



**Measuring integrated understanding of undergraduate  
 chemistry research experiences: Assessing oral and written  
 research artifacts**

Journal:	<i>Chemistry Education Research and Practice</i>
Manuscript ID	RP-ART-04-2021-000104.R2
Article Type:	Paper
Date Submitted by the Author:	31-Oct-2021
Complete List of Authors:	<p>Helix, Max; University of California Berkeley, Graduate Group in Science and Mathematics Education            Cote, Laleh; University of California Berkeley, Graduate Group in Science and Mathematics Education; Lawrence Berkeley National Laboratory, Lawrence Berkeley National Laboratory            Stachl, Christiane; University of California Berkeley, Department of Chemistry            Linn, Marcia C.; University of California Berkeley, Graduate Group in Science and Mathematics Education            Stone, Elisa; University of California Berkeley, CalTeach Program            Baranger, Anne; University of California Berkeley, Department of Chemistry; University of California Berkeley, Graduate Group in Science and Mathematics Education</p>

1  
2  
3 **Measuring integrated understanding of undergraduate chemistry research experiences:**  
4  
5 **Assessing oral and written research artifacts**  
6  
7

8  
9  
10 **Authors:** Max R. Helix, Laleh E. Coté, Christiane N. Stachl, Marcia C. Linn, Elisa M. Stone,  
11  
12 and Anne M. Baranger  
13

14 **Abstract**  
15

16  
17 Understanding the impact of undergraduate research experiences (UREs) and course-  
18 based undergraduate research experiences (CUREs) is crucial as universities debate the value of  
19 allocating scarce resources to these activities. We report on the Berkeley Undergraduate  
20 Research Evaluation Tools (BURET), designed to assess the learning outcomes of UREs and  
21 CUREs in chemistry and other sciences. To validate the tools, we administered BURET to 70  
22 undergraduate students in the College of Chemistry and 19 students from other STEM fields,  
23 comparing the performance of students who had less than one year of undergraduate research to  
24 those with more than one year of research experience. Students wrote reflections and responded  
25 to interviews during poster presentations of their research project. BURET asks students to  
26 communicate the significance of their project, analyze their experimental design, interpret their  
27 data, and propose future research. Scoring rubrics reward students for integrating disciplinary  
28 evidence into their narratives. We found that the instruments yielded reliable scores, and the  
29 results clarified the impacts of undergraduate research, specifically characterizing the strengths  
30 and weaknesses of undergraduate researchers in chemistry at our institution. Students with at  
31 least a year of research experience were able to use disciplinary evidence more effectively than  
32 those with less than one year of experience. First-year students excelled at explaining the societal  
33 relevance of their work, but they incorporated only minimal discussion of prior research into  
34  
35  
36  
37  
38  
39  
40  
41  
42  
43  
44  
45  
46  
47  
48  
49  
50  
51  
52  
53  
54  
55  
56  
57  
58  
59  
60

1  
2  
3 their reflections and presentations. Students at all levels struggled to critique their own  
4  
5 experimental design. These results have important implications for undergraduate learning,  
6  
7 suggesting areas for faculty members, graduate student research mentors, and CURE or URE  
8  
9 programs to improve undergraduate research experiences.  
10  
11  
12  
13

14  
15 **Keywords:** undergraduate research, assessment, knowledge integration, URE, CURE,  
16  
17 postsecondary chemistry education, scientific poster presentations, instrument  
18  
19

## 20 21 **Introduction**

22  
23 Opportunities to conduct research are a critical component of undergraduate education  
24  
25 for many students majoring in science, technology, engineering, and mathematics (STEM)  
26  
27 disciplines, allowing them to engage with the larger scientific enterprise while still completing  
28  
29 relevant coursework. Although research experiences vary widely in nature, they generally share  
30  
31 common goals across settings, such as developing research skills, improving understanding and  
32  
33 application of scientific content knowledge, expanding scientific reasoning skills, increasing  
34  
35 confidence for doing science, and integrating students into scientific culture (Linn *et al.*, 2015;  
36  
37 Robnett *et al.*, 2015; Rodenbusch *et al.*, 2016; National Academies of Sciences and Medicine,  
38  
39 2017). Numerous studies have focused on the relationship between participation in a research  
40  
41 experience and the development of self-efficacy, confidence, and attitudes in/toward science,  
42  
43 which are linked to academic retention and career choice (e.g., Shuster *et al.*, 2019; Ashcroft *et*  
44  
45 *al.*, 2020; Avargil *et al.*, 2020; Esparza *et al.*, 2020). Research experiences can serve as a positive  
46  
47 influence for career aspirations involving science, despite other challenges students may have  
48  
49 faced since entering their college or university (Seymour and Hunter, 2019). As a result, many  
50  
51 institutions and funding organizations across the U.S. have dedicated considerable resources to  
52  
53  
54  
55  
56  
57  
58  
59  
60

1  
2  
3 support these programs each year (Laursen *et al.*, 2010; Auchincloss *et al.*, 2014; Krim *et al.*,  
4  
5 2019).

6  
7  
8 Despite a national call to replace traditional introductory laboratory courses with  
9  
10 research-based courses, this practice is still emerging at U.S. colleges and universities (Olson and  
11  
12 Riordan, 2012; Laursen, 2019). Those course-based undergraduate research experiences  
13  
14 (CUREs) that have been developed in chemistry allow students to develop self-confidence and  
15  
16 project ownership, as well as contributing to novel research in chemistry (Kerr and Yan, 2016;  
17  
18 Ghanem *et al.*, 2018; Cruz *et al.*, 2020). Additionally, students who complete these courses  
19  
20 believe that they have learned more chemistry content than they would have in traditional lecture  
21  
22 and laboratory courses (Chase *et al.*, 2017). These research-based courses support student  
23  
24 interest in chemistry, as students find them to be more enjoyable than “cookbook” laboratories  
25  
26 with predetermined project outcomes (Clark *et al.*, 2016; Mutambuki *et al.*, 2019; Muna, 2021).  
27  
28 Several studies have explored the benefits of group-based approaches to supporting  
29  
30 undergraduates as they learn about and conduct chemistry research (Danowitz *et al.*, 2016;  
31  
32 Hauwiller *et al.*, 2019).

33  
34  
35 Due to their prevalence and potential impact, it is important to assess the effects of  
36  
37 science research experiences on student learning, in order to determine how students progress  
38  
39 over time and to identify how research experiences can be improved to better serve participants  
40  
41 (Auchincloss *et al.*, 2014). Such assessments are relevant to both CUREs and undergraduate  
42  
43 research experiences that take place in research laboratories (UREs). Most previous studies that  
44  
45 assess learning outcomes of science research experiences are limited to a description of the  
46  
47 research experience or self-report data; fewer studies validate self-reports with analysis of  
48  
49 research products, direct measures of mastery of scientific content or practice, or observations of  
50  
51  
52  
53  
54  
55  
56  
57  
58  
59  
60

1  
2  
3 student activities (Linn *et al.*, 2015; National Academies of Sciences and Medicine, 2017; Krim  
4  
5 *et al.*, 2019; Lin *et al.*, 2019). Scholars such as Pagano *et al.*, 2018 and Stone *et al.*, 2020 have  
6  
7 commented on the limited number of studies dedicated to examining the impacts of CUREs in  
8  
9 chemistry, when compared to the life sciences. Thus, there is a need for additional assessment  
10  
11 tools that can be applied to undergraduate research experiences in chemistry, both inside and  
12  
13 outside of the classroom.  
14  
15

16  
17 Literature on undergraduate research and educational policy documents have identified  
18  
19 the following scientific practices as foundational to research experiences for undergraduates  
20  
21 (Laursen *et al.*, 2010; Sadler *et al.*, 2010): formulating research questions or hypotheses,  
22  
23 designing experiments, analyzing and interpreting data, making conclusions, iteratively planning  
24  
25 next steps, and explaining the significance of the research project. Collectively, these scientific  
26  
27 reasoning skills are widely regarded as a critical component of science education; educators have  
28  
29 moved away from the idea that such skills involve a single cognitive activity, and they are most  
30  
31 often viewed as a “set of different but coordinated skills” (Opitz *et al.*, 2017). Thus, the goal of  
32  
33 this study was to develop assessment tools to be used with STEM majors, with a particular focus  
34  
35 on chemistry, that measure the extent to which they understand research as a set of connected  
36  
37 practices. A specific aim for this work was to focus on assessing students’ understanding of  
38  
39 scientific practices in the context of their own research project, rather than investigating their  
40  
41 ability to answer questions about a hypothetical scenario. In this study, we address the following  
42  
43 research questions:  
44  
45  
46  
47

- 48  
49 1. Do the tools we developed distinguish between undergraduate students with different  
50  
51 levels of prior research experience?  
52  
53  
54  
55  
56  
57  
58  
59  
60

1  
2  
3 2. What do our tools tell us about what students understand about research and what they  
4  
5 are still learning at different stages of their undergraduate careers?  
6  
7

### 8 9 Theoretical Framework

10  
11 The theoretical basis for our work comes from Knowledge Integration (KI), a framework  
12 that has been used extensively in the design of learning environments and instruments to assess  
13 K-12 student knowledge of scientific content and practices (Linn, 1995; Linn and Eylon, 2011;  
14 Ryoo and Linn, 2012; Stone, 2014; Linn *et al.*, 2018). This learning science framework  
15 emphasizes that coherent understanding occurs when students make deep connections between  
16 their prior and new ideas. KI specifies four key components that support student learning (Linn  
17 and Eylon, 2011). The first is to elicit student ideas and prior understandings about a given topic.  
18 Students already have a repertoire of knowledge to draw on, and new knowledge will ultimately  
19 be built on these existing structures. Second, as students engage more with a particular concept,  
20 they discover new, scientifically normative ideas, some of which may challenge or contradict  
21 existing ideas. Third, as students explore the ideas they discover, they begin to distinguish  
22 between competing ideas and the contexts in which they are applicable. This process leads to a  
23 more nuanced understanding of the topic. Finally, students reflect upon their new knowledge in  
24 order to consolidate it into a coherent narrative.  
25  
26  
27  
28  
29  
30  
31  
32  
33  
34  
35  
36  
37  
38  
39  
40  
41  
42

43  
44 The process of conducting research generates knowledge in a way that parallels the KI  
45 framework (Linn *et al.*, 2015). Activities such as predicting and hypothesizing allow for *eliciting*  
46 undergraduate students' initial ideas. Undergraduate researchers then begin *discovering* new  
47 ideas over time as they gather data and participate in other research practices (Linn and Eylon,  
48 2011; White and Gunstone, 2014). They gradually learn to *distinguish* between possible  
49 interpretations for their data, and *reflecting* on their research enables learners to consolidate  
50  
51  
52  
53  
54  
55  
56  
57  
58  
59  
60

1  
2  
3 knowledge and generate new ideas for future work (Brown *et al.*, 1989; Linn and Eylon, 2011).  
4  
5 KI guides our expectation that as undergraduates progress in research, they will become more  
6  
7 proficient in understanding and discussing their research project, linking their insights to relevant  
8  
9 discipline-specific content knowledge to form coherent arguments.  
10  
11  
12

## 13 Literature Review

14  
15  
16 ***Impacts of Student Participation in Science Research Experiences.*** A number of studies  
17  
18 suggest that gains related to retention in STEM (e.g., graduation rates, entry into the STEM  
19  
20 workforce, graduate school attendance) are supported through participation in research  
21  
22 experiences, especially for students from groups historically underrepresented in STEM fields  
23  
24 (Schultz *et al.*, 2011; O'Donnell *et al.*, 2015; Estrada *et al.*, 2016; Carpi *et al.*, 2017). There is  
25  
26 increasing evidence to link participation in authentic scientific research with the development of  
27  
28 science identity through immersive learning of discipline-specific practices, referred to as  
29  
30 “legitimate peripheral participation” in situated learning theory (Lave and Wenger, 1991;  
31  
32 Robnett *et al.*, 2015). Factors such as a positive science identity, self-efficacy development,  
33  
34 access to mentoring, and engagement in research at the undergraduate level are important for  
35  
36 persistence in STEM and are critical for supporting students from groups historically  
37  
38 underrepresented in STEM fields (Carlone and Johnson, 2007; Chang *et al.*, 2011; Mondisa and  
39  
40 McComb, 2018; Ortiz *et al.*, 2020). In chemistry, participation in research experiences contribute  
41  
42 to retention, enthusiasm for chemistry-related careers, and appreciation for the process of  
43  
44 engaging in research in this discipline (e.g., Kerr and Yan, 2016; Williams and Reddish, 2018;  
45  
46 Muna, 2021).  
47  
48  
49  
50  
51

52  
53 ***Measures of Student Learning in Science Research Experiences.*** Various performance  
54  
55 assessments have been developed to directly measure multiple dimensions of student knowledge  
56  
57  
58  
59  
60

1  
2  
3 and skills gained from participation in science research experiences (Butz and Branchaw, 2020).  
4  
5 For example, the Danczak–Overton–Thompson Chemistry Critical Thinking Test (DOT) was  
6  
7 designed to measure critical thinking skills in chemistry students, regardless of prior student  
8  
9 knowledge in chemistry (Danczak *et al.*, 2020). The Biological Experimental Design Concept  
10  
11 Inventory (BEDCI) measures knowledge and diagnoses non-expert-like thinking in experimental  
12  
13 design by analyzing open-ended responses to different scenarios (Deane *et al.*, 2014). The  
14  
15 Assessment of Critical Thinking Ability (ACTA) is an open-ended survey that assesses critical  
16  
17 thinking skills in biology and chemistry students (White *et al.*, 2011). The Rubric for  
18  
19 Experimental Design (RED) identifies areas of experimental design in which undergraduates  
20  
21 struggle (Dasgupta *et al.*, 2014, 2016). The Performance assessment of Undergraduate Research  
22  
23 Experiences (PURE) instrument measures experimental problem solving and quantitative  
24  
25 literacy skills in chemistry students participating in UREs through a series of multipart questions  
26  
27 about real-world scientific problems (Harsh, 2016; Harsh *et al.*, 2017). The Test of Scientific  
28  
29 Literacy Skills (TOSLS) consists of multiple-choice questions about real-world problems and  
30  
31 measures student skills related to scientific literacy (Gormally *et al.*, 2012). Crawford and  
32  
33 Kloepper, 2019 developed an exit interview involving a series of written and oral exercises that  
34  
35 assess the ways in which chemistry students connect course content to laboratory activities.  
36  
37  
38  
39  
40  
41

42 Some instruments are “authentic assessments,” which are meaningful opportunities for  
43  
44 students to integrate and apply their knowledge to novel, complex, and/or realistic situations that  
45  
46 simulate typical activities of scientists (Wiggins, 1998; Doğan and Kaya, 2009; Laungani *et al.*,  
47  
48 2018). For example, the Experimental Design Ability Test (EDAT) gives students a real-world  
49  
50 scenario and research question and tasks them with designing an appropriate experiment, and has  
51  
52 been used in chemistry and the life sciences (Sirum and Humburg, 2011; Goodey and Talgar,  
53  
54  
55  
56  
57  
58  
59  
60



1  
2  
3 2016). Several studies suggest that writing activities can support student understanding of  
4 chemistry concepts (e.g., Lewis dot structure model) and methods (e.g., spectroscopy), as well as  
5 confidence in communicating about the material (Shultz and Gere, 2015; Moon *et al.*, 2018;  
6 Watts *et al.*, 2020). The Rubric for Science Writing and the Tool to assess Interrelated  
7 Experimental Design (TIED) are two assessment tools designed for use in undergraduate science  
8 courses, which involve students in activities that scientists engage in (Timmerman *et al.*, 2011;  
9 Killpack and Fulmer, 2018).

10  
11  
12  
13  
14  
15  
16  
17  
18  
19  
20  
21  
22  
23  
24  
25  
26  
27  
28  
29  
30  
31  
32  
33  
34  
35  
36  
37  
38  
39  
40  
41  
42  
43  
44  
45  
46  
47  
48  
49  
50  
51  
52  
53  
54  
55  
56  
57  
58  
59  
60  
There is compelling evidence to suggest that participation in a CURE leads to significant gains in research skills and academic outcomes and can support the subsequent advancement to (and success in) a URE (Rodenbusch *et al.*, 2016; Krim *et al.*, 2019). Studies that measure student learning gains typically consider these gains only over the course of a semester-long research experience, though there is evidence to suggest that undergraduates need to participate in high-impact research experiences spanning more than one semester to develop their understanding of the research process (Deane *et al.*, 2014; Corwin, Runyon, *et al.*, 2015; Griffeth *et al.*, 2015; Harsh, 2016; Remich *et al.*, 2016; Hernandez *et al.*, 2018). A longitudinal study by Szteinberg and Weaver (2013) suggests that CURE students retain chemistry content knowledge longer, as compared to students in traditional laboratory courses.

***Prior Gaps Identified in Learning for Undergraduate Researchers.*** Previous studies point to a lack of mastery among undergraduate researchers in fully understanding their research projects in several areas (Airey and Linder, 2009; Coil *et al.*, 2010; Gormally *et al.*, 2012). Prior findings suggest it is possible for a student to participate in research without understanding the scientific or societal significance of their work, though this skill supports more expert-level reasoning in the discipline of the research project (Bransford *et al.*, 2000; Coil *et al.*, 2010).

1  
2  
3 Many students, and in particular those from groups underrepresented in STEM fields, choose a  
4 STEM education/career in order to make a positive contribution to their communities and/or  
5 society, and this interest is likely to influence their commitment to a career in STEM (Bonous-  
6 Hammarth, 2000; Harackiewicz and Hulleman, 2010; Chang *et al.*, 2014). However,  
7  
8 undergraduates do not always develop the ability to articulate answers to questions about the  
9  
10 context of their research project, such as: “Why is this question important to others in this  
11  
12 discipline?” or, “What is the ‘big picture?’” (Timmerman *et al.*, 2011).  
13  
14  
15  
16  
17  
18

19 In order to become independent researchers, undergraduates are also expected to develop  
20 an understanding of experimental design (Sirum and Humburg, 2011; Killpack and Fulmer,  
21 2018). Undergraduates are typically presented with narratives about previously completed  
22 experiments as part of their STEM coursework, but training in designing experiments is less  
23 common (Gormally *et al.*, 2012). When reading scientific papers, undergraduates commonly  
24 struggle with evaluating and critiquing the design elements used in the studies being discussed  
25 (Varela *et al.*, 2005; Coil *et al.*, 2010). Guided-inquiry laboratories, CUREs, and UREs, in which  
26 students design their own experiments, can be used to support experimental design skills in  
27 chemistry (Goodey and Talgar, 2016). Multiple studies make the case that instruments are  
28 needed to measure experimental design and other skills critical for the development of students  
29 as scientists as they prepare to advance in their professional career (e.g., Sirum and Humburg,  
30 2011; Dasgupta *et al.*, 2014, 2016; Danczak *et al.*, 2020).  
31  
32  
33  
34  
35  
36  
37  
38  
39  
40  
41  
42  
43  
44  
45

46 Science research experiences often require that students contribute to data interpretation,  
47 but many undergraduates enter introductory-level STEM courses with insufficient skill in  
48 understanding how to work with data (e.g., reading graphs, analyzing and interpreting data,  
49 creating data visualizations), and STEM coursework does not necessarily cover this content (Coil  
50  
51  
52  
53  
54  
55  
56  
57  
58  
59  
60

1  
2  
3 *et al.*, 2010; Maltese *et al.*, 2015). Comparing the data analysis skills of various researchers  
4  
5 showed that novices are more heavily reliant on personal beliefs, while those with more expertise  
6  
7 focus on empirical consistency to draw conclusions from their observations (Hogan and  
8  
9 Maglienti, 2001). Chemical education studies suggest that students need to be taught explicitly  
10  
11 how to generate and interpret the kinds of visualizations they will need for a particular project,  
12  
13 and instruction should be intentional about connecting data to relevant concepts and addressing  
14  
15 misconceptions (Connor *et al.*, 2019; Rodriguez *et al.*, 2019). Relatively few studies in chemistry  
16  
17 have focused on assessing student skill level in this area, though there is consensus that data  
18  
19 interpretation is critical for developing chemists (Maltese *et al.*, 2015; Peteroy-Kelly *et al.*,  
20  
21 2017).  
22  
23  
24  
25

26 Undergraduate students should also be able to develop hypotheses and conduct  
27  
28 appropriate experiments to test these hypotheses by the time they graduate with a STEM  
29  
30 bachelor's degree (White *et al.*, 2013). When students are provided with the space to encounter  
31  
32 challenges, revise their research goals, and repeat their work, this iterative process can have a  
33  
34 powerful impact on their sense of ownership as they learn to navigate obstacles in their scientific  
35  
36 discipline (Corwin *et al.*, 2018; Gin *et al.*, 2018). CUREs focused on chemistry have been shown  
37  
38 to improve students' project ownership in lower-division, upper-division, and large-enrollment  
39  
40 undergraduate courses (Williams and Reddish, 2018; Cruz *et al.*, 2020; Heller *et al.*, 2020).  
41  
42  
43  
44

#### 45 The BURET Study

46  
47  
48 We have drawn from this literature to develop four **Berkeley Undergraduate Research**  
49  
50 **Evaluation Tools (BURET) Indicators** that describe areas where undergraduates are expected  
51  
52 to integrate their understanding of foundational scientific practices:  
53  
54  
55  
56  
57  
58  
59  
60

- I. Communicate the significance of their specific project to the overarching research questions of the laboratory and the broader scientific field
- II. Justify their experimental design as appropriate for their research question
- III. Analyze and interpret data in order to construct explanations and models that are relevant to their research question
- IV. Generate hypotheses and plan future experiments relevant to their research question in response to their analysis and interpretation of data

These Indicators provide the focus for the new instruments described in this study. We designed an interview protocol and reflective prompts to assess how undergraduates develop an integrated knowledge of these dimensions as they engage in research. The first is the Reflection instrument (BURET-R), which prompts written student reflections about the progress of their research project. The second is the Poster interview (BURET-P), which is administered at capstone poster sessions. Both tools were administered to students in a variety of research settings.

## Methods

### Participants and Context

Participants were recruited from CUREs and UREs at our institution. Undergraduate researcher volunteers came from five different populations (Table 1), with 78% majoring in chemistry disciplines (chemistry, chemical engineering, or chemical biology). It should be noted that nearly all students in UREs had previously taken a CURE, as is typical for URE students in many science departments at this institution. Our procedures were approved by the University of California, Berkeley Committee for Protection of Human Subjects, Protocol #2016-02-8360.

**Table 1.** Study populations.

Group	n	Response Rate	Type	Duration	Prior Research Experience	Student Description
1	35	58%	CURE	Semester	Mostly none	Freshman chemistry students
2	6	90%	CURE	Summer	None	New chemistry transfer students
3	28	59%	URE	Ongoing	Variable	Department of chemistry students
4	5	65%	CURE	Semester	None	Pre-service STEM teachers
5	15	88%	URE	Summer+	Variable	Pre-service STEM teachers

Students in Groups 1 and 4 (Table 1) were enrolled in CUREs in which the student was responsible for developing their own research question and choosing the methods used to investigate that question. Students in Group 2 chose from possible research projects that could be investigated using computational chemistry approaches. Students in Groups 3 and 5 had typical apprentice-style research experiences in faculty labs, where the projects varied but fit into the overarching goals of their faculty advisor and were generally related to the projects of their graduate student mentor. The level of independence in designing their own work varied as well, generally according to the amount of time each undergraduate had spent working in their research group.

The students participating in the study ranged from having zero to four or more semesters of research experience prior to study participation. The study population of 89 undergraduates contained a mixture of identities, including gender, race, ethnicity, and first language. Our study participants were 58% female, and 24% stated that English is not their first language. Students who self-identified as American Indian/Alaska Native, Black/African-American, Hispanic/Latinx, or other Pacific Islander, collectively referred to as underrepresented minorities (URM), were intentionally oversampled and comprised 19% of our study population.

## Instrument Development

### *Expert Review of the BURET Indicators*

To confirm that the BURET Indicators were aligned with the goals of faculty advisors, all chemistry faculty at this institution working with undergraduate researchers were invited to participate in an interview. A total of 21 faculty agreed to be interviewed, for a response rate of 41%. The faculty in this study ranged from assistant professors to full professors and had a wide variety of research group sizes. During a 1-hour interview, faculty were asked to describe their goals for their undergraduate researchers, discuss mentoring practices, and review the BURET Indicators. They commented on whether these were appropriate goals for their undergraduates. Nearly all of the responses were positive, with some faculty members expressing that the Indicators “exactly” described their overall goals for undergraduate researchers. An additional 12 faculty members from other STEM departments were also interviewed, and they gave largely similar responses.

### *Assessment Design*

We sought data collection and assessment approaches that would both support student learning and allow for direct measures of the Indicators across both CUREs and UREs. Many undergraduate researchers create a poster and present their research project as a capstone requirement, providing an opportunity to assess student integration of scientific content and practices in the context of their own work. A set of interview questions targeting the BURET Indicators were developed to ask at the end of each student’s prepared presentation. This interview protocol coupled with a rubric to assess several aspects of these verbal presentations make up the BURET Poster Presentation instrument (BURET-P, Fig. 1).

1. Can you please summarize why your research project and what you've learned is important?
2. Can you explain more about why you (and your lab) chose this general strategy for your research project?
3. Can you choose one experimental technique that is central to this work and say why you used it, rather than other options?
  - 3b. What are the limitations of this technique?
4. Could you expand on how you interpret these results?
  - 4b. How confident are you in your data and your conclusions?
5. What would you do if you had another year to work on this project, and why?

**Fig. 1.** BURET-P interview protocol.

A pair of reflective prompts (BURET-R, Fig. 2) were developed to complement the poster presentation assessment. These prompts can be administered at different points in the research experience to provide information on students' developing progress on the Indicators. In this study, BURET-R was administered a few weeks prior to their poster session. These prompts targeted Indicators 3 and 4, respectively, but many students also incorporate discussions of Indicators 1 and 2 in their responses.

1. *Data Analysis Prompt:* Think about the ways you have analyzed data recently.
  - (a) Describe one example of data analysis you have done.
  - (b) Reflect on how you used this data analysis to create or change an explanation or a model. Frame your response for an experienced scientist who is unfamiliar with your project.
2. *Next Steps Prompt:*
  - (a) If you had another month or two to work, what would be your next steps and why?
  - (b) What about if you had another year?

**Fig. 2.** BURET-R reflective prompts.

### *Items and Scoring Rubrics*

Preliminary rubric development was conducted with a small group of 7 undergraduate and 5 graduate students. All participants responded to BURET-R prompts, and a few also presented posters to the research team. Four undergraduates were also interviewed, during which they were asked to expand on their BURET-R and BURET-P responses. Written responses and audio recordings were reviewed to develop the rubrics. The emergent themes from initial rounds

of coding and a review of the relevant literature were used to develop an overlapping set of specific items aligned with the BURET Indicators, resulting in a set of 6 items for the BURET-R and 11 items for the BURET-P scoring instruments (Table 2). Although not all of these items were explicitly elicited by our prompting questions, they were all commonly discussed in student answers. Our assumption was not necessarily that every response would address every item, but that on average, more expert-like responses would integrate more different types of content into the overall narrative, as suggested by the Knowledge Integration theoretical framework. It should be emphasized that the rubrics are not a direct measure of everything a student knows, but rather of what they choose to present.

**Table 2.** Scoring rubric items for the BURET instruments align to BURET Indicators.

Item	BURET-R	BURET-P
Placing their work in a broader context	X	
Placing their work in a broader scientific context		X
Placing their work in a broader societal context		X
Providing rationale for an experimental design choice	X	X
Addressing limitations of an experimental design choice		X
Comparing alternatives to an experimental design choice		X
Number of experimental design choices with some rationale		X
Identifying and discussing the key variables	X	
Describing their data analysis procedures OR Interpreting their data	X	
Interpreting their data		X
Analyzing sources of error and uncertainty		X
Proposing next steps for the project	X	X
Incorporating references to previous work		X
Integrating additional content knowledge	X	X

Although many of the research projects assessed during this study involved experiments, this was not universally the case. To account for other types of research, “experimental design choice” was defined as any approaches, strategies, techniques, or other decisions made during the study design process. This definition was sufficiently broad to encompass the work described by all students who participated in this study.



1  
2  
3 To develop a KI rubric for each item, we were guided by prior KI rubrics for similar  
4 items. The KI rubrics score items on a 0-5 scale. Each score represents a progressively more  
5 integrated and connected response. Applying the KI framework to BURET scoring, descriptions  
6 were written for each possible score on each item (see Appendix A for the complete rubric).  
7  
8 These were anchored by the idea that a 2 should be a correct statement about an isolated part of  
9 the research process, and a 4 should be a clear, basic link from a relevant part of the research  
10 process to evidence from scientific content and research practices. For example, a score of 2 on  
11 “addressing the limitations of an experimental design choice” could be obtained by simply  
12 noting a drawback for a particular technique, whereas a score of 4 would require that students  
13 explain that limitation by integrating underlying scientific principles into their discussion or by  
14 making a clear reference to the research question. The remaining levels were defined as follows:  
15  
16 0 indicates that responses relevant to the item are absent, 1 indicates a vague statement, 3  
17 indicates a partial link between an assertion and relevant scientific content or practices, and 5  
18 indicates a complex link of 3 or more isolated concepts. The highest level descriptions were  
19 informed in part by the graduate students who responded to the BURET-R and BURET-P  
20 assessments as scoring categories were being refined. A partial scoring rubric with example  
21 responses can be found in Appendix B.  
22  
23  
24  
25  
26  
27  
28  
29  
30  
31  
32  
33  
34  
35  
36  
37  
38  
39  
40  
41  
42

### 43 Instrument Testing

#### 44 *Data Collection*

45  
46 To determine whether the BURET instruments could detect a difference between novice  
47 and advanced researchers, students enrolled in the three target CUREs (Groups 1, 2, and 4 in  
48 Table 1) were invited to be part of this study. A few weeks before their corresponding poster  
49 session, student responses to BURET-R were collected in class from all who agreed to  
50  
51  
52  
53  
54  
55  
56  
57  
58  
59  
60

1  
2  
3 participate (see Table 1 for response rates). At the final poster session for those courses, a sample  
4  
5 of the consenting students was interviewed using the BURET-P protocol (see Appendix C for  
6  
7 sampling procedures). Students were interviewed by one of 10 different interviewers by first  
8  
9 allowing the student to give their prepared presentation uninterrupted, then using the  
10  
11 standardized protocol in Fig. 1 to elicit specific elaborations. Poster interviews were recorded  
12  
13 and transcribed before coding.  
14  
15

16  
17 Additionally, students presenting at one of the two target URE poster sessions (Groups 3  
18  
19 and 5 in Table 1) were invited by email to participate in this study. Responses to BURET-R were  
20  
21 collected via Qualtrics a few weeks prior to the poster sessions, and all consenting students who  
22  
23 provided responses to BURET-R were interviewed at their poster session, using the same  
24  
25 protocol that was used with the CURE students. From our full dataset, 80 BURET-R responses  
26  
27 and 55 BURET-P interviews were found to be complete and fully legible or audible, and these  
28  
29 were used in our subsequent analysis.  
30  
31  
32  
33  
34

### 35 *Coding and Rubric Reliability*

36  
37 Data from BURET-R and BURET-P were scored for each study participant according to  
38  
39 the corresponding rubrics. 60% of the written responses to BURET-R were coded by two  
40  
41 different researchers, and discrepancies were resolved through subsequent discussion. A  
42  
43 weighted Cohen's kappa of 0.73 was deemed acceptable, and subsequent coding was completed  
44  
45 individually. Each poster transcript was deidentified with respect to student experience level and  
46  
47 other characteristics. Transcripts were then coded independently by at least two people, and any  
48  
49 discrepancies between coders were discussed and resolved. Two different pairs of coders  
50  
51 assessed each transcript. The coders varied in their level of chemistry expertise, and we found  
52  
53  
54  
55  
56  
57  
58  
59  
60

1  
2  
3 that our scoring rubrics could be used by any coders with a basic understanding of different  
4 research topics and expertise in scientific research practices. Rather than assessing “correctness,”  
5 these rubrics focused on our primary goal of assessing integration of content and practices with  
6 respect to various parts of each research project. A weighted Cohen’s kappa of 0.65 was  
7 achieved between coding pairs, using posters that were coded by all researchers. This is  
8 considered to be a substantial level of agreement according to Landis and Koch (1977).  
9  
10  
11  
12  
13  
14  
15  
16  
17

### 18 *Quantitative Analysis*

19  
20

21 We established the experimental validity of our instruments based on their ability to  
22 successfully distinguish between responses from undergraduates with more or less prior research  
23 experience. Participants were divided into novice (0-1 semester) and advanced (2+ semesters)  
24 groups based on how many semesters of research they had completed prior to the one in which  
25 they were presenting a poster. For each instrument, all items were averaged to produce a single  
26 test statistic. We compared novice and advanced participants using a t-test. Additionally, student  
27 scores on each item were collapsed into either low (KI score of 0-3) or high scores (KI score of  
28 4-5), and chi-squared tests were performed to determine whether high scores were significantly  
29 associated with increased research experience for individual items. Further psychometric  
30 analysis was performed to establish the internal structure at the instrument level and to  
31 investigate the dimensionality of our construct (see Appendix D for details). As a measure of  
32 internal consistency, Cronbach’s alpha was calculated for each instrument. All statistical analysis  
33 was conducted on Stata.  
34  
35  
36  
37  
38  
39  
40  
41  
42  
43  
44  
45  
46  
47  
48  
49  
50  
51

52 Item-response theory (IRT) analysis was conducted to gather validity evidence based on  
53 internal structure at the instrument level. Because our sample was not sufficient to run the  
54  
55  
56  
57  
58  
59  
60

analysis using all thresholds from our rubric, data were collapsed into scores of low (0-2), moderate (3), or high (4-5), and Wright maps for each instrument were generated from the collapsed data. Additionally, exploratory factor analysis was performed and item-test correlations were calculated to determine whether the construct we are measuring is uni- or multi-dimensional. All statistical analysis was conducted on Stata except for the IRT analysis, which was performed on Conquest.

## Results

### Do the BURET Instruments Distinguish Between Undergraduate Students with Different Levels of Prior Research Experience?

An analysis of student responses showed that both the BURET-R ( $n = 80$ ) and BURET-P ( $n = 55$ ) instruments are able to distinguish between more and less experienced undergraduate researchers. Total scores for each instrument revealed statistically significant differences between students with 2 or more semesters of prior research experience and students with less experience ( $p < 0.001$ ; Tables 3 and 4). Average scores on each item also increased with more research experience, with 9 of the 17 items showing statistically significant gains.

**Table 3.** Mean scores on BURET-R rubric items.

Semesters of Previous Research Experience	0-1	2+	Sig
Sample size (n)	42	38	
Placing work in a broader context	2.9	3.6	*
Providing rationale for expt. design choice	1.9	2.8	**
Identifying and discussing the key variables	2.0	2.4	
Describing OR interpreting data analysis	2.6	3.5	**
Proposing next steps for the project	2.6	3.2	
Integrating additional content knowledge	0.8	2.4	**
Average Score	2.1	3.0	***

\*\*\*  $p < 0.001$ ; \*\*  $p < 0.01$ ; \*  $p < 0.05$

**Table 4.** Mean scores on BURET-P rubric items.

Semesters of Previous Research Experience	0-1	2+	Sig
Sample size (n)	24	31	
Placing work in broader scientific context	2.3	3.5	*
Placing work in broader societal context	3.6	3.7	
Providing rationale for expt. design choice	3.5	3.9	
Addressing limitations of expt. design choice	2.8	3.3	
Comparing alternatives to expt. design choice	2.7	3.4	
Expt. design choices with some rationale (max. 5)	2.5	3.3	*
Interpreting their data	3.1	3.5	*
Analyzing sources of error and uncertainty	2.3	2.5	
Proposing next steps for the project	3.1	3.2	
Incorporating references to previous work	1.9	2.9	*
Integrating additional content knowledge	2.3	3.5	*
Average Score	2.7	3.3	***
*** $p < 0.001$ ; ** $p < 0.01$ ; * $p < 0.05$			

As a measure of reliability, Cronbach's alpha is calculated to be 0.78 for both instruments, which is in the range considered acceptable for science education research instruments (Taber, 2018). Further psychometric analysis suggests an acceptable consistency of the items to measure respondent performance and provides evidence that a unidimensional construct is being measured (see Appendix D for more information).

Two other variables that are highly correlated with increased research experience are year in school and whether the research experience was part of a course. As previously mentioned, most of the novice researchers in our sample were enrolled in a CURE, while most of the advanced researchers were participating in a URE in a faculty lab and had previously completed a CURE. To determine which of these variables was the best predictor of total score on our instruments, factorial ANOVAs were run using year in school, semesters of research experience, and URE/CURE as the independent variables. For the BURET-R, only URE/CURE was a significant predictor ( $p < 0.05$ ), whereas the duration of time spent in college or in undergraduate

1  
2  
3 research were not. For the BURET-P, only semesters of research experience was a significant  
4 predictor ( $p < 0.05$ ). On this instrument, the type of research experience did not have a  
5 significant effect on student performance, and the URE students with minimal total research  
6 experience performed similarly to other novice researchers. No interaction terms were significant  
7 for either instrument. We were unable to examine whether there were differential effects for  
8 students who identified as a URM because we recruited too few of these students with at least 2  
9 semesters of research experience. Future work will be needed to investigate this aspect of our  
10 instrument.  
11  
12  
13  
14  
15  
16  
17  
18  
19  
20  
21  
22

### 23 Comparison of the BURET-R and BURET-P

24  
25 Overall, largely similar final results were obtained from the BURET-R and the BURET-P  
26 instruments, which were generally administered several weeks apart. Student scores on the  
27 BURET-R and BURET-P instruments were significantly correlated with one another ( $r = 0.4$ ,  $p$   
28  $< 0.01$ , see Appendix E for a scatterplot). The items on which students tended to excel or  
29 struggle were similar across the two instruments, with some variations based on the exact  
30 relationships between the items assessed and the specific prompt or interview questions being  
31 answered. Targeted questions asked during the poster presentations generally elicit more specific  
32 information than the broader reflective prompts, resulting in more items being coded when  
33 assessing poster presentations. Poster presentations were also much longer than the written  
34 responses to BURET-R; on average, written responses were 248 words in length, while poster  
35 presentations (including answers to questions) were 1,682 words in length. In general, students  
36 scored higher on BURET-P (average score = 3.1) than on BURET-R (average score = 2.3). This  
37 can also be seen by looking at individual participants; 85% of the participants scored higher on  
38 their poster presentations than on their written responses.  
39  
40  
41  
42  
43  
44  
45  
46  
47  
48  
49  
50  
51  
52  
53  
54  
55  
56  
57  
58  
59  
60

1  
2  
3 What do the BURET instruments tell us about what undergraduate chemistry students  
4 understand about research and what they are still learning at different stages of their  
5 undergraduate careers?  
6

7 The BURET instruments provide information about the progress students make on each  
8 of the BURET Indicators as they gain in research experience. The following sections describe  
9 the characteristics of student progression along each Indicator using the KI framework, including  
10 undergraduate student performance on each item and the items that most differentiate novice and  
11 advanced study participants. Items are grouped by which Indicator they are most closely  
12 associated with to provide a more holistic picture of each primary area of assessment and for  
13 clarity of interpretation. (Table 5).  
14  
15  
16  
17  
18  
19  
20  
21  
22

23 **Table 5.** BURET-P items grouped by most related BURET Indicator.<sup>a</sup>  
24

Indicator	Items
I	Placing their work in a broader scientific context Placing their work in a broader societal context Incorporating references to previous work Integrating additional content knowledge
II	Providing rationale for an experimental design choice Addressing limitations of an experimental design choice Comparing alternatives to an experimental design choice Number of experimental design choices with some rationale
III	Interpreting their data Analyzing sources of error and uncertainty
IV	Proposing next steps for the project

25  
26  
27  
28  
29  
30  
31  
32  
33  
34  
35  
36  
37  
38 <sup>a</sup> The items are grouped by Indicator for clarity of interpretation. Note that there is no evidence from the internal  
39 structure of the instruments for the items to be grouped in this way  
40  
41  
42  
43  
44  
45

#### 46 Indicator 1: Communicating Significance

47  
48 The first BURET Indicator assesses how well students can communicate the significance  
49 of their specific project to the overarching research questions of the laboratory and the broader  
50 scientific field. Three of the four items corresponding to this Indicator for the BURET-P  
51 instrument showed statistically significant growth between novice and advanced students.  
52  
53  
54  
55  
56  
57  
58  
59  
60

1  
2  
3 Advanced students demonstrated a more sophisticated understanding of their project's scientific  
4 context ( $p < 0.05$ ), referred to previous work more often ( $p < 0.05$ ), and integrated more content  
5  
6 knowledge into their presentations ( $p < 0.05$ ) when compared to less experienced students (Table  
7  
8 4). Analysis of students' responses to the BURET-R instrument also provided evidence that they  
9  
10 develop in their ability to place their project into a broader context (Table 3).  
11  
12  
13

14  
15 Advanced students often demonstrated a more integrated understanding of scientific  
16  
17 context by explaining the current state of the field or how their research might affect projects in  
18  
19 other labs. One chemistry student who received a score of 4 stated, "I was ... working on  
20  
21 investigating ... the mechanical properties of polycarbonate urethane. Our research is particularly  
22  
23 relevant to joint implants and joint replacements, ... the current industry standard polymer is  
24  
25 called ultra high molecular wave polyethylene. ... Polycarbonate urethane or PCU is being  
26  
27 pioneered as a new material. ... But it's pretty new so we're still doing research on the very  
28  
29 mechanical properties and how it will react to being in the body and in an ionic environment  
30  
31 where there is salts and stuff like that, that can affect its microstructure." In this response, the  
32  
33 student clearly connects their work on the mechanical properties of PCU to the broader field of  
34  
35 material science, particularly in the area of artificial joints.  
36  
37  
38  
39

40 A student would receive a 2 on the "Incorporating references to previous work" item by  
41  
42 clearly referring to previous research but failing to explicitly link that research to their  
43  
44 experimental design or compare it to their own results. For example, a student was scored 2 for  
45  
46 the following vague reference to previous work, 'A lot of it was help from literature that we've  
47  
48 seen online, especially the solvents. I wouldn't have known where to start without using some of  
49  
50 these.' As an example of a higher scoring discussion, one student stated that, "There'd been, not  
51  
52 a consensus, but almost every single study that we had read previously looking for these heavy  
53  
54  
55  
56  
57  
58  
59  
60



1  
2  
3 metals in chocolate, but also in other candy, had focused on the cocoa, then being the source and  
4 maybe mentioned other possible sources in passing.” The student then compared this body of  
5  
6 previous work with their own work, which found a possible alternate source of heavy metals,  
7  
8 resulting in a score of 4.  
9  
10

11  
12 Additionally, advanced students scored higher on providing context by integrating more  
13 additional content knowledge into their presentation and answers. Additional content knowledge  
14 was defined as “exhibiting scientific content knowledge beyond what is required to describe the  
15 project.” Students received a 2 by simply providing some additional clarification, or a 4 by  
16 providing multiple examples or extensive discussions of relevant information. It should be noted  
17 that this does not directly measure the content knowledge of a student, but rather the extent to  
18 which students have *integrated* that content knowledge into discussions of their research.  
19  
20  
21  
22  
23  
24  
25  
26  
27  
28

### 29 Indicator 2: Justifying the Experimental Design 30

31  
32 The second BURET Indicator was assessed with three items that focused on how students  
33 discussed their experimental design choices, which were defined as approaches, strategies,  
34 techniques, or other decisions made during the study design process. As mentioned previously,  
35 these design choices did not necessarily have to be strictly “experimental,” since not all research  
36 projects encountered in this study were based on experiments. When asked to provide a rationale  
37 for an experimental design choice, the difference between novice and advanced student  
38 responses on the BURET-R instrument is significant ( $p < 0.001$ ). Although more advanced  
39 students generally scored higher than novices on the BURET-P instrument on providing  
40 rationale, addressing limitations, and comparing alternatives to their experimental design  
41 choices, none of these differences were statistically significant.  
42  
43  
44  
45  
46  
47  
48  
49  
50  
51  
52  
53  
54  
55  
56  
57  
58  
59  
60

1  
2  
3  
4  
5  
6  
7  
8  
9  
10  
11  
12  
13  
14  
15  
16  
17  
18  
19  
20  
21  
22  
23  
24  
25  
26  
27  
28  
29  
30  
31  
32  
33  
34  
35  
36  
37  
38  
39  
40  
41  
42  
43  
44  
45  
46  
47  
48  
49  
50  
51  
52  
53  
54  
55  
56  
57  
58  
59  
60

Examples varied broadly, from why the research group chose to study a certain topic to the specific instruments used to collect raw data. In the BURET-P interview protocol, students were asked *why* they made a given design choice *instead of* something else, and they were also asked about the *limitations* of that choice. In general, both novice and advanced students scored highly on providing a rationale for an experimental design choice related to their project for the BURET-P instrument; over half of the students scored 4 or higher, which requires a clear description of the design choice *and* an explicit rationale that integrates domain-specific content knowledge. For example, “We chose to use micro plasma atomic emissions spectroscopy because of its wide dynamic range. While there were many other instruments that would have worked similarly well, but not within this large range. And we were very uncertain as to whether we were over diluting or under diluting our samples.... We only had rough EPA guidelines to kind of guide our choices.” The marginally higher average score for advanced students compared to novices was not significant. However, advanced students did explain a greater number of their decisions than novices. To reflect this, an item was included that simply counted the number of design choices for which the student provided some rationale. This number was significantly higher ( $p < 0.05$ ) for advanced students, reflecting the greater detail in which they described their experimental design.

Students were less proficient at discussing the limitations of experimental design choices. A representative response is, “So if the standards aren't prepared correctly or if they're too high on concentration, it may negatively, it definitely will negatively affect our data. So I think that's a big limitation.” which received a 2 for only identifying user error as a possible limitation. However, some students were able to discuss limitations more fluently; for example, the following excerpt scored a 4: “The limitations of that technique are that bringing it under PBS,

1  
2  
3 which is phosphate buffered saline, only mimics the ionic concentrations. It doesn't mimic the  
4 chemical function. So [what] we'd like to do for further research is hydrate it in [inaudible],  
5  
6 which ... mimics in vivo synovial fluid.” Both novice and advanced students showed moderate  
7  
8 levels of sophistication on the “comparing alternatives” item but rarely scored as high as 4, for  
9  
10 which they needed to make a clear comparison between their choice and the alternative,  
11  
12 explaining why their choice was superior. For example, “We decided to use MPAS instead of  
13  
14 graphite furnace atomic absorption spectroscopy, even though both measure lead very well.  
15  
16 Because MPAS has a larger dynamic range, and we were very uncertain as to the concentration  
17  
18 we were gonna get.”  
19  
20  
21  
22  
23  
24

### 25 Indicator 3: Interpreting the Data

26  
27 Items for the third BURET Indicator measured the extent to which students were able to  
28  
29 analyze and interpret data in order to construct explanations and models relevant to their research  
30  
31 question. The data interpretation item for BURET-P focused on constructing explanations that  
32  
33 demonstrated domain-specific content knowledge; advanced students were significantly more  
34  
35 likely ( $p < 0.05$ ) to score higher on this item. This is consistent with results from the BURET-R  
36  
37 instrument, on which advanced students scored higher on describing or interpreting their data  
38  
39 analysis ( $p < 0.01$ ).  
40  
41  
42

43 For example, “And we found that with lower concentrations of silver, we get the same  
44  
45 amount of silver conductivity” scored a 2 because there was a clear statement about the  
46  
47 experimental results but no additional comments were made about the data or their conclusions.  
48  
49 A score of 4 required students to *explain* what they observed: “We stained the plates, which  
50  
51 contained the cellulose media, with Congo red, which is a dye that binds to cellulase. So what  
52  
53 that allowed us to do is once we washed the excess dye away, we got results that looked like this:  
54  
55  
56  
57  
58  
59  
60

1  
2  
3 the bacterial colonies that didn't produce any cellulase show no halo, and the whole plate is red,  
4 because the cellulose is still there, the dye binds, it's all still there. The ones that you see here  
5 have a halo of white, are positive results. They produce cellulase, and we know that because  
6 around the bacterial colony, is a halo where the cellulose has been degraded, and the dye doesn't  
7 bind.” This chemical biology student describes the underlying mechanism of the assay,  
8 explaining what is happening on a molecular and cellular level to justify their interpretation.  
9

10  
11  
12 While scores for data interpretation were generally relatively high, students performed  
13 less well on analyzing sources of error and uncertainty. Most students identified a clear potential  
14 source of error or expressed skepticism about their results, but less than half of the students  
15 elaborated on their answer or connected that source of error to either their experimental design or  
16 their conclusions. A more complete response might explain how the experiment was designed to  
17 control for possible sources of error. For example, “And also, to avoid error we wanted to use  
18 NMR. First, we dissolve our wristbands using deuterated chloroform, and then running that  
19 through NMR, and seeing if there are any errors that we can possibly encounter for  
20 contamination. We just wanted to make sure the wristbands were mostly silicon. We had a  
21 positive control and negative control in just the chemical that we tested.” However, most  
22 students did not discuss sources of error at this level and there were no significant differences  
23 between novice and advanced students.  
24  
25  
26  
27  
28  
29  
30  
31  
32  
33  
34  
35  
36  
37  
38  
39  
40  
41  
42  
43

#### 44 45 Indicator 4: Proposing Future Investigations

46  
47  
48 A single assessment item aligned with the final BURET Indicator measures the extent to  
49 which students are able to generate hypotheses and plan future experiments relevant to their  
50 research question and in response to their analysis and interpretation of data. This item primarily  
51 evaluates the rationale given along with the next steps for the research project proposed by the  
52  
53  
54  
55  
56  
57  
58  
59  
60

1  
2  
3 student. Interestingly, advanced students did not score significantly higher on this item than  
4  
5 novice students for either the BURET-R or -P instrument.  
6

7  
8 Students who received a 2 on this item typically suggested “more”-based continuations of  
9  
10 their work with no rationale: more trials, more substrates, more different temperatures, and so on.  
11  
12 In contrast, students who received a 4 would include a rationale that integrates domain-specific  
13  
14 content knowledge; for example, one chemistry student said, “In the future we hope to perform  
15  
16 confocal microscopy to determine the depth of infiltration, that's another common problem with  
17  
18 current scaffolds is that they'll grow in an x-y plane and spread out in a nice flat layer, but they  
19  
20 don't go into the bi-layer membrane. So that's what we're hoping to get with these fiber mats  
21  
22 later, when you spin onto a mesh collector plate you get these really nice nodes, and we're  
23  
24 hoping that cells could easily fit into those pores and infiltrate deeper into the membrane.” Most  
25  
26 students fell in between these two points; over 50% of participants scored a 3 on this item.  
27  
28  
29  
30

### 31 32 **Discussion** 33

34  
35 We have introduced two novel instruments for assessing how undergraduate researchers  
36  
37 grow in their understanding of scientific research. Our instruments assess student discussions of  
38  
39 their own research project, complementing previously published instruments that assess the  
40  
41 ability of undergraduates to answer questions about completely different research scenarios (e.g.,  
42  
43 Harsh, 2016) or specific components of the research process like experimental design (e.g.,  
44  
45 Deane et al., 2014). The BURET instruments use the Knowledge Integration framework to  
46  
47 evaluate how undergraduates develop an integrated understanding of the different components of  
48  
49 research and the scientific practices and content of their projects. Though developed primarily  
50  
51 for chemistry researchers, they can be applied to different types of research situations and across  
52  
53 various scientific disciplines, in contrast to most existing tools. Our instruments are able to  
54  
55  
56  
57  
58  
59  
60

1  
2  
3 distinguish between students at different levels of research experience, and evidence is presented  
4  
5 for their validity and reliability. Excerpts from our rich datasets of student responses provided a  
6  
7 detailed picture of how students progress in their understanding of research and helped us to  
8  
9 identify specific areas where they need more support to fully develop as researchers.  
10  
11

12  
13 Novice undergraduate students require more guidance to place their research into a larger  
14  
15 scientific context  
16  
17

18 We found that the largest difference between novice and advanced undergraduates was  
19  
20 for providing a scientific context for their work. Even though chemistry faculty reported that  
21  
22 undergraduates are sometimes given key papers to read when beginning work in a new  
23  
24 laboratory, our results show that their understanding of the connection between their  
25  
26 experimental work and the broader scientific context is often weak. The faculty interviews we  
27  
28 conducted suggest that, at least in some research groups, minimal emphasis is placed on teaching  
29  
30 novice undergraduates the scientific context of their research projects. Multiple faculty singled  
31  
32 out the first Indicator as important but “hard in some cases for undergrads, they don’t necessarily  
33  
34 see the big picture at this time.” Several faculty also mentioned that having their undergraduates  
35  
36 read the literature was a weak point in their mentoring. The lower priority given to these areas by  
37  
38 faculty mentors, particularly for novice students, may help explain why there is such an increase  
39  
40 in performance once students have been participating in research for at least two semesters.  
41  
42

43  
44 Reading the scientific literature has been shown to be challenging for novice students, but these  
45  
46 skills develop over time as they work with their graduate student mentors to read more papers  
47  
48 (Nelms and Segura-Totten, 2019). Reisner and Stewart (2020) make the case that incorporating  
49  
50 activities to engage students in reading and discussing literature is of critical importance in  
51  
52 chemistry research settings, to support students “to think more like disciplinary experts.”  
53  
54  
55  
56  
57  
58  
59  
60

1  
2  
3 However, faculty members may vary in terms of when they feel it is appropriate to better  
4 acquaint their undergraduates with the scientific literature. When they are ready, we suggest that  
5 mentors use published approaches for teaching students to read the literature (Hoskins *et al.*,  
6 2011; Krontiris-Litowitz, 2013; Sato *et al.*, 2014) in order to help their undergraduates  
7 understand the scientific context of their project more rapidly. Additionally, curriculum  
8 developers could help facilitate this process by including more interaction with the primary  
9 literature in undergraduate coursework.

10  
11 In contrast, students at all levels performed well on providing an integrated societal  
12 context for their work, and more advanced students did not receive higher average scores on this  
13 item. The ability to discuss the broader impacts of a research project is a valued skill, with some  
14 institutions offering courses explicitly aimed at training students in this area (MacFadden, 2009;  
15 Heath *et al.*, 2014). In two of the CUREs included in this study, students developed research  
16 questions, often addressing a societal issue of interest to them, and as a result, they could fluently  
17 discuss the societal relevance of their project. Because novice students were strong on this item,  
18 there was little growth with more research experience.

19  
20 Support is needed for beginning undergraduate researchers to better justify their experimental  
21 design and interpret their data

22  
23 The extent to which undergraduates are exposed to experimental design in chemistry  
24 coursework and research experiences varies widely, and students struggle with this critical skill  
25 (Espinosa, 2011; Gormally *et al.*, 2012; Laursen, 2019). Previous attempts to assess gains in  
26 experimental design ability during scientific research experiences showed a general trend that  
27 participation in a CURE or URE improves student reasoning in this area (Sirum and Humburg,  
28 2011; Dasgupta *et al.*, 2014; Harsh, 2016; Harsh *et al.*, 2017; Shanks *et al.*, 2017). However,  
29  
30

1  
2  
3 identifying the limitations of an experimental design has been found to be a weak point, even for  
4 graduate students (Gilmore *et al.*, 2015). In our work, we found that both novice and advanced  
5 students scored relatively high on their ability to rationalize experimental design choices. We  
6 also found that more advanced students recognized that rationalizing experimental design,  
7 including providing the limitations of and alternatives to their experimental design choices, is an  
8 important component of talking about their research. The difference between novice and  
9 advanced students' rationalizations of experimental design choices for BURET-R is significant  
10 ( $p < 0.001$ ). Similarly, advanced students were more likely to include limitations and alternatives  
11 as components of their presentation of their research.  
12  
13  
14  
15  
16  
17  
18  
19  
20  
21  
22  
23

24         The general trends we observed for experimental design also hold for data interpretation;  
25 we showed that students generally performed well on giving straightforward interpretations of  
26 their data but were less likely to provide a richer description unless specifically prompted. Scores  
27 on the combined data analysis and interpretation items on both instruments were relatively high,  
28 with advanced students scoring significantly higher than novice students. This is consistent with  
29 other studies showing that data interpretation skills correlate with increased research experience  
30 (White *et al.*, 2011; Harsh *et al.*, 2017). In contrast, one of the lowest scoring items for both  
31 novice and advanced students was their ability to identify and discuss potential sources of error  
32 in their work. Students may deliberately focus on more positive aspects of their project, or the  
33 low scores may reveal a genuine deficit among undergraduates, who have been shown to  
34 struggle with critically analyzing experimental designs, generating data visualizations, and  
35 interpreting chemical data (Varela *et al.*, 2005; White *et al.*, 2011). Our study suggests that  
36 students may benefit from targeted interventions in these areas throughout their undergraduate  
37 career.  
38  
39  
40  
41  
42  
43  
44  
45  
46  
47  
48  
49  
50  
51  
52  
53  
54  
55  
56  
57  
58  
59  
60



1  
2  
3 Novice and advanced students were equally proficient at proposing future work for their projects  
4  
5

6 The advanced and novice students in our sample were equally successful at proposing  
7  
8 next steps for their research projects. This was surprising, because when faculty were asked what  
9  
10 specifically they look for as signs of progress in their undergraduate researchers, many focused  
11  
12 on day-to-day independence, including “thinking about what’s next, what would be the next  
13  
14 experiment after this one.” The faculty interviewed by Laursen *et al.* (2010) also identified  
15  
16 taking initiative, making decisions, and acting independently as markers of student progress.  
17  
18

19 A concept from the literature that is closely related to the item on proposing future work  
20  
21 is that of iteration, as students scored higher when the proposed work was linked in some way to  
22  
23 their most recent results. Authentic research is an iterative process, where the data from one  
24  
25 experiment helps inform the next. Some have suggested that iteration is an essential part of an  
26  
27 undergraduate research experience (e.g., Auchincloss *et al.*, 2014), and efforts have been made to  
28  
29 explicitly include iteration in CUREs (Light *et al.*, 2019). Although there are instruments that  
30  
31 measure whether a student perceives iteration to be a part of their research experience (Corwin,  
32  
33 Graham, *et al.*, 2015), to our knowledge, there are no instruments that assess student proficiency  
34  
35 in proposing next steps for an ongoing research project.  
36  
37  
38

39 We anticipated that advanced students would be more experienced at proposing future  
40  
41 experiments and would therefore be able to more fluently discuss them in their written responses  
42  
43 and poster presentations. Although this was not reflected in the average scores, we observed that  
44  
45 only advanced students received the highest possible score for proposing next steps on either the  
46  
47 BURET-R or BURET-P instrument. Additionally, most of the advanced graduate students who  
48  
49 were interviewed during the development of the instrument (see Methods), scored at the highest  
50  
51 level on the BURET-P for this item.  
52  
53  
54  
55  
56  
57  
58  
59  
60

1  
2  
3 One potential explanation for the discrepancy between expectations and observed results  
4 for undergraduate researchers on average is that many of the advanced undergraduate presenters  
5 were weeks away from graduation. Those students were likely in the process of concluding their  
6 research and were not planning longer-term directions of the project. As a result, their scores on  
7 proposing future work might be lower than if we had interviewed them earlier. In contrast, many  
8 novice students were enrolled in a one-semester CURE in which they were explicitly instructed  
9 to talk about future work as part of their poster presentation. Their relative success in this area  
10 suggests that, contrary to faculty expectations, even novice students can be expected to propose  
11 the next steps of their research project, and this expectation should be more explicitly integrated  
12 into UREs.  
13  
14  
15  
16  
17  
18  
19  
20  
21  
22  
23  
24  
25  
26

27 The BURET instruments apply to a range of different chemistry subdisciplines  
28  
29

30 One advantage of the BURET instruments is that they can be applied to very disparate  
31 projects spanning the wide range of subfields that fall under the larger domain of chemistry. The  
32 BURET instruments attempt to account for subdiscipline-specific knowledge without being  
33 restrictive, taking into account the fact that the understanding developed by working on a  
34 synthetic organic project is quite different from what one learns doing biophysical chemistry or  
35 atmospheric chemistry. To receive a higher score of 4 on the BURET data interpretation item, a  
36 student must explain what they observed in a way that demonstrates domain-specific content  
37 knowledge relevant to their research project. Because content knowledge for a diverse sampling  
38 of undergraduate researchers can be from a variety of disciplines, it has previously been difficult  
39 to measure with existing instruments about a single hypothetical scenario. We showed that such  
40 knowledge can generally be identified using the BURET instruments. For example, all of the  
41 excerpts in Table 6 scored a 4 on data interpretation except the atmospheric chemistry passage,  
42  
43  
44  
45  
46  
47  
48  
49  
50  
51  
52  
53  
54  
55  
56  
57  
58  
59  
60

which scored a 3 because domain-specific content knowledge was vaguely alluded to instead of explicitly stated. We envision the BURET instruments being used by educational researchers to monitor student progress in a unified way across UREs and CUREs.

**Table 6.** Excerpts from poster presentation transcripts: Interpretation of observed results across various chemistry subdisciplines.

Chemistry Sub-Discipline	Excerpt from Student Poster Presentation
Biochemistry	My interpretation of these results is that the R-pal is utilizing the thiosulfate to grow and produce ammonia, so that's the main takeaway of this experiment and that if we took out thiosulfate and replaced it with another electron donor then they would grow with those electrons donated from that.
Inorganic Chemistry	What I've done here is I've synthesized a magnet that targets the lanthanide that has a strongly axial crystal field, but also a radical bridge, and this works very well because the 2,2'-bipyrimidine, that is substituted with chlorines, is a very weak epineural donor and so the crystal field becomes more axial because you have such a weak epineural donor even though you still have a radical lanthanide bridge.
Materials Chemistry	But the decrease is that prevention of growth that I was talking about, [due to] the charge neutralization of the bromide ions on the ends of the surfactant. So, if the surfactant is more packed, no more gaps are available for precipitation to occur, and so you can't grow any nanorods per se. All you're gonna be left with is a bunch of spherical nanoparticles, no growth curve. So that's the reason for this decrease.
Physical Chemistry	We first ran them on the mass spec to know that we know, that it's working as a control. So we can tell there's one peak for the full rotaxane, and then over time we can see the cleaved product come off, and that peak grows in over time. So then after 16 hours it works. So we go to do it with xenon NMR we can see the same thing after 16 hours you have a pretty full peak come in for CB-6. Here, this is the water peak for CB-6, we always see that in the xenon NMR experiment, and you see CB-6 peak, that's the xenon going in and out there.
Atmospheric Chemistry	What is shown here is the VOC reactivity to show it's relatively constant, and then the NO <sub>x</sub> concentrations, and the ozone concentrations. So the NO <sub>x</sub> decreases from weekday to weekend because there are less giant trucks driving. Then this is showing that ozone decreases, but it doesn't really decrease that much, it's basically the same.

## Limitations

We identify four potential limitations of our study.

*Self-selection bias* is a limitation of undergraduate research studies, because those who participate are likely to be among the most highly motivated and high performing students. We expect selection bias to be minimal in our case, as approximately 70% of chemistry majors, who

1  
2  
3 make up the majority of our sample, participate in undergraduate research, giving them an  
4  
5 opportunity to participate in the poster session from which we recruited our participants.  
6

7  
8 *Non-uniform experimental conditions.* Data was necessarily collected from a variety of  
9  
10 CURE and URE contexts, such that some students typed their responses, while others submitted  
11  
12 hand-written responses, leading to differences in length of response. We have attempted to  
13  
14 counteract these issues by designing the BURET instruments specifically for different contexts  
15  
16 and to deeply probe the quality of the responses. However, to ensure that accurate comparisons  
17  
18 can be made between students, it is suggested that users of the BURET-R instrument provide an  
19  
20 additional statement clarifying the desired length of the responses.  
21  
22

23  
24 *Low numbers of advanced URM participants.* Although we attempted to oversample  
25  
26 students who identify as a URM, we were only able to recruit four such students who had  
27  
28 completed at least two semesters of undergraduate research. With so few advanced researchers,  
29  
30 we were unable to determine whether our instrument has any intrinsic bias regarding URM  
31  
32 students. This will be an important feature to assess in subsequent studies. The goals identified in  
33  
34 faculty interviews were consistent with the categories we used in our coding rubrics, but  
35  
36 researchers may assume certain cultural norms about the “correct” way to answer a question that  
37  
38 is posed in a scientific setting. Although enculturation into a research program will likely result  
39  
40 in more homogeneity over time, novice students in particular may differ in their interpretation of  
41  
42 the questions purely as a result of their demographic background.  
43  
44

45  
46  
47 *Preliminary CURE Experience.* The most common trajectory for undergraduates in  
48  
49 chemistry at our institution is to take a CURE prior to starting a URE in a faculty research group.  
50  
51 At some institutions, students may start a URE without any prior CURE experience or enroll in a  
52  
53 CURE concurrently with or after participating in a URE. Because we found that BURET-R  
54  
55  
56  
57  
58  
59  
60

1  
2  
3 scores appear to be sensitive to the type of research experience, more work will be needed for its  
4 use in different universities and for comparisons across different sequencing of CURE and URE  
5 experiences.  
6  
7  
8  
9

## 10 Implications

11  
12  
13  
14 In recent years, there has been increasing attention on the development and assessment of  
15 research experiences for undergraduate STEM majors. However, as compared to the life  
16 sciences, chemistry has produced fewer studies about the development and assessment of these  
17 opportunities, and fewer instruments have been developed to support student learning of research  
18 skills such as experimental design, data analysis, and reading the primary literature in chemistry.  
19 We have used the BURET-R and BURET-P instruments to characterize the progression of  
20 student expertise and reveal weaknesses in the learning outcomes of undergraduate researchers.  
21 While time-intensive to code, we envision the use of the BURET instruments to be highly  
22 valuable in the contexts of mentoring and future research on undergraduate research experiences  
23 in chemistry. The BURET-R instrument is more generally applicable, as it can be quickly  
24 administered. For assessing poster presentations in a variety of educational contexts, the  
25 corresponding BURET-P instrument can provide a more detailed picture of student knowledge  
26 integration. These instruments offer a method of assessing student learning in relationship with  
27 students' own chemistry research projects. Because the focus is on the student's project and not  
28 on answering questions about a hypothetical scenario, the instruments are authentic and can be  
29 used across the breadth of chemistry subdisciplines. Following initial development, many  
30 surveys, interview protocols, and performance-based instruments designed to measure student  
31 learning in the sciences have been applied as stand-alone assessments or in combination with  
32 other instruments in subsequent studies. We envision that some may find it helpful to use the  
33  
34  
35  
36  
37  
38  
39  
40  
41  
42  
43  
44  
45  
46  
47  
48  
49  
50  
51  
52  
53  
54  
55  
56  
57  
58  
59  
60

1  
2  
3 BURET instruments (both protocols and rubrics together) in their entirety, as an undergraduate  
4 research inventory, while others may administer a selection to their students.  
5  
6

7  
8       Moreover, the BURET instruments provide an informal, low-stakes method for mentors  
9  
10 to check on the progression of their students. Research mentors can regularly observe students  
11 setting up and analyzing the results of experiments, but they often have fewer opportunities to  
12 probe how their undergraduate students think about the research project more broadly. The  
13 BURET-R can be used as a way to quickly gauge how the student discusses their project in  
14 response to open-ended questions. This data can serve to guide research mentors to initiate  
15 conversations with the student to strengthen their understanding of the project and to consider  
16 how to better turn what they know into an integrated narrative about their project. Additionally,  
17 the act of responding to the BURET-R prompts is itself a useful opportunity for the student to  
18 reflect on their project, which may not be a regular feature of their research experience.  
19  
20 Similarly, answering the BURET-P protocol questions is an inherently useful activity, as it can  
21 help students to strengthen their poster talks and provide practice taking questions from the  
22 audience.  
23  
24  
25  
26  
27  
28  
29  
30  
31  
32  
33  
34  
35  
36

37  
38       At a departmental or institutional level, the BURET instruments can be used at regular  
39 intervals to assess how well a particular research experience is supporting student learning as  
40 they progress from novice to advanced researchers. The BURET instruments complement self-  
41 report survey data by enabling educational researchers to directly measure student learning with  
42 respect to knowledge and skills that are critical for their development as scientists. In the event  
43 that certain BURET Indicators are of greater importance with a particular student group, specific  
44 probes, like the interview questions in the BURET-P instrument, can be used to further explore  
45 student thinking for different components of the research process. Both of the BURET  
46  
47  
48  
49  
50  
51  
52  
53  
54  
55  
56  
57  
58  
59  
60

instruments can be used to provide students with feedback about their strengths and knowledge gaps with respect to the research project they are working on in a CURE or URE. These instruments can also be used to compare different research experiences, providing individual CUREs or UREs with information about the areas in which students need additional instruction or training from their research mentors.

## Conflicts of Interest

There are no conflicts to declare.

## Appendices

### Appendix A – Abbreviated Coding Rubric

---

#### CONTEXT – SCIENTIFIC OR SOCIETAL

Indicator 1: Place the research questions or goals of their laboratory and/or project in the context of the larger field.

Score	BURET-R	BURET-P – Scientific	BURET-P – Societal
0	Does not explain goals of experiment or project	Does not explain goals of project	Does not explain goals of project
1	Partial or unclear description of experiment and/or project goals	Partial or unclear description of project goals	Only discusses “personal” goals, but does not mention a societally relevant topic
2	States goal of experiment OR States goal of project	Clearly states goal of project OR States a very limited scientific application of their work	Collecting data with no further connection to societal importance OR Reader can infer societal importance or application of the data collected (i.e. mentions a societally-relevant topic like semiconductor or cancer)
3	Clearly states goal of experiment AND States goal of project (vagueness allowed)	States a general area of science that their work contributes to OR Vague or implied version of below	Implies societal importance OR Vague statement about the possible benefits or use of results
4	Partial Link (3) AND (Explains how expt advances larger project OR Explicit link of project to broader significance (scientific or societal))	Discusses how future projects (by other labs) might be affected by current project OR Suggests new research paths or projects that could be based on this work OR Provides sufficient	Explicitly connects project to specific societal need OR Explicit statement about the possible benefits or use of results (Accurate content knowledge and coherent argument should be present.

		background for reader to understand current state of field	However, exact mechanism of connection does not need to be stated.)
5	Partial Link (3) AND Explains how expt advances larger project or the portion of the project they are working on AND Explicit link of project to broader significance (scientific or societal)	Basic Link (4) with two out of three of the criteria present OR Two different scientific contexts explained for the project - both at Basic Links (4)	Explicit comparison between current project goals and existing solutions to those problems. Exact mechanism of connection does need to be stated. OR Explicit and specific statement about the possible benefits or use of results, including statement of existing societal issue or need

---

## DESIGN CHOICES

Indicator 2: Justify their experimental design as appropriate for their research question and scientific content of their project.

Score	BURET-R & P – Rationale	BURET-P – Limitation	BURET-P – Comparison
0	Coder cannot identify any design choice discussed	Does not discuss any limitations of design choice	Does not mention any alternative design choices
1	Partial or unclear description of one design choice	Vague statement of a very generic limitation or logistical issue OR Vague or implied description of a thoughtful limitation - implied in the description of results	Mentions the fact that there are alternatives, but doesn't mention what these are.
2	Clear description of one design choice, but rationale is poor or absent	Clear statement of a very generic limitation or logistical issue OR Vague description of a thoughtful limitation	Mentions specific alternative, but no comparison OR Compares to alternative because alternative is, in their opinion "not possible"
3	Clear description of one design choice AND Gives reasonable (sounding) rationale but vague, implied, or invokes little to no content knowledge	One or more thoughtful limitations mentioned, but content knowledge only implied	Compares design choice to an alternative, but is somewhat vague or implied OR Compares to alternative because alternative is, in their opinion "not possible", plus <i>why</i> it wouldn't be possible
4	Clear description of one design choice AND Gives explicit rationale for choice of instrument or experiment that integrates domain-specific content knowledge	Gives at least one explicit limitation that integrates domain-specific content knowledge	Comparison to an alternative design choice on a single facet with a clear statement of difference or advantage or reason to use one or the other
5	Basic Link (4) but multiple distinct reasons for design choice are discussed AND Strong evidence of	Basic Link (4) AND (Discusses how limitations affect conclusions OR	Clear comparison to an alternative design choice on



extensive content knowledge that supports their choices	Discusses how limitation was addressed, minimized, avoided, etc.)	more than one facet OR 3 or more Basic Links (4)
---	---	--

---

### VARIABLES - WHAT IS BEING MEASURED, MANIPULATED, OR COMPARED?

Indicator 2: Justify their experimental design as appropriate for their research question and scientific content of their project.

Score	BURET-R Definition
0	Does not indicate what type of data is being collected or discuss any other relevant variables
1	Isolated Concept but vague or implied (unclear what they are actually measuring, manipulating, comparing) OR Basic instrument verification on a standard
2	Clearly identifies what is being measured (raw OR analyzed) OR Clearly identifies one or more variables being manipulated, compared, or held constant
3	Isolated Concept (2) AND (Provides basic rationale for choice of variables and/or range being investigated OR Gives details on how or to what extent the variables are manipulated) OR Basic link (4), but rationale or predictions are vague or questionable
4	Clearly identifies what is being measured (raw OR analyzed) AND Clearly states one or more variables being manipulated, controlled or compared AND (Provides rationale (clear, but slightly generic okay) for why manipulated variables would affect measurements/output OR Provides reasonable prediction of how manipulated variables will affect output)
5	Basic Link (4) AND Rationale and/or predictions are strong and integrate content knowledge

---

### DATA MANIPULATION AND INTERPRETATION

Indicator 3: Analyze and interpret data in order to construct explanations and models that are relevant to their research question.

Score	BURET-R – Manipulation	BURET-R and P – Interpretation
0	Does not describe any analysis of raw data	Does not describe results OR Has not collected data yet
1	States that no data analysis was performed OR States that results are inconclusive with no elaboration	Unclear how conclusion is supported by results OR Implies data interpretation but does not sufficiently describe
2	States a procedure for analyzing or manipulating data with no elaboration	Summarizes results without interpretation OR Pre-packaged conclusion OR States an interpretation with no connection to data
3	Links raw data to analyzed results, but discussion of data or analysis method/procedure is vague	Summarizes results and links to content knowledge or compares to expectations, but vague or minimal insights
4	Clearly links raw data to analyzed results, including (clear) description of the analysis process	Gives plausible explanation for results (or compares results to expectations in a way) that integrates clear content knowledge

5	Basic Link (4), plus discusses at least one assumption or consequential decision made during analysis	Basic Link (4), but integrates extensive content knowledge OR Discusses alternate interpretations
---	---	---

---

### CONFIDENCE/ERROR ANALYSIS

Indicator 3: Analyze and interpret data in order to construct explanations and models that are relevant to their research question.

Score	BURET-P Definition
0	Does not identify any potential sources of error
1	States that the experiment (or a large part of it) “didn’t work” without any elaboration as to why OR Describes confidence in the ability of methods to answer the RQ
2	Identifies a clear “error” in what was done OR Vague reference to limitation of method/technique when discussing confidence in results OR Vague “doubts” about data
3	Identifies potential sources of error that are less “obvious” OR Clear reference to limitation of method/technique when discussing confidence in results
4	Clearly identifies potential reasonable source(s) of error AND Mentions how these connect to at least <i>one</i> of the following: 1. Research questions; 2. Experimental design (current or future); 3. Their conclusions
5	Clearly identifies multiple distinct potential reasonable source(s) of error at the level of a Basic Link (4)

---

### NEXT STEPS

Indicator 4: Generate hypotheses and plan future experiments in response to their analysis and interpretation of data and research question.

Score	BURET-R and P Definition
0	Does not discuss any potential future work
1	Completely different goals for future work with no/minimal relationship to current work OR Implies that they will “continue with the plan” but does not sufficiently describe
2	Simple quantitative extension, modification, or new experiment with no or poor rationale OR “Continue with the plan” OR Repeat experiment with simple issue fixed
3	Simple quantitative extension with good rationale OR Modification or new experiment with credible but vague rationale OR Repeat experiment after difficult-to-predict issue fixed (troubleshooting), link to content knowledge is vague or absent
4	Modification, troubleshooting, or new experiment with clear rationale that integrates content knowledge
5	Multiple Basic Links (4), at least one of which is not a borderline Partial Link (3) OR (Basic Link AND Explicitly links new choices to the <i>results</i> of current work)

---

### PREVIOUS WORK

Indicator 1: Place the research questions or goals of their laboratory and/or project in the context of the larger field.

Score	BURET-P Definition
0	Does not mention any prior work
1	Vague references to “other studies” without any specific designs/results or clear specification of how this informs part of project
2	Clear reference to previous work, but no stated connection to current work OR Vague reference to previous work with connection to current project

- 1  
2  
3 3 Clear description of previous design or results AND (Vague connection to/influence on current work  
4 OR Vague comparison b/w old and new design or results)  
5 4 Summarizes previous work (specific design or results) AND (Explicitly states how it connects  
6 to/influenced current work OR Compares to current results)  
7 5 Basic Link (4) AND (Explanation of how current work is different or novel OR Attempts to interpret  
8 sim/diff between current and previous results)  
9

---

### 11 INTEGRATION OF (ADDITIONAL) CONTENT KNOWLEDGE

12  
13 Indicator 1: Place the research questions or goals of their laboratory and/or project in the context of the larger  
14 field.

15 Score	BURET-R	BURET-P
16 0	Response does not integrate any scientific content 17 knowledge beyond what is necessary to describe 18 the project	Response does not integrate any scientific content knowledge beyond what is necessary to describe 19 the project
20 1	(not used)	(not used)
21 2	Weak example of a Partial Link (3)	Weak example of a Partial Link (3)
22 3	Exhibits scientific content knowledge beyond what is required to describe project	Exhibits scientific content knowledge beyond what is required to describe project
23 4	(not used)	Exhibits <i>extensive</i> scientific content knowledge beyond what is required to describe project
24 5	Exhibits <i>extensive</i> scientific content knowledge beyond what is required to describe project	Multiple Basic Links (4)

### 29 Appendix B – Partial Coding Rubric with Examples

31 Score	Description	Examples
32 0	- Does not discuss any 33 limitations of design choice	
34 1	- Vague reference to limitations	- “Again, part of the main problem is that graphite furnace is really 35 temperamental.”
36 2	- Clear statement of a generic 37 limitation, OR 38 - Vague description of 39 thoughtful limitation	- “In terms of that technique, I think it depends on the accuracy in 40 which the solutions are prepared. So if the standards aren't prepared 41 correctly or if they're too high on concentration, it may negatively, it 42 definitely will negatively affect our data. So I think that's a big 43 limitation. And also you have to produce a lot of different samples, 44 which can be time consuming.”
45 3	- One or more thoughtful 46 limitations mentioned, but 47 content knowledge only 48 implied	- “The limitations of Congo Red is that it is visual. So it is qualitative 49 even though we can't measure the radius. The radius isn't really going 50 to tell us anything numerical about how much cellulose the bacteria 51 digests.”
52 4	- Gives at least one explicit 53 limitation that integrates 54 domain-specific content 55 knowledge	- “One experimental technique that we use is hydrating the sample 56 and then putting them under nanoindentation. ... So the limitations of 57 that technique are that you're running it under PBS, which is 58 phosphate buffered saline, and that only mimics the ionic 59 concentrations, it doesn't mimic the chemical functionality you'd 60 encounter in in vivo synovial fluid.”
61 5	- Basic Link (4) AND	- “The main limitation is that the scaled particle theory ignores the entropic consideration in the energy of interaction here, so it's hard to say what would happen at different temperatures. In order to predict

<p>1</p> <p>2</p> <p>3</p> <p>4</p> <p>5</p> <p>6</p> <p>7</p> <p>8</p> <p>9</p> <p>10</p> <p>11</p> <p>12</p> <p>13</p> <p>14</p> <p>15</p> <p>16</p> <p>17</p> <p>18</p> <p>19</p> <p>20</p> <p>21</p> <p>22</p> <p>23</p> <p>24</p> <p>25</p> <p>26</p> <p>27</p> <p>28</p> <p>29</p> <p>30</p> <p>31</p> <p>32</p> <p>33</p> <p>34</p> <p>35</p> <p>36</p> <p>37</p> <p>38</p> <p>39</p> <p>40</p> <p>41</p> <p>42</p> <p>43</p> <p>44</p> <p>45</p> <p>46</p> <p>47</p> <p>48</p> <p>49</p> <p>50</p> <p>51</p> <p>52</p> <p>53</p> <p>54</p> <p>55</p> <p>56</p> <p>57</p> <p>58</p> <p>59</p> <p>60</p>	<p>(- Discusses how limitations affected conclusions OR - Discusses how limitation was addressed, minimized, avoided, etc.)</p>	<p>the temperature dependence, you need an approximate value of the entropy of dissolution, which isn't known for a lot of these molecules. However, we found that that's actually very easy to predict. For each group of molecules it's approximately constant for a certain chlorination number so you know that if you have a PCB and it has three chlorines that you will know the entropy very well."</p> <hr/>
--	---	---

## Appendix C – Complete Sampling Procedures

All of the students enrolled in the three target CUREs were invited to be part of this study in person by one of the researchers. A few weeks before their corresponding poster session, students were asked to respond to the two reflective prompts, and answers were collected from all students who agreed to participate. This included 135 members of the chemistry CURE, for a response rate of 58%, 11 students in the pre-service teacher CURE, for a response rate of 65%, and 9 participants from the transfer student CURE, for a response rate of 90%. All of the consenting students in the pre-service teacher and transfer student CUREs were then interviewed at the final poster session for those courses. For the chemistry CURE, a subset of 35 consenting students were interviewed at the final poster session. To choose a representative sample, students were stratified by major and prior research experience. After intentionally oversampling 7 URM students, a random sample of 28 was chosen from among the remaining 128 students. From this pool of CURE participants, 6 URM students and a random sample of 9 other students were chosen for further analysis. An additional 20 students for whom we had prompt responses but not poster session interviews were also randomly selected for analysis.

Additionally, all students presenting at one of the two target URE poster sessions were invited by email to participate in this study. For the chemistry poster session, 112 students were invited to participate and 66 consented, for a response rate of 59%. For the pre-service teacher poster session, 23 of the 26 students (88%) responded affirmatively to the invitation. Responses to the reflective prompts were collected via Qualtrics a few weeks prior to the poster sessions.

1  
2  
3 All consenting students who provided answers to our reflective prompts (30 from the chemistry  
4 poster session and 23 from the pre-service teacher session) were interviewed at these poster  
5 sessions, using the same protocol. From this pool of URE participants, 6 URM students and a  
6 random sample of 24 other students were chosen for further analysis. An additional 15 students  
7 for whom we had prompt responses but not poster session interviews were also randomly  
8 selected for analysis. In total, the dataset we analyzed included 80 responses to reflective  
9 prompts and 55 poster session interviews.  
10  
11  
12  
13  
14  
15  
16  
17  
18  
19  
20  
21

## 22 Appendix D – Psychometric Analysis

23  
24 Factor analysis provided evidence that a unidimensional construct is being measured. For  
25 each instrument, only one factor had an eigenvalue greater than 1, and the ratio of the first two  
26 eigenvalues was well above 4. Additionally, all items except one had a correlation of at least  $r =$   
27 0.55 with the overall score on the corresponding instrument. The sample size is expected to be  
28 sufficient for the one factor solution that used 6 variables for BURET-R and the one factor  
29 solution that used 11 variables for BURET-P (Mundfrom *et al.*, 2005).  
30  
31  
32  
33  
34  
35  
36  
37

38 Item-response theory (IRT) analysis was then conducted to establish the internal structure  
39 at the instrument level (Wilson, 2005). Because the sample size was not sufficient to run the  
40 analysis using all thresholds from the rubrics, data were collapsed into scores of low (0-2),  
41 moderate (3), or high (4-5), and Wright maps for each instrument were generated from the  
42 collapsed data, and there was at least one response for each possible answer choice in order to fit  
43 the data to an item response model. The resulting Wright maps (see Figs. D1 and D2 below)  
44 show that the range of instrument item logit values span nearly the entire distribution of  
45 respondent logit values, with only a few students falling below all item thresholds on the  
46  
47  
48  
49  
50  
51  
52  
53  
54  
55  
56  
57  
58  
59  
60

BURET-R instrument, and a few Thurstonian thresholds located below the lowest respondent logit value for the BURET-P instrument. The reliability of partial credit model analysis carried out on the data is 0.77 for BURET-R and 0.76 for BURET-P. These values indicate an acceptable consistency of the items to measure respondent performance (Bond & Fox, 2007; Wright & Masters, 1982).

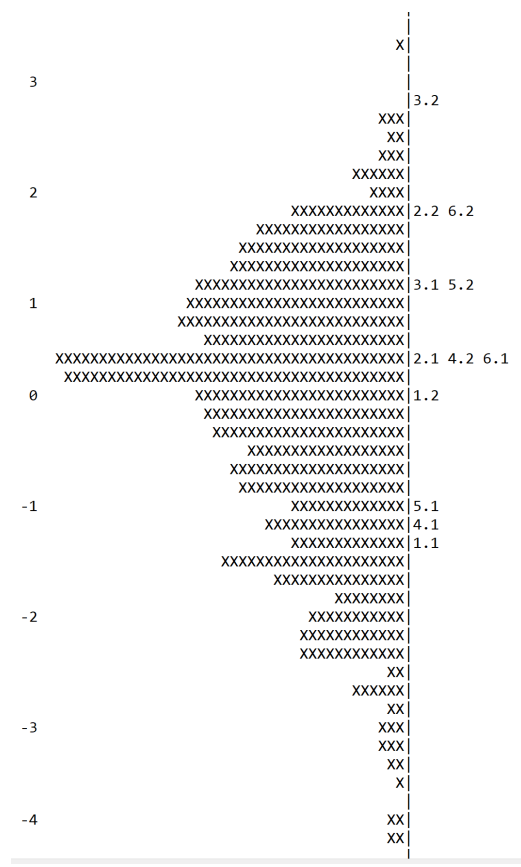


Fig. D1. Wright map for BURET-R

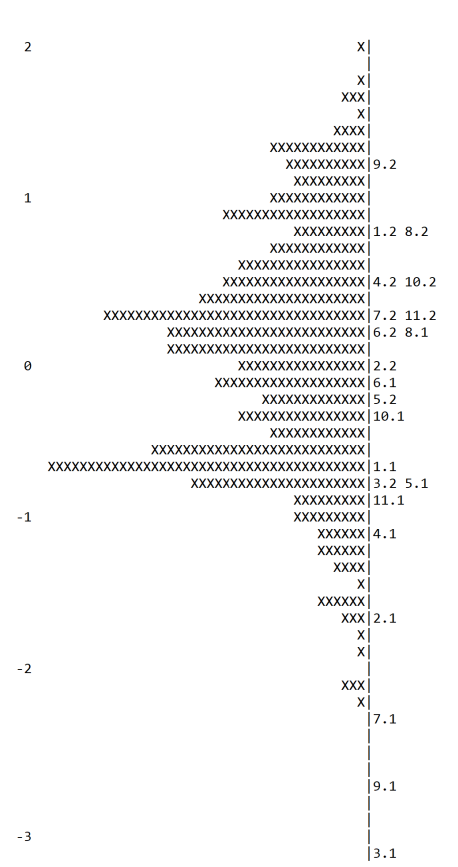
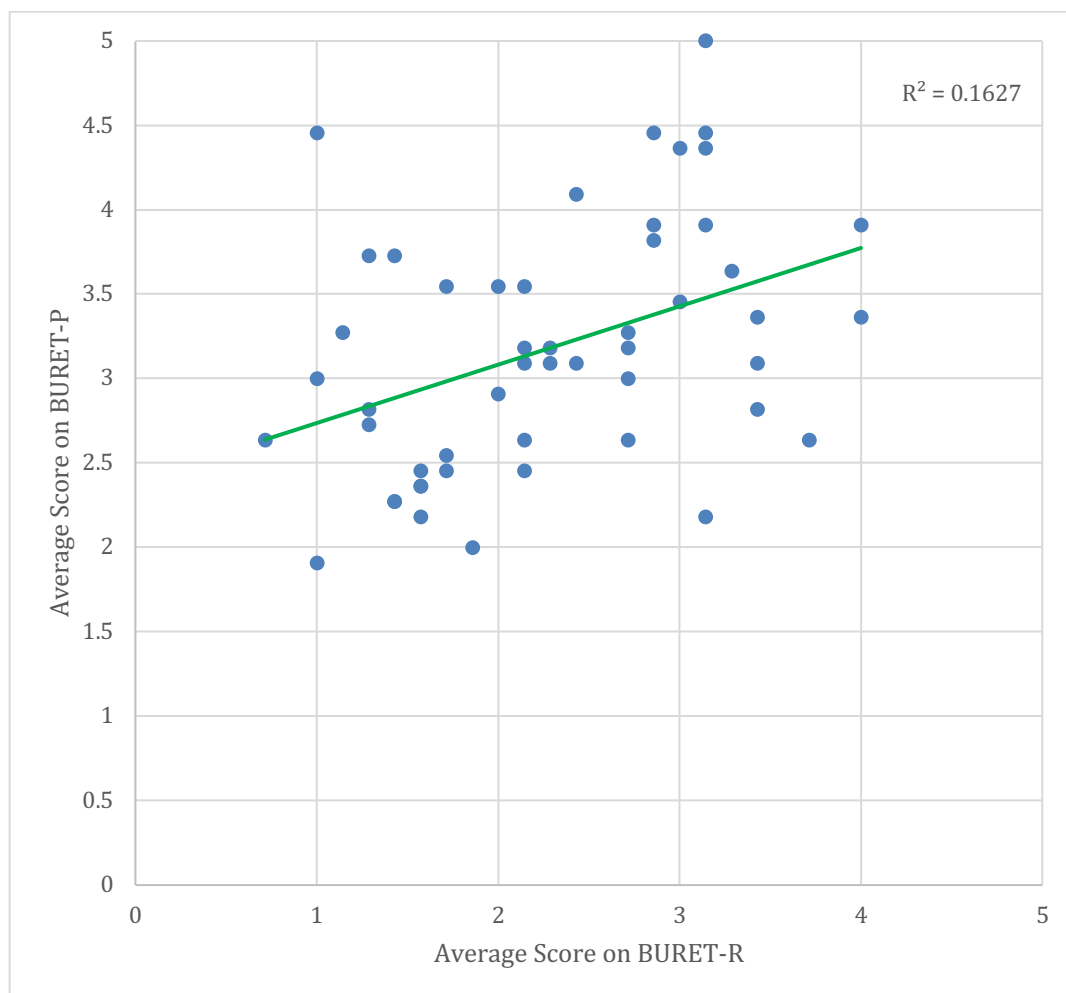


Fig. D2. Wright map for BURET-P

Appendix E – Scatterplot for BURET-R vs. BURET-P scores



## Acknowledgements

This work was funded by the National Science Foundation Improving Undergraduate STEM Education initiative (NSF IUSE 1712001), and the Alfred P. Sloan Foundation (G-2016-7112). Additional support comes from the Berkeley Graduate School of Education Barbara Y. White Fund and a National Science Foundation Research Experiences for Teachers (RET) Site Award (NSF EEC-1542471). L.E.C. was supported by the National Science Foundation Graduate Research Fellowship Program under Grant No. DGE 1106400. Thanks to Gabe Otero, Erica Dettmer-Radtke, Ibrahim Hajar, Seth Van Doren, Jiho Kim, Yuki Watanabe, Michelle Douskey, Zachary Firestein, Eddy Ham, Michelle Wilkerson, Vicky Laina for assistance with

1  
2  
3 data collection and Laura Armstrong for assistance with data analysis. Any opinions, findings,  
4 conclusions, or recommendations expressed in this article are those of the authors and do not  
5 necessarily represent the views of these agencies.  
6  
7  
8  
9

## 11 References

- 12  
13  
14 Ashcroft, J., Blatti, J., & Jaramillo, V. (2020), Early career undergraduate research as a meaningful  
15 academic experience in which students develop professional workforce skills: A community  
16 college perspective. In *Integrating Professional Skills into Undergraduate Chemistry Curricula*  
17 (pp. 281-299). American Chemical Society. DOI: 10.1021/bk-2020-1365.ch016  
18 Airey J. and Linder C., (2009), A disciplinary discourse perspective on university science learning:  
19 Achieving fluency in a critical constellation of modes. *J. Res. Sci. Teach. Off. J. Natl. Assoc. Res.*  
20 *Sci. Teach.*, **46**(1), 27–49.  
21 Auchincloss L. C., Laursen S. L., Branchaw J. L., Eagan K., Graham M., Hanauer D. I., et al., (2014),  
22 *Assessment of course-based undergraduate research experiences: A meeting report. CBE-Life*  
23 *Sci. Educ.*, **13**(1), 29-40.  
24 Avargil, S., Kohen, Z., & Dori, Y. J. (2020), Trends and perceptions of choosing chemistry as a major  
25 and a career. *Chem. Educ. Res. Pract.*, **21**(2), 668-684.  
26 Bond, T. G., & Fox, C. M. (2007). *Fundamental measurement in the human sciences. Chicago, IL:*  
27 *Institute for Objective Measurement.*  
28 Bonous-Hammarth M., (2000), Pathways to success: Affirming opportunities for science, mathematics,  
29 and engineering majors. *J. Negro Educ.*, 92–111.  
30 Bransford J. D., Brown A. L., and Cocking R. R., (2000), *How people learn*, Washington, DC: National  
31 academy press.  
32 Brown J. S., Collins A., and Duguid P., (1989), Situated cognition and the culture of learning. *Educ. Res.*,  
33 **18**(1), 32–42.  
34 Butz A. R. and Branchaw J. L., (2020), Entering Research Learning Assessment (ERLA): Validity  
35 evidence for an instrument to measure undergraduate and graduate research trainee development.  
36 *CBE—Life Sci. Educ.*, **19**(2), ar18.  
37 Carlone H. B. and Johnson A., (2007), Understanding the science experiences of successful women of  
38 color: Science identity as an analytic lens. *J. Res. Sci. Teach. Off. J. Natl. Assoc. Res. Sci. Teach.*,  
39 **44**(8), 1187–1218.  
40 Carpi A., Ronan D. M., Falconer H. M., and Lents N. H., (2017), Cultivating minority scientists:  
41 Undergraduate research increases self-efficacy and career ambitions for underrepresented  
42 students in STEM. *J. Res. Sci. Teach.*, **54**(2), 169–194.  
43 Chang M. J., Eagan M. K., Lin M. H., and Hurtado S., (2011), Considering the impact of racial stigmas  
44 and science identity: Persistence among biomedical and behavioral science aspirants. *J. High.*  
45 *Educ.*, **82**(5), 564–596.  
46 Chang M. J., Sharkness J., Hurtado S., and Newman C. B., (2014), What matters in college for retaining  
47 aspiring scientists and engineers from underrepresented racial groups. *J. Res. Sci. Teach.*, **51**(5),  
48 555–580.  
49 Chase A. M., Clancy H. A., Lachance R. P., Mathison B. M., Chiu M. M., and Weaver G. C., (2017),  
50 Improving critical thinking via authenticity: the CASPiE research experience in a military  
51 academy chemistry course. *Chem. Educ. Res. Pract.*, **18**(1), 55–63.  
52 Clark T. M., Ricciardo R., and Weaver T., (2016), Transitioning from expository laboratory experiments  
53 to course-based undergraduate research in general chemistry. *J. Chem. Educ.*, **93**(1), 56–63.  
54  
55  
56  
57  
58  
59  
60



- 1  
2  
3 Coil D., Wenderoth M. P., Cunningham M., and Dirks C., (2010), Teaching the process of science:  
4 faculty perceptions and an effective methodology. *CBE—Life Sci. Educ.*, **9**(4), 524–535.
- 5 Connor M. C., Finkenstaedt-Quinn S. A., and Shultz G. V., (2019), Constraints on organic chemistry  
6 students' reasoning during IR and <sup>1</sup>H NMR spectral interpretation. *Chem. Educ. Res. Pract.*,  
7 **20**(3), 522–541.
- 8 Corwin L. A., Graham M. J., and Dolan E. L., (2015), Modeling course-based undergraduate research  
9 experiences: An agenda for future research and evaluation. *CBE—Life Sci. Educ.*, **14**(1), es1.
- 10 Corwin L. A., Runyon C. R., Ghanem E., Sandy M., Clark G., Palmer G. C., et al., (2018), Effects of  
11 discovery, iteration, and collaboration in laboratory courses on undergraduates' research career  
12 intentions fully mediated by student ownership. *CBE—Life Sci. Educ.*, **17**(2), ar20.
- 13 Corwin L. A., Runyon C., Robinson A., and Dolan E. L., (2015), The laboratory course assessment  
14 survey: a tool to measure three dimensions of research-course design. *CBE—Life Sci. Educ.*,  
15 **14**(4), ar37.
- 16 Crawford G. L. and Kloepper K. D., (2019), Exit interviews: laboratory assessment incorporating written  
17 and oral communication. *J. Chem. Educ.*, **96**(5), 880–887.
- 18 Cruz C. L., Holmberg-Douglas N., Onuska N. P., McManus J. B., MacKenzie I. A., Hutson B. L., et al.,  
19 (2020), Development of a Large-Enrollment Course-Based Research Experience in an  
20 Undergraduate Organic Chemistry Laboratory: Structure–Function Relationships in Pyrylium  
21 Photoredox Catalysts. *J. Chem. Educ.*, **97**(6), 1572–1578.
- 22 Danczak S. M., Thompson C. D., and Overton T. L., (2020), Development and validation of an  
23 instrument to measure undergraduate chemistry students' critical thinking skills. *Chem. Educ.*  
24 *Res. Pract.*, **21**(1), 62–78.
- 25 Danowitz A. M., Brown R. C., Jones C. D., Diegelman-Parente A., and Taylor C. E., (2016), A  
26 combination course and lab-based approach to teaching research skills to undergraduates. *J.*  
27 *Chem. Educ.*, **93**(3), 434–438.
- 28 Dasgupta A. P., Anderson T. R., and Pelaez N., (2014), Development and validation of a rubric for  
29 diagnosing students' experimental design knowledge and difficulties. *CBE—Life Sci. Educ.*,  
30 **13**(2), 265–284.
- 31 Dasgupta A. P., Anderson T. R., and Pelaez N. J., (2016), Development of the neuron assessment for  
32 measuring biology students' use of experimental design concepts and representations. *CBE—Life*  
33 *Sci. Educ.*, **15**(2), ar10.
- 34 Deane T., Nomme K., Jeffery E., Pollock C., and Birol G., (2014), Development of the biological  
35 experimental design concept inventory (BEDCI). *CBE—Life Sci. Educ.*, **13**(3), 540–551.
- 36 Doğan A. and Kaya O. N., (2009), Poster sessions as an authentic assessment approach in an open-Ended  
37 University general chemistry laboratory. *Procedia-Soc. Behav. Sci.*, **1**(1), 829–833.
- 38 Esparza, D., Wagler, A. E., & Olimpo, J. T. (2020), Characterization of instructor and student behaviors  
39 in CURE and Non-CURE learning environments: Impacts on student motivation, science identity  
40 development, and perceptions of the laboratory experience. *CBE—Life Sci. Educ.*, **19**(1), ar10.
- 41 Espinosa L., (2011), Pipelines and pathways: Women of color in undergraduate STEM majors and the  
42 college experiences that contribute to persistence. *Harv. Educ. Rev.*, **81**(2), 209–241.
- 43 Estrada M., Burnett M., Campbell A. G., Campbell P. B., Denetclaw W. F., Gutiérrez C. G., et al., (2016),  
44 Improving underrepresented minority student persistence in STEM. *CBE—Life Sci. Educ.*, **15**(3),  
45 es5.
- 46 Ghanem E., Long S. R., Rodenbusch S. E., Shear R. I., Beckham J. T., Procko K., et al., (2018), Teaching  
47 through Research: Alignment of Core Chemistry Competencies and Skills within a  
48 Multidisciplinary Research Framework. *J. Chem. Educ.*, **95**(2), 248–258.
- 49 Gilmore J., Vieyra M., Timmerman B., Feldon D., and Maher M., (2015), The Relationship between  
50 Undergraduate Research Participation and Subsequent Research Performance of Early Career  
51 STEM Graduate Students. *J. High. Educ.*, **86**(6), 834–863.
- 52 Gin L. E., Rowland A. A., Steinwand B., Bruno J., and Corwin L. A., (2018), Students who fail to  
53 achieve predefined research goals may still experience many positive outcomes as a result of  
54  
55  
56  
57  
58  
59  
60

- CURE participation. *CBE—Life Sci. Educ.*, **17**(4), ar57.
- Goodey N. M. and Talgar C. P., (2016), Guided inquiry in a biochemistry laboratory course improves experimental design ability. *Chem. Educ. Res. Pract.*, **17**(4), 1127–1144.
- Gormally C., Brickman P., and Lutz M., (2012), Developing a test of scientific literacy skills (TOSLS): Measuring undergraduates' evaluation of scientific information and arguments. *CBE—Life Sci. Educ.*, **11**(4), 364–377.
- Griffeth N., Batista N., Grosso T., Arianna G., Bhatia R., Boukerche F., et al., (2015), An Undergraduate Research Experience Studying Ras and Ras Mutants. *IEEE Trans. Educ.*, **59**(2), 91–97.
- Harackiewicz J. M. and Hulleman C. S., (2010), The importance of interest: The role of achievement goals and task values in promoting the development of interest. *Soc. Personal. Psychol. Compass*, **4**(1), 42–52.
- Harsh J. A., (2016), Designing performance-based measures to assess the scientific thinking skills of chemistry undergraduate researchers. *Chem. Educ. Res. Pract.*, **17**(4), 808–817.
- Harsh J., Esteb J. J., and Maltese A. V., (2017), Evaluating the development of chemistry undergraduate researchers' scientific thinking skills using performance-data: first findings from the performance assessment of undergraduate research (PURE) instrument. *Chem. Educ. Res. Pract.*, **18**(3), 472–485.
- Hauwiller M. R., Ondry J. C., Calvin J. J., Baranger A. M., and Alivisatos A. P., (2019), Translatable Research Group-Based Undergraduate Research Program for Lower-Division Students. *J. Chem. Educ.*, **96**(9), 1881–1890.
- Heath K. D., Bagley E., Berkey A. J. M., Birlenbach D. M., Carr-Markell M. K., Crawford J. W., et al., (2014), Amplify the Signal: Graduate Training in Broader Impacts of Scientific Research. *BioScience*, **64**(6), 517–523.
- Heller S. T., Duncan A. P., Moy C. L., and Kirk S. R., (2020), The Value of Failure: A Student-Driven Course-Based Research Experience in an Undergraduate Organic Chemistry Lab Inspired by an Unexpected Result. *J. Chem. Educ.*, **97**(10), 3609–3616.
- Hernandez P. R., Woodcock A., Estrada M., and Schultz P. W., (2018), Undergraduate research experiences broaden diversity in the scientific workforce. *BioScience*, **68**(3), 204–211.
- Hogan K. and Maglienti M., (2001), Comparing the epistemological underpinnings of students' and scientists' reasoning about conclusions. *J. Res. Sci. Teach. Off. J. Natl. Assoc. Res. Sci. Teach.*, **38**(6), 663–687.
- Hoskins S. G., Lopatto D., and Stevens L. M., (2011), The CREATE approach to primary literature shifts undergraduates' self-assessed ability to read and analyze journal articles, attitudes about science, and epistemological beliefs. *CBE—Life Sci. Educ.*, **10**(4), 368–378.
- Kerr M. A. and Yan F., (2016), Incorporating course-based undergraduate research experiences into analytical chemistry laboratory curricula. *J. Chem. Educ.*, **93**(4), 658–662.
- Killpack T. L. and Fulmer S. M., (2018), Development of a Tool to Assess Interrelated Experimental Design in Introductory Biology. *J. Microbiol. Biol. Educ.*, **19**(3).
- Krim J. S., Coté L. E., Schwartz R. S., Stone E. M., Cleeves J. J., Barry K. J., et al., (2019), Models and Impacts of Science Research Experiences: A Review of the Literature of CUREs, UREs, and TREs. *CBE—Life Sci. Educ.*, **18**(4), ar65.
- Krontiris-Litowitz J., (2013), Using Primary Literature to Teach Science Literacy to Introductory Biology Students. *J. Microbiol. Biol. Educ.*, **14**(1), 66–77.
- Landis J. R. and Koch G. G., (1977), The measurement of observer agreement for categorical data. *biometrics*, 159–174.
- Laungani R., Tanner C., Brooks T. D., Clement B., Clouse M., Doyle E., et al., (2018), Finding some good in an invasive species: introduction and assessment of a novel CURE to improve experimental design in undergraduate biology classrooms. *J. Microbiol. Biol. Educ.*, **19**(2).
- Laursen S., (2019), Levers for Change: An Assessment of Progress on Changing STEM Instruction: Executive Summary, American Association for the Advancement of Science.  
[https://www.aaas.org/sites/default/files/2019-07/levers-for-change-WEB100\\_2019.pdf](https://www.aaas.org/sites/default/files/2019-07/levers-for-change-WEB100_2019.pdf)

- 1  
2  
3 Laursen S., Hunter A.-B., Seymour E., Thiry H., and Melton G., (2010), *Undergraduate research in the*  
4 *sciences: Engaging students in real science*, John Wiley & Sons.
- 5 Lave J. and Wenger E., (1991), *Situated learning: Legitimate peripheral participation*, Cambridge  
6 university press.
- 7 Light C. J., Fegley M., and Stamp N., (2019), Emphasizing iterative practices for a sequential course-  
8 based undergraduate research experience in microbial biofilms. *FEMS Microbiol. Lett.*, **366**(23),  
9 fnaa001.
- 10 Lin T.-J., Lin T.-C., Potvin P., and Tsai C.-C., (2019), Research trends in science education from 2013 to  
11 2017: a systematic content analysis of publications in selected journals. *Int. J. Sci. Educ.*, **41**(3),  
12 367–387.
- 13 Linn M. C., (1995), Designing computer learning environments for engineering and computer science:  
14 The scaffolded knowledge integration framework. *J. Sci. Educ. Technol.*, **4**(2), 103–126.
- 15 Linn M. C. and Eylon B.-S., (2011), *Science learning and instruction: Taking advantage of technology to*  
16 *promote knowledge integration*, Routledge.
- 17 Linn M. C., Palmer E., Baranger A., Gerard E., and Stone E., (2015), Undergraduate research  
18 experiences: Impacts and opportunities. *Science*, **347**(6222), 1261757–1261757.
- 19 Linn M., Eylon B.-S., Kidron A., Gerard L., Toutkoushian E., Ryoo K., et al., (2018), Knowledge  
20 integration in the digital age: Trajectories, opportunities and future directions, International  
21 Society of the Learning Sciences, Inc.[ISLS].
- 22 MacFadden B. J., (2009), Training the Next Generation of Scientists about Broader Impacts. *Soc.*  
23 *Epistemol.*, **23**(3–4), 239–248.
- 24 Maltese A. V., Harsh J. A., and Svetina D., (2015), Data visualization literacy: Investigating data  
25 interpretation along the novice—expert continuum. *J. Coll. Sci. Teach.*, **45**(1), 84–90.
- 26 Mondisa J.-L. and McComb S. A., (2018), The role of social community and individual differences in  
27 minority mentoring programs. *Mentor. Tutoring Partnersh. Learn.*, **26**(1), 91–113.
- 28 Moon A., Zotos E., Finkenstaedt-Quinn S., Gere A. R., and Shultz G., (2018), Investigation of the role of  
29 writing-to-learn in promoting student understanding of light–matter interactions. *Chem. Educ.*  
30 *Res. Pract.*, **19**(3), 807–818.
- 31 Muna G. W., (2021), Stimulating Students’ Learning in Analytical Chemistry through an Environmental-  
32 Based CURE Project. *J. Chem. Educ.*, **98**(4), 1221-1226.
- 33 Mundfrom D.J., Shaw D.J., Ke T.L. (2005) Minimum Sample Size Recommendations for Conducting  
34 Factor Analyses, *International Journal of Testing*, **5**(2), 159-168.
- 35 Mutambuki J. M., Fyneweever H., Douglass K., Cobern W. W., and Obare S. O., (2019), Integrating  
36 Authentic Research Experiences into the Quantitative Analysis Chemistry Laboratory Course:  
37 STEM Majors’ Self-Reported Perceptions and Experiences. *J. Chem. Educ.*, **96**(8), 1591–1599.
- 38 National Academies of Sciences and Medicine, (2017), *Undergraduate research experiences for STEM*  
39 *students: Successes, challenges, and opportunities*, National Academies Press.
- 40 Nelms A. A. and Segura-Totten M., (2019), Expert–Novice Comparison Reveals Pedagogical  
41 Implications for Students’ Analysis of Primary Literature. *CBE—Life Sci. Educ.*, **18**(4), ar56.
- 42 O’Donnell K., Botelho J., Brown J., González G. M., and Head W., (2015), Undergraduate research and  
43 its impact on student success for underrepresented students. *New Dir. High. Educ.*, **2015**(169),  
44 27–38.
- 45 Olson S. and Riordan D. G., (2012), Engage to Excel: Producing One Million Additional College  
46 Graduates with Degrees in Science, Technology, Engineering, and Mathematics. Report to the  
47 President. *Exec. Off. Pres.* <https://files.eric.ed.gov/fulltext/ED541511.pdf>
- 48 Opitz A., Heene M., and Fischer F., (2017), Measuring scientific reasoning—a review of test instruments.  
49 *Educ. Res. Eval.*, **23**(3–4), 78–101.
- 50 Ortiz N. A., Morton T. R., Miles M. L., and Roby R. S., (2020), What About Us? Exploring the  
51 Challenges and Sources of Support Influencing Black Students’ STEM Identity Development in  
52 Postsecondary Education. *J. Negro Educ.*, **88**(3), 311–326.
- 53 Pagano J. K., Jaworski L., Lopatto D., and Waterman R., (2018), An inorganic chemistry laboratory  
54  
55  
56  
57  
58  
59  
60

- course as research. *J. Chem. Educ.*, **95**(9), 1520–1525.
- Peteroy-Kelly M. A., Marcello M. R., Crispo E., Buraei Z., Strahs D., Isaacson M., et al., (2017), Participation in a year-long CURE embedded into major core genetics and cellular and molecular biology laboratory courses results in gains in foundational biological concepts and experimental design skills by novice undergraduate researchers. *J. Microbiol. Biol. Educ.*, **18**(1).
- Reisner B. A. and Stewart J. L., (2020), The Literature Discussion: A Signature Pedagogy for Chemistry, in *Advances in Teaching Inorganic Chemistry Volume 1: Classroom Innovations and Faculty Development*, ACS Publications, pp. 3–20.
- Remich R., Naffziger-Hirsch M. E., Gazley J. L., and McGee R., (2016), Scientific growth and identity development during a postbaccalaureate program: Results from a multisite qualitative study. *CBE—Life Sci. Educ.*, **15**(3), ar25.
- Robnett R. D., Chemers M. M., and Zurbriggen E. L., (2015), Longitudinal associations among undergraduates' research experience, self-efficacy, and identity. *J. Res. Sci. Teach.*, **52**(6), 847–867.
- Rodenbusch S. E., Hernandez P. R., Simmons S. L., and Dolan E. L., (2016), Early engagement in course-based research increases graduation rates and completion of science, engineering, and mathematics degrees. *CBE—Life Sci. Educ.*, **15**(2), ar20.
- Rodriguez J.-M. G., Bain K., Towns M. H., Elmgren M., and Ho F. M., (2019), Covariational reasoning and mathematical narratives: investigating students' understanding of graphs in chemical kinetics. *Chem. Educ. Res. Pract.*, **20**(1), 107–119.
- Ryoo K. and Linn M. C., (2012), Can dynamic visualizations improve middle school students' understanding of energy in photosynthesis? *J. Res. Sci. Teach.*, **49**(2), 218–243.
- Sadler T. D., Burgin S., McKinney L., and Ponjuan L., (2010), Learning science through research apprenticeships: A critical review of the literature. *J. Res. Sci. Teach. Off. J. Natl. Assoc. Res. Sci. Teach.*, **47**(3), 235–256.
- Sato B. K., Kadandale P., He W., Murata P. M. N., Latif Y., and Warschauer M., (2014), Practice Makes Pretty Good: Assessment of Primary Literature Reading Abilities across Multiple Large-Enrollment Biology Laboratory Courses. *CBE—Life Sci. Educ.*, **13**(4), 677–686.
- Schultz P. W., Hernandez P. R., Woodcock A., Estrada M., Chance R. C., Aguilar M., and Serpe R. T., (2011), Patching the pipeline: Reducing educational disparities in the sciences through minority training programs. *Educ. Eval. Policy Anal.*, **33**(1), 95–114.
- Seymour E. and Hunter A.-B., (2019), *Talking about leaving revisited*, Springer.
- Shanks R. A., Robertson C. L., Haygood C. S., Herdliksa A. M., Herdliksa H. R., and Lloyd S. A., (2017), Measuring and Advancing Experimental Design Ability in an Introductory Course without Altering Existing Lab Curriculum. *J. Microbiol. Biol. Educ.*, **18**(1), 1–8.
- Shultz G. V. and Gere A. R., (2015), Writing-to-learn the nature of science in the context of the Lewis dot structure model. *J. Chem. Educ.*, **92**(8), 1325–1329.
- Shuster, M. I., Curtiss, J., Wright, T. F., Champion, C., Sharifi, M., & Bosland, J. (2019), Implementing and evaluating a course-based undergraduate research experience (CURE) at a Hispanic-Serving institution. *Interdiscip. J. Probl.*, **13**(2), 1.
- Sirum K. and Humburg J., (2011), The Experimental Design Ability Test (EDAT). *Bioscene J. Coll. Biol. Teach.*, **37**(1), 8–16.
- Stone E. M., (2014), Guiding Students to Develop an Understanding of Scientific Inquiry: A Science Skills Approach to Instruction and Assessment. *Cell Biol. Educ.*, **13**(1), 90–101.
- Stone K. L., Kissel D. S., Shaner S. E., Grice K. A., and van Opstal M. T., (2020), Forming a Community of Practice to Support Faculty in Implementing Course-Based Undergraduate Research Experiences, in *Advances in Teaching Inorganic Chemistry Volume 2: Laboratory Enrichment and Faculty Community*, ACS Publications, pp. 35–55.
- Szteinberg G. A. and Weaver G. C., (2013), Participants' reflections two and three years after an introductory chemistry course-embedded research experience. *Chem. Educ. Res. Pract.*, **14**(1), 23–35.

- 1  
2  
3 Taber K. S., (2018), The use of Cronbach's alpha when developing and reporting research instruments in  
4 science education. *Res. Sci. Educ.*, **48**(6), 1273–1296.
- 5 Timmerman B. E. C., Strickland D. C., Johnson R. L., and Payne J. R., (2011), Development of a  
6 'universal' rubric for assessing undergraduates' scientific reasoning skills using scientific writing.  
7 *Assess. Eval. High. Educ.*, **36**(5), 509–547.
- 8 Varela M. F., Lutnesky M. M. F., and Osgood M. P., (2005), Assessment of Student Skills for Critiquing  
9 Published Primary Scientific Literature Using a Primary Trait Analysis Scale. *Microbiol. Educ.*,  
10 **6**, 20–27.
- 11 Watts F. M., Spencer J. L., and Shultz G. V., (2020), Writing Assignments to Support the Learning Goals  
12 of a CURE. *J. Chem. Educ.*, **98**(2), 510-514.
- 13 White B., Stains M., Escriu-Sune M., Medaglia E., Rostamnjad L., Chinn C., and Sevian H., (2011), A  
14 Novel Instrument for Assessing Students' Critical Thinking Abilities. *J. Coll. Sci. Teach.*, **40**,  
15 102–107.
- 16 White H. B., Benore M. A., Sumter T. F., Caldwell B. D., and Bell E., (2013), What skills should students  
17 of undergraduate biochemistry and molecular biology programs have upon graduation? *Biochem.*  
18 *Mol. Biol. Educ.*, **41**(5), 297–301.
- 19 White R. and Gunstone R., (2014), *Probing understanding*, Routledge.
- 20 Wiggins G., (1998), *Educative Assessment. Designing Assessments To Inform and Improve Student*  
21 *Performance.*, San Francisco, CA: Jossey-Bass Publishers.
- 22 Williams L. C. and Reddish M. J., (2018), Integrating primary research into the teaching lab: benefits and  
23 impacts of a one-semester CURE for physical chemistry. *J. Chem. Educ.*, **95**(6), 928–938.
- 24 Wilson M. Constructing measures: An item response modeling approach. Mahwah, NJ, US: Lawrence  
25 Erlbaum Associates Publishers; 2005.
- 26 Wright, B. D., & Masters, G. N. (1982). *Rating scale analysis*. MESA press.
- 27  
28  
29  
30  
31  
32  
33  
34  
35  
36  
37  
38  
39  
40  
41  
42  
43  
44  
45  
46  
47  
48  
49  
50  
51  
52  
53  
54  
55  
56  
57  
58  
59  
60