Value (baseline & post average)	0.17	0.03	2,486	<.0001

PARLO x Intrinsic Value Interaction	0.09	0.04	2,483	.026
<i>Random Effects</i>	<i># Observations</i>	<i>Variance</i>		
School	29	0.03		
Course x Teacher x Year	81	0.13		
Residual	2,529	0.52		

Notes: This analysis utilized data from 2,529 students, 62 teachers, and 29 schools.

Table A16

PARLO x Utility Value Interaction Effect on Content Knowledge

<i>Fixed Effects</i>	<i>Estimate</i>	<i>Standard Error</i>	<i>Degrees of Freedom</i>	<i>p-value</i>
Intercept	-0.65	0.13	18	<.0001
PARLO Treatment	0.33	0.12	13	.019
Catholic Girls' School	0.28	0.22	20	.21
Charter School	0.01	0.29	30	.96
Cohort 2 Public School, 2011-2012 school year	0.27	0.15	13	.10
Cohort 2 Public School, 2012-2013 school year	0.38	0.15	13	.026
Geometry Student	0.09	0.18	78	.62
White or Asian	0.17	0.04	2,437	<.0001
Female	0.08	0.03	2,467	.017
Pretest	0.29	0.02	2,508	<.0001
Pretest-squared	0.03	0.01	2,470	.034
Geometry x Pretest	0.25	0.06	2,515	<.0001
Geometry x Pretest-squared	-0.03	0.04	2,492	.47
School % Disadvantaged	-0.59	0.23	11	.027
Utility Value (baseline & post average)	0.10	0.04	2,485	.0069
PARLO x Utility Value Interaction	0.11	0.05	2,483	.028
<i>Random Effects</i>	<i># Observations</i>	<i>Variance</i>		
School	29	0.03		
Course x Teacher x Year	81	0.14		
Residual	2,532	0.54		

Notes: This analysis utilized data from 2,532 students, 62 teachers, and 29 schools.

Table A17

PARLO x Expectancy of Success Interaction Effect on Content Knowledge

<i>Fixed Effects</i>	<i>Estimate</i>	<i>Standard Error</i>	<i>Degrees of Freedom</i>	<i>p-value</i>
Intercept	-0.66	0.13	18	<.0001
PARLO Treatment	0.33	0.12	13	.020
Catholic Girls' School	0.34	0.22	19	.13
Charter School	-0.00	0.28	29	.99
Cohort 2 Public School, 2011-2012 school year	0.24	0.15	13	.13
Cohort 2 Public School, 2012-2013 school year	0.35	0.15	13	.039
Geometry Student	0.05	0.18	78	.79
White or Asian	0.17	0.04	2,441	<.0001
Female	0.14	0.03	2,464	<.0001
Pretest	0.25	0.02	2,504	<.0001
Pretest-squared	0.02	0.01	2,466	.041
Geometry x Pretest	0.23	0.06	2,513	.0001
Geometry x Pretest-squared	-0.02	0.03	2,489	.53
School % Disadvantaged	-0.56	0.23	12	.035
Expectancy (baseline & post average)	0.24	0.03	2,484	<.0001
PARLO x Expectancy Interaction	0.13	0.04	2,488	.002
<i>Random Effects</i>	<i># Observations</i>	<i>Variance</i>		
School	29	0.03		
Course x Teacher x Year	81	0.13		
Residual	2,529	0.51		

Notes: This analysis utilized data from 2,529 students, 62 teachers, and 29 schools

Table A18

PARLO x Long-term Motivation Interaction Effect on Content Knowledge

<i>Fixed Effects</i>	<i>Estimate</i>	<i>Standard Error</i>	<i>Degrees of Freedom</i>	<i>p-value</i>
Intercept	-0.65	0.13	19	<.0001
PARLO Treatment	0.33	0.12	14	.018
Catholic Girls' School	0.28	0.22	20	.22
Charter School	0.00	0.29	30	.99
Cohort 2 Public School, 2011-2012 school year	0.26	0.15	13	0.11
Cohort 2 Public School, 2012-2013 school year	0.37	0.15	13	0.029
Geometry Student	0.07	0.18	78	.70
White or Asian	0.17	0.04	2,445	<.0001
Female	0.08	0.03	2,468	.0079
Pretest	0.29	0.02	2,510	<.0001
Pretest-squared	0.03	0.01	2,471	.033
Geometry x Pretest	0.23	0.06	2,517	<.0001
Geometry x Pretest-squared	-0.03	0.04	2,4984	0.47
School % Disadvantaged	-0.59	0.24	12	.028
Long-term Motivation (baseline & post average)	0.13	0.03	2,489	.0002
PARLO x Long-term Motivation Interaction	0.07	0.04	2,485	.13
<i>Random Effects</i>	<i># Observations</i>		<i>Variance</i>	
School	29		0.03	
Course x Teacher x Year	81		0.14	
Residual	2,533		0.54	

Notes: This analysis utilized data from 2,564 students, 62 teachers, and 29 schools

Benjamini-Hochberg (BH) Corrections to Account for Multiple Comparisons

The statistical significance of findings and the number of Type I errors (incorrect rejections of the null hypotheses) may be inflated when an analysis conducts multiple hypothesis tests simultaneously. The traditional way of addressing this issue is to employ the Bonferroni adjustment, which modifies the critical p value for the hypothesis test to p/m , where m is the number of hypothesis tests being made. The Bonferroni adjustment controls the *family-wise error rate*, which keeps the probability of at least one Type I error to less than the critical value, p .

Following DOE (2020) we used the less conservative Benjamini-Hochberg (BH) procedure, which controls *the false discovery rate*, which keeps the *expected proportion of “statistically significant” results that are Type I errors* (i.e., false rejections of the null) to less than the critical value, p . In many situations, the BH procedure is more intuitively appealing than the Bonferroni. For example, using Bonferroni to control for familywise error rate, if researcher A conducted one significance test of a null hypothesis at the .05 level and achieved $p = .04$, researcher A would report significant results. Meanwhile, if researcher B conducted five significance tests of null hypotheses at the .05 level and all five results achieved p -values between $p = .01$ and $p = .04$, then using a Bonferroni adjustment researcher B would report “no significant results”. The BH procedure, by controlling for false discovery rate instead of familywise error rate, avoids this type of problem.

To perform the BH procedure on N hypothesis tests, one first ranks the hypothesis tests in ascending order of statistical significance. For each of our N tests, one then adjusts the critical value to be $0.05 * m/N$, where m was the rank-order of the hypothesis test’s p -value. If the p -value for any of the tests is less than its critical value, we conclude that that test and all tests with lower p -values are statistically significant. Doing this ensures that the expected proportion of Type I errors (false discovery rate) will be less than 5%.

We conducted three BH analyses. The first investigated the significance of Treatment x Covariate Interaction effects on student mathematics achievement. These analyses were done by adding interactions, one at a time, to the analysis reported in Table 1 of the main article. See Table A19 below. The second BH analysis investigated the significance of Treatment x Covariate interactions on our four motivation-related variables. These analyses were done by

adding interactions, one at a time, to the analyses reported in Table A9 above and summarized in Table 2 of the main article. See the discussion below.

The third BH analysis investigated the significance of the four moderation effects reported in Tables A15 through A18 above and summarized in Table 3 of the main article. See Table A20 below.

BH Analysis of Interaction Effects on Content Post-test Scores. There were five possible interaction effects on student mathematics achievement. The p-values for each are reported in Table A19 below. As can be seen, after controlling for the false discovery rate, none of the interactions achieved significance at the 0.05 level, or indeed at the 0.10 level.

Table A19
BH Adjustment to Test Treatment x Covariate Interaction Effects on Our Achievement Measures

<i>Interaction variable used to test moderation</i>	<i>p-value from Table 3 of main article</i>	<i>BH critical value .05 level</i>	<i>BH critical value .10 level</i>
Baseline Expectancy of Success x PARLO	0.052	.01	.02
Geometry Student x PARLO	0.056	.02	.04
Female x PARLO	0.183	.03	.06
White-or-Asian x PARLO	0.586	.04	.08
Pretest x PARLO	0.866	.05	.10

BH Analysis of Interaction Effects on Intrinsic Value, Utility Value, Expectancy of Success, and Long-term Motivation in Mathematics. There were 32 possible interaction effects on student measures of long-term motivation or motivational antecedents, 8 per dependent variable. Of the 32 interactions, the most significant had a p-value of .054, so clearly none of the interactions came close to meeting the BH criterion for significance after controlling for false discovery. We did not prepare a table reporting the 32 p-values, because this particular set of interactions was so clearly non-significant.

BH Analysis of Interaction Effects on Content Post-test Scores. There were four possible interaction effects on student mathematics achievement. The p-values for each are reported in Table A20 below. As can be seen in the table, the p-value for the Utility x PARLO interaction, at .028, was smaller than the BH critical value for significance while maintaining a false-discovery rate of $p = .05$. Consequently, we concluded that Utility Value moderated PARLO effects on academic achievement to a statistically significant extent, as did all moderation effects with smaller p-values, i.e., Intrinsic Value x PARLO and Expectancy of Success x PARLO.

Table A20
BH Adjustment to Test Moderation Effects

<i>Interaction variable used to test moderation</i>	<i>p-value from Table 3 of main article</i>	<i>BH critical value, .05 level</i>
Expectancy of Success x PARLO	.0019	.0125
Intrinsic Value x PARLO	.026	.025
Utility Value x PARLO	.028	.0375
Long-term Motivation x PARLO	.13	.05

Addendum 1: Algebra Post Test

Post-Test

Algebra 1- End of Course Assessment*

*Adapted from *Virginia Standards of Learning Assessments – Spring 2008
Released Tests* © 2008 by the Commonwealth of Virginia Department of Education.
Reproduced by permission.

Directions: Read each question and choose the best answer. Then fill in the bubble on the answer sheet for the answer you have chosen. Please use pencil only. Do not write on the printed test.

1. Which is equivalent to the following expression?

$$(3x + 1)(4x - 1)$$

- A. $12x^2 + x - 1$ B. $12x^2 - x - 1$ C. $12x^2 - 1$ D. $12x^2 + 7x - 1$

2. Which is equivalent to the following expression?

$$(3x^2y^3z^{-2})^3$$

- A. $9x^5y^3z^6$ B. $\frac{27x^5y^3}{z^6}$ C. $\frac{9x^6y^6}{z^6}$ D. $\frac{27x^6y^6}{z^6}$

3. Andrea has 37 coins, all nickels and dimes. The value of the 37 coins is \$3.10. How many dimes does Andrea have?

- A. 12 B. 19 C. 25 D. 31

4. Jerri wrote these steps when solving an equation.

$$17(x + 3) = 6 - 4$$

Step 1: $17x + 51 = 6 - 4$

Step 2: $17x + 51 = 2$

Step 3: $17x = -49$

Step 4: $x = -\frac{49}{17}$

Which property justifies Step 1?

- A. Associative property for addition B. Commutative property for addition
C. Distributive property D. Additive identity property

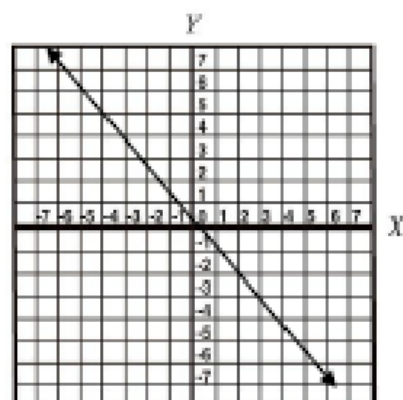
5. Which of the following represents the domain for the table of values listed here?

- A. 2, 4, 6, 8 B. 3, 4, 5, 6
C. 2, 8 D. 6

x	y
2	3
4	4
6	5
8	6

Go on →

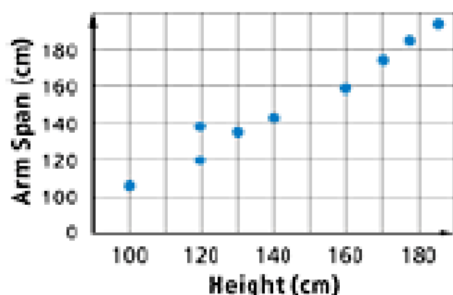
6. Which is closest to the slope of the line graphed below?



- A. $-\frac{1}{2}$ B. 1 C. -1 D. $\frac{1}{2}$
7. A dresser drawer contains one pair of socks of each of the following colors: blue, brown, red, white and black. Each pair is folded together in matching pairs. You reach into the sock drawer and choose a pair of socks without looking. The first pair you pull out is red – the wrong color. You put this pair back in the drawer then pull out a random pair. What is the probability that you will choose the red pair of socks twice?
- A. $\frac{1}{5}$ B. $\frac{1}{25}$ C. $\frac{2}{5}$ D. $\frac{1}{2}$
8. What is the slope of the line that passes through $(-3, -5)$ and $(4, -2)$?
- A. 1 B. $\frac{3}{7}$ C. $-\frac{3}{7}$ D. -1
9. Which set of real numbers is in order from least to greatest?
- A. π , -3 , -8 , 0.5 , $\frac{7}{8}$, $\sqrt{25}$
B. -8 , -3 , π , 0.5 , $\frac{7}{8}$, $\sqrt{25}$
C. -3 , -8 , 0.5 , $\frac{7}{8}$, π , $\sqrt{25}$
D. -8 , -3 , 0.5 , $\frac{7}{8}$, π , $\sqrt{25}$
10. Line l has slope 2 and goes through the point $(1, 3)$. Which is one form of the equation?
- A. $y = x + 2$ B. $y = 2x + 1$ C. $y = 3x + 2$ D. $y = 2x + 5$

Go on \rightarrow

11. The graph shows the height and arm span for a group of 10 people. Which of the following equations is closest to the best line of fit?



- A. $y = 3x$ B. $y = x$ C. $y = x - 1$ D. $y = -2x + 1$
12. Suppose you can work a total of no more than 20 hours per week at your two jobs. Babysitting pays \$5 per hour and your cashier job pays \$6 per hour. You need to earn at least \$90 per week to cover your expenses. Which system of linear inequalities best represents your job situation?

(Let x represent the number of hours you babysit and y represent the number of hours you work as a cashier. Assume x and y are both ≥ 0).

- A. $x + y < 90$
 $5x + 6y > 20$ B. $x + y \leq 90$
 $5x + 6y \geq 20$ C. $x + y \leq 20$
 $5x + 6y \geq 90$ D. $x + y < 20$
 $5x + 6y > 90$
13. Written in simplest form, $\sqrt{32}$ is equal to –
- A. $2\sqrt{4}$ B. $2\sqrt{16}$ C. $4\sqrt{2}$ D. $8\sqrt{2}$
14. If $x \neq 0$, which is equivalent to the following expression?

$$\frac{2x^4 - 6x^3 + 4x^2 + 10x}{2x}$$

- A. $x^3 - 3x^2 + 2x + 5$ B. $x^3 - 6x^3 + 4x^2 + 5x$
C. $2x^3 - 6x^2 + 4x^2 + 5$ D. $2x^4 - 6x^3 + 4x^2 + 5x$

Go on \longrightarrow

15. Which equation fits the pattern in the table?

A. $y = \frac{1}{2}x + 3$

B. $y = 2x - 1$

C. $y = x + 1$

D. $y = \frac{1}{2}x + 2$

x	y
2	3
4	4
6	5
8	6

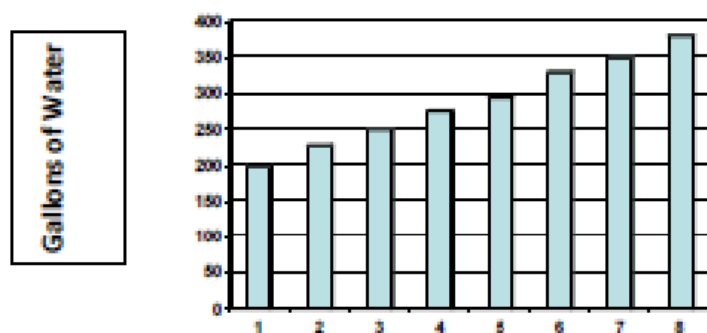
16. The bar graph below shows the number of gallons of water sold each month during an 8 month period. Based on the trend show on the graph, which is the most reasonable estimate of the number of gallons of water that will be sold during the twelfth month?

A. 380

B. 400

C. 450

D. 480



17. What is the Greatest Common Factor (GFC) for the following set of monomials?

xy^3z ,

$14y^2z^2$,

$7x^2y^3z^2$

A. $7xyz$

B. y^2z

C. xyz^2

D. $14z$

18. Solve the following equation. $|x - 2| = 5$

A. $-7, 3$

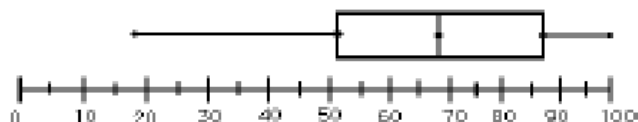
B. $7, -3$

C. $8, -1$

D. $-8, 1$

Go on →

19. Based on the box-and-whisker plot shown below, what is the interquartile range?



- A. 35 B. 18 C. 52 D. 83
20. If 112 children sign up for a field trip and each vehicle carries x children, which expression could be used to determine the number of vehicles needed for the trip?
- A. $112 - x$ B. $112x$ C. $\frac{112}{x}$ D. $\frac{x}{112}$
21. Each of the following tables contains elements of an (x, y) relationship. Which table contains four points that *cannot* lie on the graph of a function of x ?

A.

x	0	2	3	4
y	-1	-2	-3	-4

B.

x	1	2	3	2
y	4	2	2	4

C.

x	-1	-2	3	4
y	2	4	6	8

D.

x	0	1	5	6
y	5	9	2	-1

22. What is the Least Common Multiple (LCM) for the following set of monomials?

$$x^2, \quad 2x^3, \quad 3y^3$$

- A. $5x^6y^3$ B. $6x^6y^3$ C. $6x^3y^3$ D. $5x^3y^3$

23. Which is the factored form of the following equation?

$$18x^2 + 12x + 2$$

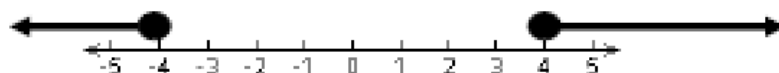
- A. $2(3x + 1)^2$ B. $(9x + 3)(2x + 4)$
C. $2(3x + 1)(3x - 1)$ D. $2(3x - 1)^2$

Go on →

24. The number of water bottles used during a team's football practice varies directly with the temperature (Fahrenheit). If a team uses 75 bottles when the temperature is 60°F , what is the temperature if they use 120 bottles?

A. 96°F B. 92°F C. 84°F D. 80°

25. Which of the following best describes the numbers shaded on this number line?

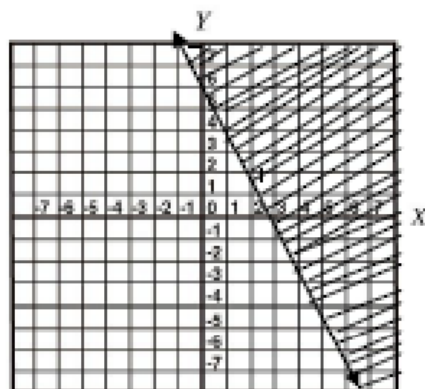


A. $-4 < x < 4$ B. $-4 \leq x \leq 4$ C. $x > -4$ and $x > 4$ D. $x \leq -4$ or $x \geq 4$

26. Find the y-intercept of the graph of the equation below.
 $4x - 5y = -35$

A. -7 B. 35 C. 7 D. $y = -35$

27. Which of the following best describes the shaded area on this graph?



A. $y \geq 2x + 5$ B. $y > 2x + 5$ C. $y > -2x + 5$ D. $y \geq -2x + 5$

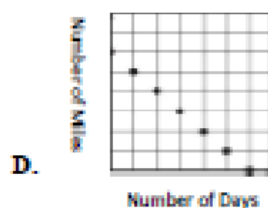
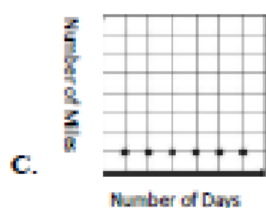
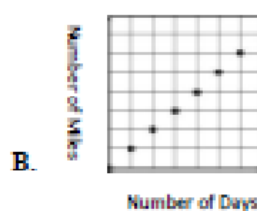
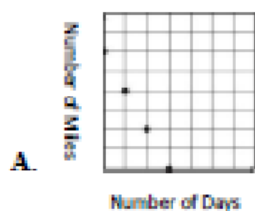
28. Find the following sum. Write the answer in standard form.

$$(2x^2 + x - 5) + (x + x^2 + 6)$$

A. $x + 3x^2 - 1$ B. $3x^3 + 2x^2 + 1$ C. $3x^2 + 2x + 1$ D. $3x^2 + x - 1$

Go on \longrightarrow

29. To train for a bicycle road race, Enrique needs to ride 150 miles per week at an average rate of 25 miles per day. The equation $M = 150 - 25d$ gives the number of miles, M , left to ride after d days. Which graph shows the number of miles Enrique has left to ride after d days?

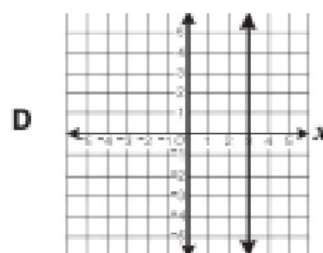
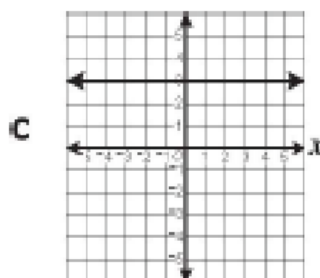
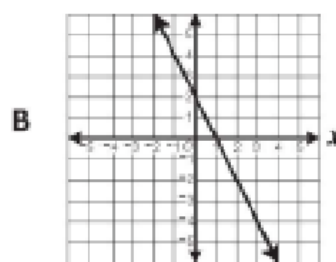
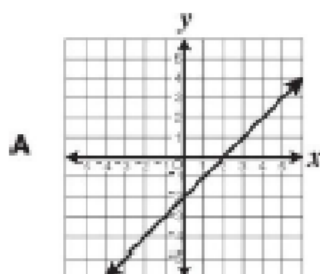


30. What is the solution set for the following quadratic equation?

$$x^2 - 4x + 4 = 0$$

- A. $\{2\}$ B. $\{-2\}$ C. $\{-2, 2\}$ D. $\{1, 3\}$

31. Which is most likely the graph of a line with a positive slope?



End of Exam.

Geometry

End of Course Assessment*

* Adapted from *Virginia Standards of Learning Assessments—Spring 2008 Released Tests* © 2008 by the Commonwealth of Virginia Department of Education. Reproduced by permission.

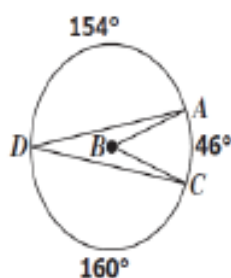
Directions: Read each question and choose the best answer. Then fill in the bubble on the answer sheet for the answer you have chosen. Please use pencil only. Do not write on the printed test.

- 1.) Calculate the surface area, in square inches, of a ball with a diameter of 4 inches.
Use 3.14 as an approximation of π .

A. 12.56 in.² B. 25.12 in.² C. 33.49 in.² D. 50.24 in.²

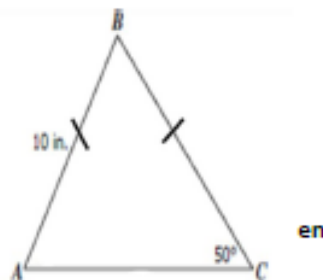
- 2.) Given: $\odot B$.
What is the $m\angle ADC$?

A. 23° B. 46°
C. 77° D. 80°



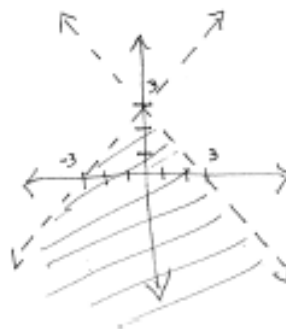
- 3.) Jennifer made these measurements on $\triangle ABC$. AC must be —

A. less than 10 inches B. equal to 10 inches
C. greater than 10 inches D. not enough information given



- 4.) Select the system of inequalities for the given graph.

A. $y > x + 3$, $y < x + 3$
B. $y < 3 - x$, $y < 3 + x$
C. $y > x - 3$, $y > x + 3$
D. $y < x - 3$, $y < x + 3$



Go on to next page.

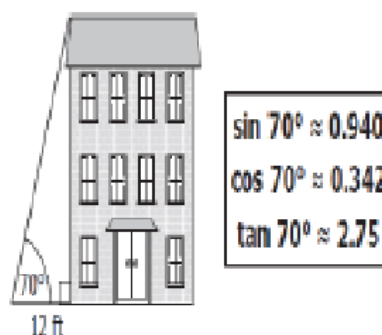
- 5.) If line Q has a slope of $-3/5$ and is parallel to line R , which statement must be true?

A. The slope of line R is $3/5$.
 B. The slope of line R is $-5/3$.
 C. The two lines have the same slope.
 D. The slope of line R is $5/3$.

- 6.) From a point 12 feet from the base of a building, the angle of elevation from the ground to the top of the building is 70° .

Which is closest to the height of the building?

A. 24 ft
 B. 33 ft
 C. 35 ft
 D. 41 ft



- 7.) The diameter of a circle has endpoints $(-3, 2)$ and $(3, -2)$. Which is closest to the length of the diameter of the circle?

A. 1.4
 B. 3.2
 C. 7.2
 D. 10.0

- 8.) A rectangular place mat is similar to the table upon which it is placed.

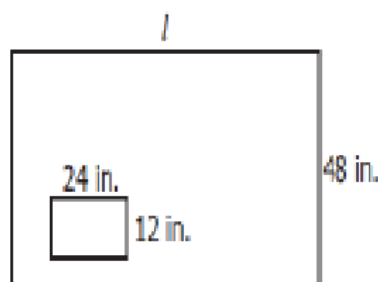
According to the diagram, which proportion can be used to determine the length of the table, l ?

A. $\frac{12}{48} = \frac{24}{l}$

B. $\frac{12}{24} = \frac{l}{48}$

C. $\frac{12}{l} = \frac{24}{48}$

D. $12l = 48$



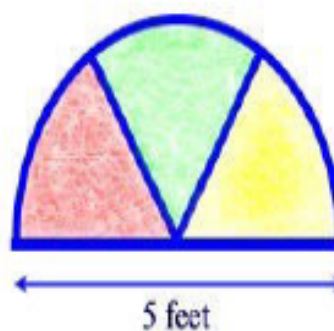
- 9.) Which pipe lengths could be joined to form a triangle?

A. 15 ft, 6 ft, 5 ft
 B. 13 ft, 12 ft, 5 ft
 C. 40 ft, 20 ft, 10 ft
 D. 19 ft, 16 ft, 2 ft

Go on to next page.

- 10.) A cathedral window is built in the shape of a semicircle. If the window is to contain three stained glass sections of equal size, what is the area of each stained glass section? Express answer to the nearest square foot.

- A. 1 sq. ft. B. 3 sq.ft.
C. 13 sq. ft. D. 26 sq. ft.

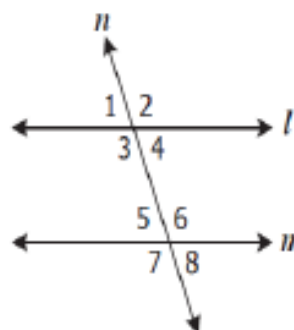


- 11.) A rectangle has a width of 6 ft. and an area of 54 square ft. The width and the length of the rectangle are then doubled. What is the new area of the rectangle?

- A. 108 sq. ft. B. 144 sq. ft. C. 216 sq. ft. D. 648 sq.ft.

- 12.) Lines l and m are cut by transversal n . Which statement would prove $l \parallel m$?

- A. $m \angle 2 = m \angle 6$ B. $m \angle 2 = m \angle 3$
C. $m \angle 7 + m \angle 8 = 180^\circ$ D. $m \angle 3 + m \angle 5 = 90^\circ$



- 13.) When assembled, the figure shown at the right will create a:

- A. Cube B. Cylinder
C. Cone D. Prism

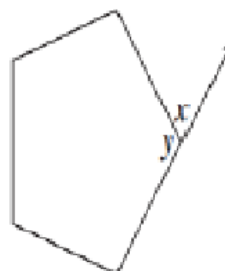


Go on to next page.

- 14.) This is a regular polygon. What are the values of x and y ?

A. $78^\circ, 102^\circ$ B. $72^\circ, 108^\circ$

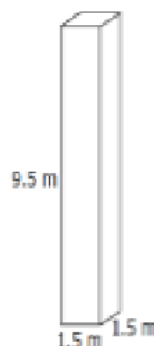
C. $60^\circ, 120^\circ$ D. $45^\circ, 135^\circ$



- 15.) A concrete pillar shaped as a rectangular prism is designed as follows. Which is closest to the volume of concrete needed to fill the pillar?

A. 12.5m^3 B. 14.3m^3

C. 21.4m^3 D. 28.5m^3



- 16.) The perimeter of a rectangle is 38m. The base is four more than two times the height. What is the height?

A. 5.67m

B. 5m

C. 11.33m

D. 8m

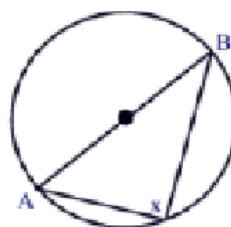
- 17.) Given the circle at the right with diameter \overline{AB} , find x .

A. 30°

B. 45°

C. 60°

D. 90°



- 18.) If the area of the base of a square pyramid is 49 square centimeters, and the lateral area is 210 square centimeters, what is the slant height?

A. 15 cm.

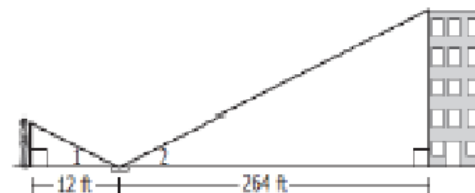
B. 7.5 cm.

C. 161 cm.

D. 12.86 cm.

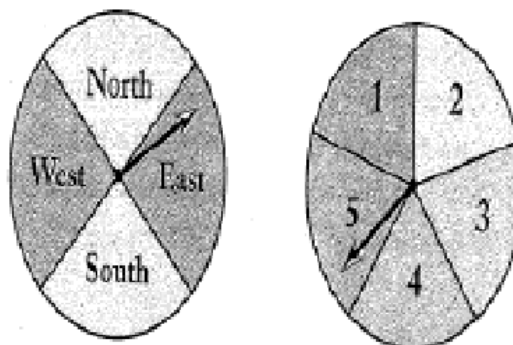
Go on to next page.

- 19.) Joseph is standing 12 feet from a mirror lying on the ground, and his eyes are 5 feet above the ground. The line-of-sight reflection on the mirror makes $\triangle 1$ congruent to $\triangle 2$.



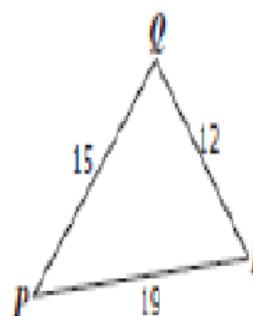
If the building is 264 feet from the mirror, which is closest to the height of the building?

- A. 100 ft B. 110 ft C. 130 ft D. 145 ft
- 20.) Kangmin and Roland are playing a game in which they spin two spinners. The first spinner tells what direction the player has to move his playing piece, and the second spinner tells how many spaces to move. What is the probability that Kangmin will move 5 spaces to the east?



- A. $1/5$ B. $1/9$ C. $2/9$ D. $1/20$
- 21.) Which lists the angles of the triangle in order from least to greatest?

- A. $R, \angle Q, \angle P$ B. $\angle Q, \angle P, \angle R$
- C. $\angle P, \angle R, \angle Q$ D. $\angle P, \angle Q, \angle R$



Go on to next page.

- 22.) Which rule of congruence is illustrated by the two triangles shown below?

A. Angle-Angle-Side

B. Side-Side-Angle



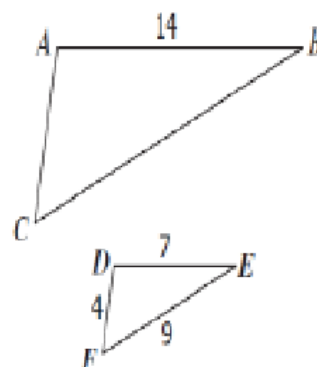
- 23.) In addition to the information given in the drawing, which statement would be sufficient to prove that $\triangle ABC \sim \triangle DEF$?

A. $\frac{BC}{AC} = \frac{2}{1}$

B. $\frac{BC}{AC} = \frac{9}{4}$

C. $AC = 18$ and $BC = 8$

D. $AC = 8$ and $BC = 18$



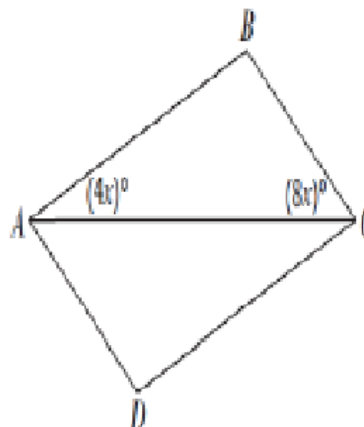
- 24.) If $ABCD$ is a parallelogram and $x = 5$, what is $m\angle D$?

A. 100°

B. 120°

C. 140°

D. 160°



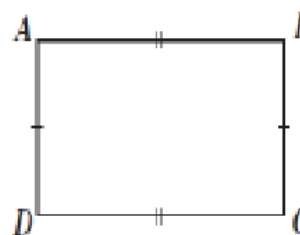
- 25.) The opposite sides of a window frame are congruent. Which additional piece of information would verify that the frame is a rectangle?

A. $\angle B \cong \angle D$

B. $\overline{AC} \perp \overline{BD}$

C. $\overline{AC} \cong \overline{BD}$

D. $m\angle A + m\angle D = 180^\circ$

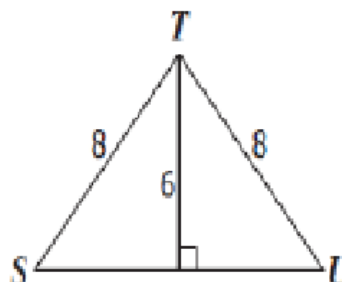


Window Frame

Go on to next page.

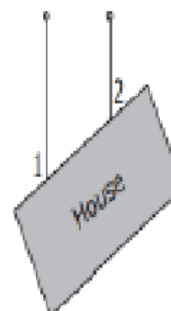
26.) What is the length of \overline{SU} ?

- A. $2\sqrt{7}$ cm
- B. 7 cm
- C. $4\sqrt{7}$ cm
- D. 20 cm



27.) Sally is using strings to mark parallel rows for a vegetable garden behind her house.
If the measure of $\angle 1$ is 115° , what should be the measure of $\angle 2$?

- A. 25°
- B. 65°
- C. 75°
- D. 115°



28.) Faith made two square pyramids out of clay. Each side of the base of the first pyramid measures 4 inches, and the height of the pyramid is 6 inches. Each dimension of the second pyramid is three times larger than the same dimension of the original pyramid.
Which of the following statements is true?

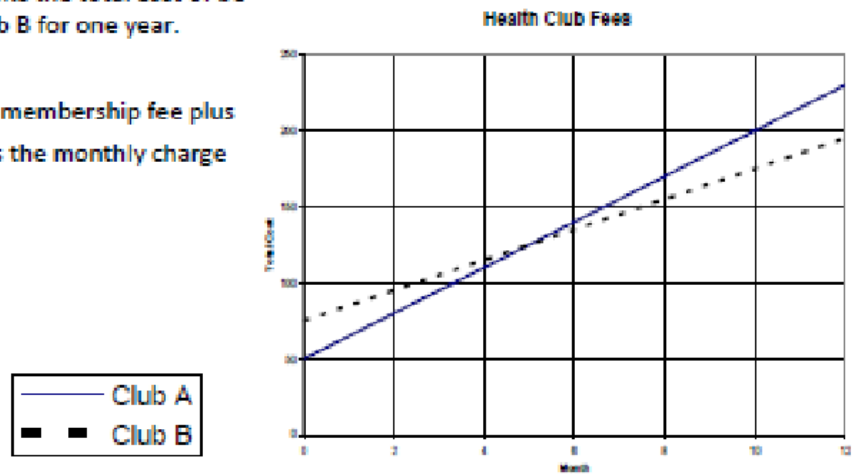
- A. The volume of the second pyramid is 3 times that of the first.
- B. The volume of the second pyramid is 9 times that of the first.
- C. The volume of the second pyramid is 18 times that of the first.
- D. The volume of the second pyramid is 27 times that of the first.

Go on to next page.

- 29.) Two health clubs offer different membership plans. The graph below represents the total cost of belonging to Club A and Club B for one year.

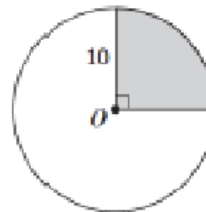
If a yearly cost includes a membership fee plus a monthly charge, what is the monthly charge for Club A?

- A. \$5 per month
B. \$10 per month
C. \$15 per month
D. \$20 per month



- 30.) When proving that a quadrilateral is a parallelogram by using slopes, you must find:
- A. The slopes of all four sides B. The slopes of two opposite sides.
- C. The lengths of all four sides. D. Both the lengths and slopes of all four sides.
- 31.) With a fixed perimeter of 60 inches, which regular shape will have the most area?
- A. a triangle B. a rectangle C. a pentagon D. an octagon
- 32.) The area of the *shaded* sector of circle *O* is —

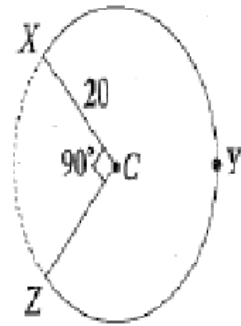
- A. 5π B. 20π C. 25π D. 50π



Go on to next page.

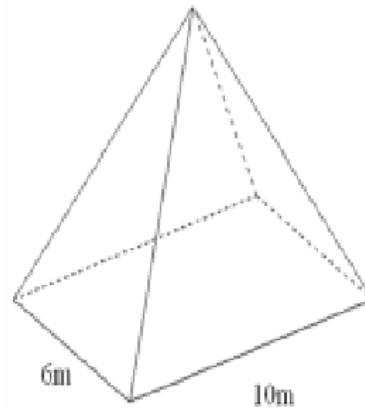
- 33.) If line segment CX of circle C is 20 units long, how long is arc XYZ ? Use 3.14 as an approximation of π . Round your answer to the nearest unit.

- A. 67 units B. 94 units
C. 2,000 units D. 44 units



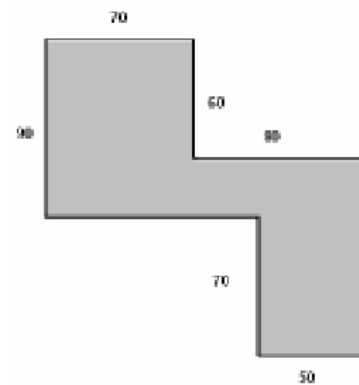
- 34.) A rectangular pyramid with the base dimensions as shown as a height of 12 meters and a slant height of 13 meters. The base diagonal is $2\sqrt{34}$ and the edge is $\sqrt{178}$ meters. What is the volume?

- A. 720 m^3 B. 240 m^3
C. 260 m^3 D. $20\sqrt{178} \text{ m}^3$



- 35.) Find the area of the figure to the right:

- A. 6,300 B. 24,000
C. 12,200 D. 11,800



End of exam.

Addendum 3: ATMI Subscales

ENJOYMENT (measures “intrinsic value”)

- 9. Mathematics is one of my most dreaded subjects. (R)
- 12. Mathematics makes me feel uncomfortable. (R)
- 13. I am always under a terrible strain in a math class. (R)
- 14. When I hear the word mathematics, I have a feeling of dislike. (R)
- 24. I have usually enjoyed studying mathematics in school.
- 26. I like to solve new problems in mathematics.
- 27. I would prefer to do an assignment in math than to write an essay.
- 29. I really like mathematics.
- 30. I am happier in a math class than in any other class.
- 31. Mathematics is a very interesting subject.

SELF CONFIDENCE (measures “expectancy”)

- 10. My mind goes blank and I am unable to think clearly when working with mathematics. (R)
- 11. Studying mathematics makes me feel nervous. (R)
- 15. It makes me nervous to even think about having to do a mathematics problem. (R)
- 16. Mathematics does not scare me at all.
- 17. I have a lot of self-confidence when it comes to mathematics.
- 18. I am able to solve mathematics problems without too much difficulty.
- 19. I expect to do fairly well in any math class I take.
- 20. I am always confused in my mathematics class. (R)
- 21. I feel a sense of insecurity when attempting mathematics. (R)
- 22. I learn mathematics easily.
- 23. I am confident that I could learn advanced mathematics.
- 28. I would like to avoid using mathematics in college. (R)
- 37. I am comfortable expressing my own ideas on how to look for solutions to a difficult problem in math.
- 38. I am comfortable answering questions in math class.
- 40. I believe I am good at solving math problems.

VALUE (measures “utility value”)

- 1. Mathematics is a very worthwhile and necessary subject.
- 4. Mathematics helps develop the mind and teaches a person to think.
- 5. Mathematics is important in everyday life.
- 6. Mathematics is one of the most important subjects for people to study.
- 7. High school math courses would be very helpful no matter what I decide to study.
- 8. I can think of many ways that I use math outside of school.
- 25. Mathematics is dull and boring. (R)
- 35. I think studying advanced mathematics is useful.
- 36. I believe studying math helps me with problem solving in other areas.
- 39. A strong math background could help me in my professional life.

MOTIVATION (measures “long term motivation”)

- 2. I want to develop my mathematical skills.
- 3. I get a great deal of satisfaction out of solving a mathematics problem.
- 32. I am willing to take more than the required amount of mathematics.
- 33. I plan to take as much mathematics as I can during my education.
- 34. The challenge of math appeals to me.

Addendum 4A: PARLO Fall 2010 Teacher interview protocol

Observed math class

1. What were your goals for today's math class? Going over slope ideas
2. How did you/do you assess whether or not you were meeting/met your goals? You can't- it is tough to figure out with a group like this
3. Were today's goals in any way informed by yesterday's class? If so, in what ways? – we had a short class yesterday [due to ice] so we had to finish yesterday's class today
4. Did you have different goals for different students? If so, describe that to me?
5. Did you adjust your instruction in any way today? If so, tell me about it.

PARLO check-in

6. With regards to PARLO, what is going well?
7. What is challenging and/or frustrating you about PARLO? [Probe to see if what kind of “challenge” it is – procedural, technical, philosophical, structural, etc., and whether it has been resolved.]

Students and PARLO

8. How are your students responding to PARLO?
9. What has challenged your students?
What has gone smoothly for them?
10. Does student engagement “look” different? By that, I mean are students interacting with math, with you, with each other differently under PARLO? Do you sense that there is more student ownership for both learning the material and turning “NYPs” into “Ps” or greater?
11. Think of one student in the class in the class I observed that you think has really benefitted from PARLO. Could you describe that child to me, and how you think PARLO has aided him/her in learning math? Conversely, is there a child in your class who has not benefitted from PARLO? Describe that child to me.
12. Do you know if students are logging onto EASE*?
13. Are your students using PARLO language? That is, do you hear them talking about being “NYP” or wanting to move from “P” to “HP”? If they are using PARLO language, when did you first notice it?

Instructional Practices

14. In what ways has your instruction changed as a result of PARLO? [Or has it changed?
15. In what way has the feedback you give to students changed as a result of PARLO?

Parents/Guardians and PARLO

16. Have you had any feedback from parents/guardians? If so, what are you hearing from parents/guardians and how have you responded?
17. Do you know if parents/guardians are logging onto EASE? *(This assumes that logging in is possible.)* NO

Professional Learning Communities

18. A big component of the study is the creation of PLCs. What are your expectations for these monthly meetings?
19. Do teachers in your math department routinely collaborate together? If so, what does that collaboration look like? Do you find yourself interacting more with your colleagues, about the same, or less under PARLO?
20. Do math teachers share a common planning time?

Summer Training

21. What was helpful?
 22. What was not helpful?
 23. What would you have liked to have learned that wasn't provided?
 24. Would you rearrange anything in the sequence of the training?
- [Other comments?]

* EASE was the online performance tracker used in fall, 2010. Later interviews asked about the replacement system, PARLO Tracker.

Addendum 4B: PARLO Exit Interview for Teachers

Questions A and B were administered only in Year 2.

- A. You've been in the PARLO study for two years now, and this year marks the end of your obligation. Why did you decide to continue (not to continue) with PARLO for a third year?
- B. (if school will be implementing in Year 3). What reservations if any did you have about signing up for a third year?

Reflections

- 1. Thinking back to the professional development you received, including the summer professional development, and the monthly PLCs, do you feel that this was sufficient or not to implement PARLO? And if not, what do you think was missing from it?
- 2. Can you describe for me the routines or practices you've implemented or adopted to support PARLO that you think you will continue to use in the future?
- 3. Do you use any of the PARLO techniques in your non-PARLO classrooms?

Students

- 4. What changes have you made to the feedback that you give students? And can you share an example? Are you writing more on tests?
- 5. How about challenges you had implementing PARLO. Can you tell me about some of them and how you have addressed them?
- 6. How about successes? Can you think of any student this year that particularly benefited from PARLO?

7. In what way has participating in PARLO changed your thinking about how 9th graders learn math?
8. What about student learning? Do you think being in a PARLO classroom has deepened your students' mathematical learning and understanding? If so, can you describe a couple of examples?
9. Do you think PARLO has made your students more responsible? If so, in what ways? Can you share some examples?
10. What aspects of PARLO have been the most frustrating or challenging for the students?

Parents/Guardians

11. Tell me about your parents or guardians, how have they responded this year?
12. In PARLO classrooms, parents and guardians are supposed to have more information on what their child knows and where they need to focus, do you think your parents have acted on this additional information?
13. How useful was the progress tracker during this project?

School Leadership

14. What kinds of support have you had from your principal and other school leaders in implementing PARLO?
15. So is there anything ideally that you think principals should be doing?
16. We gave a copy of the book Embedded Formative Assessment, has she mentioned it to you? Do you know if she read it?
17. In thinking about how to sustain this work, how to keep it going after the project ends, what do you think is needed at the classroom, school, and district levels?

Addendum 5: Open-Ended Student Survey

1. How would you describe PARLO to next year's 9th grade students?
2. What do you like about PARLO?
3. What do you not like about PARLO?
4. Would you like to have PARLO next year?
5. If yes, would you like it for just Math or All subjects?
6. About how often did you log on to Tracker?
7. If you logged on to Tracker, where did you log on?
8. Anything else you want to tell us about Tracker?

References

- Allison, P. D. (2012). Handling missing data by maximum likelihood. *SAS Global Forum, Paper 312-212*. Retrieved Online Sep 8, 2021, from <https://statisticalhorizons.com/wp-content/uploads/MissingDataByML.pdf>.
- Baker, J., & Boruch, R. (2015). *Ambient Positional Instability Among Ohio Math and Science Teachers: 2008 to 2014*. Retrieved from https://repository.upenn.edu/gse_pubs/269
- Boruch, R., Merlino, F. J., & Porter, A. (2012). *Golfing in a Hurricane: Education System Instability, Randomized Controlled Trials, and Children's Achievement*. University of Pennsylvania Center for Researcher and Evaluation in Social Policy. Retrieved from <http://cogscied.org/wp-content/uploads/2017/07/GolfingInAHurricane.pdf>
- Bun, Maurice J. G. & Harrison, Theresa D. (2018). OLS and IV estimation of regression models including endogenous interaction terms. *Econometric Reviews*, 38 (7) 814-27. <https://doi.org/10.1080/07474938.2018.1427486>
- Duckworth, A. L., Peterson, C., Matthews, M. D., & Kelly, D. R. (2007). Grit: Perseverance and passion for long-term goals. *Journal of Personality and Social Psychology*, 92(6), 1087–1101. <https://doi.org/10.1037/0022-3514.92.6.1087>

Koenker, Roger (1981). A Note on Studentizing a Test for Heteroscedasticity. *Journal of Econometrics*. 17: 107–112. [https://doi.org/10.1016/0304-4076\(81\)90062-2](https://doi.org/10.1016/0304-4076(81)90062-2)

Petreson, T. (2001). Endogeneity: Methodology. In M. J. Smelser and P. J. Baltes (eds.) *International Encyclopedia of the Social & Behavioral Sciences*, pp. 4511-4513.

Retrieved from <https://www.sciencedirect.com/topics/nursing-and-health-professions/endogenous-variable>

Raudenbush, S. W., & Bryk, A. S. (2002). *Hierarchical linear models: Applications and data analysis methods* (2nd ed.). Sage Publications.

U.S. Department of Education (2020). What Works Clearinghouse Standards Handbook, v4.1.

Washington, DC: Institute of Education Sciences, National Center for Education

Evaluation and Regional Assistance, What Works Clearinghouse. Downloaded 8/30/2021

from <https://ies.ed.gov/ncee/wwc/handbooks>