



**Application and testing of a framework for characterizing
the quality of scientific reasoning in chemistry students'
writing on ocean acidification**

Journal:	<i>Chemistry Education Research and Practice</i>
Manuscript ID	RP-ART-01-2019-000005.R1
Article Type:	Paper
Date Submitted by the Author:	15-Mar-2019
Complete List of Authors:	Moon, Alena; University of Michigan, Chemistry Moeller, Robert; University of Michigan, Department of Chemistry Gere, Anne; University of Michigan, Sweetland Center for Writing Shultz, Ginger; University of Michigan, Department of Chemistry

1
2
3 **Application and testing of a framework for characterizing the quality**
4 **of scientific reasoning in chemistry students' writing on ocean**
5 **acidification**
6
7
8
9

10 Alena Moon, Robert Moeller, Anne Ruggles Gere†, and Ginger V. Shultz

11 *Department of Chemistry, University of Michigan, Ann Arbor, MI, USA*

12 *†Sweetland Center for Writing, University of Michigan, Ann Arbor, MI, USA*

13
14
15 Corresponding author: Ginger V. Shultz, gshultz@umich.edu, 2521 Chemistry,
16
17 University of Michigan, Ann Arbor, MI, 48109.
18
19

20 Provide short biographical notes on all contributors here if the journal requires them.
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

Development and testing of a framework for characterizing the quality of scientific reasoning in students' writing on ocean acidification

Science educators recognize the need to teach scientific ways of knowing and reasoning in addition to scientific knowledge. However, characterizing and assessing scientific ways of knowing and reasoning is challenging. Writing-to-learn offers one way of eliciting and supporting students' reasoning; further, writing serves to externalize and make traceable students' reasoning. For this reason, it is a useful formative assessment of scientific reasoning. The utility hinges on researchers' ability to understand what students can do and think from their writing. Given the challenges in assessing students' writing, this research offers an adapted framework for assessing students' scientific reasoning evident in writing. This work will introduce the adapted framework and show an application to general chemistry students' argumentative writing about ocean acidification. We provide evidence that this framework can be used to validly estimate the quality of students' reasoning. We argue that this framework offers some affordances that overcome challenges reported in the literature. It serves to define scientific reasoning in a domain-general way by breaking it down into its components, but in a way that can produce a composite score that tells us about how students reason using chemistry content. Further, the framework provides a way to characterize the scientific accuracy of students' reasoning that can inform instructors' treatment of alternative conceptions.

Keywords: Writing-to-learn; assessing writing; scientific reasoning

Background

Science educators recognize that it is insufficient to only teach students' scientific knowledge as a collection of concepts and topics. Rather, to enable students to use scientific knowledge, we must support the development of reasoning and thinking skills that scientists use (NRC, 2012; Sevian & Talanquer, 2014; Bulte, Westbroek, De Jong, & Pilot, 2006). Writing-to-learn (WTL) is one way of supporting the development of this skill by activating deep thinking and reasoning in students (Keys, 1999) and, more importantly, making that reasoning visible and traceable (Emig, 1977; Kelly & Takao,

1
2
3 2002; Kelly, Regev, & Prothero, 2007). From an assessment perspective, this evidence
4 of student reasoning is valuable in so far as researchers and practitioners can use it to
5 make an argument about students' abilities to reason scientifically (Laverty et al., 2016;
6 NRC, 2001). However, there are challenges that currently limit the utility of this
7 evidence. There are few widely agreed upon epistemic criteria for characterizing the
8 quality of students' reasoning (i.e., what makes one students' reasoning better than
9 another's). Further, actually applying these criteria to understand and evaluate students'
10 writing is difficult as writing requires the researcher to make choices about grain size,
11 whether to evaluate structure or content or both, and what the presence or absence of a
12 quality criterion actually looks like in students' writing (Kelly & Takao, 2002; Takao &
13 Kelly, 2003a). To address these challenges, we have modified and applied a framework
14 for characterizing and evaluating reasoning in students' argumentative writing. This
15 framework contributes meaningfully to efforts to conceptualize and evaluate scientific
16 reasoning, as well as to efforts to analyse writing, which poses unique challenges.

35 *Writing to Learn*

37 Writing-to-learn refers to the kind of informal writing about science that facilitates
38 learning and ownership of scientific ideas. This informal writing is distinct in that its
39 primary aim is not to communicate or display mastery to an instructor, but to actually
40 facilitate sense-making by activating deep thinking and interaction with the concepts
41 (Keys, 1994). A secondary benefit of writing-to-learn, then, is promoting engagement
42 with disciplinary norms of writing and thinking (Prain & Hand, 2016). There is quite a
43 bit of variation around this primary aim, however; WTL assignments take a variety of
44 forms, lengths, methods of text production, audiences, and genres (Keys, 1994).

56 A secondary analysis of six writing-to-learn studies revealed some promising
57 gains as a result of writing-to-learn—the treatment condition outperformed comparison
58
59
60

1
2
3 groups on a total test scores and conceptual question scores and this effect was largely
4
5 due to the treatment (Gunel, Hand, & Prain, 2007). All six studies followed a similar
6
7 design including a pre-test/post-test design with the test having multiple-choice and
8
9 conceptual extended response questions. More importantly, all writing interventions
10
11 were grounded in the same theoretical considerations that have been identified as key
12
13 for successful learning from writing: 1) opportunities for brainstorming, 2) provision of
14
15 authentic audiences, 3) drafting and redrafting with feedback, 4) explicit instruction of
16
17 genre specifications, 5) focus on big ideas, 6) use of rubrics, and 7) diverse
18
19 opportunities to plan and draft writing (Gere, Limlamai, Wilson, MacDougall Saylor, &
20
21 Pugh, 2019; Gunel et al., 2007; Klein, 1999, 2015). The theoretical grounding afforded
22
23 comparisons across domains and writing assignment types and served to reveal the
24
25 benefits of WTL more broadly (Gunel et al., 2007; Prain & Hand, 2016). However, at
26
27 the undergraduate STEM level specifically, more work is needed to understand the
28
29 mechanism of effect for WTL assignments (Reynolds, Thaiss, Katkin, & Thompson,
30
31 2012) and we argue that to undertake investigations into the mechanism of effect, we
32
33 need a reliable and meaningful framework for interpreting and evaluating students'
34
35 written work.
36
37
38
39
40
41
42
43

Characterizing Students' Reasoning in Written Products

44
45 Constructed responses reveal rich insight into students' ideas and the coherence of those
46
47 ideas, but evaluating open responses remains a barrier to implementing such rich
48
49 assessments (Liu, Rios, Heilman, Gerard, & Linn, 2016). This barrier consists of two
50
51 distinct but interdependent challenges: characterizing the quality usually in some sort of
52
53 rubric (Kelly & Bazerman, 2003; Sandoval, 2003; Sandoval & Millwood, 2005) and
54
55 consistently and reliably applying the quality criteria (Ha, Nehm, Urban-Lurain, &
56
57 Merrill, 2011; Liu et al., 2016). Additionally, the nature and difficulty of these
58
59
60

1
2
3 challenges vary with the length of constructed response; for example, extended
4
5 arguments in the form of research reports in oceanography tended to have long and
6
7 complex chains of reasoning that are difficult to characterize with a single rubric (Kelly,
8
9 Regev, & Prothero, 2007; Kelly & Takao, 2002). Researchers have sought to overcome
10
11 these challenges with a variety of approaches for a variety of written products, ranging
12
13 from short written explanations (Ha et al., 2011; Liu et al., 2016; Moreira, Marzabal, &
14
15 Talanquer, 2018; Sandoval, 2003; Sandoval & Millwood, 2005; Sandoval & Reiser,
16
17 2004) to more extensive writing like laboratory reports or research reports (Grimberg &
18
19 Hand, 2009; Kelly, Chen, & Prothero, 2000; Kelly et al., 2007; Kelly & Takao, 2002;
20
21 Takao & Kelly, 2003b). A few of these approaches are distinct in that they break down
22
23 students' responses into smaller units to then identify patterns in students' reasoning, as
24
25 opposed to a rubric that considers the quality of the response as a whole (Grimberg &
26
27 Hand, 2009; Kelly et al., 2007; Kelly & Takao, 2002; Moreira et al., 2018). These
28
29 approaches will be the focus of this literature review.
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

Moreira et al. (2018) specifically sought to characterize the causal reasoning of 10th grade chemistry students' explanations of freezing point depression. To do so, the authors modified and applied a discourse analysis framework developed by Russ et al. (2008) to students' written explanations and drawings. The final form of the analysis scheme included four components—entities, properties, activities, and organisation—and the relationships between the components and the students' drawings that were identified in students' responses. Entities are the 'things' in the system that are being considered. Properties are characteristics of those entities and activities are actions of those entities. Organisation refers to the spatial-temporal relationship between the entities and activities or properties of the system. By coding the explanations for these components and relationships, they were able to elucidate patterns in students'

1
2
3 explanations and organize these patterns into four levels according to the quality of
4 causal reasoning. These levels increased in sophistication from descriptive to relational
5 to simple causal, culminating in emerging mechanistic. Unsurprisingly given the
6 authors' previous findings (Sevian & Talanquer, 2014), the majority of students (45%)
7 used relational causal reasoning (Moreira et al., 2018). Explanations at this level could
8 be modelled to show that students generally identified two entities and the properties of
9 one or both entities, and then related either the properties to each other or related
10 entities to properties. Such complex modelling of students' explanations can hopefully
11 equip teachers with more sophisticated approaches to understanding, interpreting, and
12 developing students' reasoning abilities (Moreira et al., 2018).

13
14
15
16
17
18
19
20
21
22
23
24
25
26 Grimberg and Hand (2009) similarly identified the presence or absence of
27 dimensions of reasoning and determined *patterns* in students' reasoning evident in their
28 laboratory reports. In this study, Grimberg and Hand (2009) identify cognitive
29 operations used in writing laboratory reports and then construct what they term
30 'cognitive pathways'—the sequence of cognitive operations used by author(s) of a lab
31 report. The authors argued that because writing a laboratory report was a meaning-
32 making activity, considering the sequence of cognitive operations revealed how students
33 were constructing meaning. Using a list of 11 cognitive operations derived partially
34 from the literature and from the students' data, authors coded students' writing for use
35 of cognitive operations. The cognitive operations included observation, measurement,
36 comparison, analogy, clarifications, claim, cause/effect, induction/generalization,
37 deduction, investigation design, and argumentation (Grimberg & Hand, 2009). When
38 comparing cognitive pathways of low achievers to high achievers, as determined by a
39 standardized skills test, both high and low achievers used the same range of operations,
40 but with a different structure. Though the cognitive structure was partially determined
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60

1
2
3 by the structure of the Science Writing Heuristic (SWH) activity (i.e., clarification
4 questions were posed during the research question portion of SWH activity), high
5
6
7
8 achievers began using complex operations earlier in the text than low achievers. This
9
10 ultimately demonstrated that SWH scaffolding supported all students in using high-
11
12 complexity operations, albeit at different rates (Grimberg & Hand, 2009).
13

14
15 Grimberg and Hand's (2009) work demonstrates the capacity of writing to make
16
17 students' thinking visible and traceable. Emig (1977) argues that writing is unique in its
18
19 capacity to do this. Kelly and Takao (2002) demonstrate the utility of students' writing
20
21 for understanding their reasoning by characterizing how undergraduate students use
22
23 evidence to construct arguments. To make this characterization, they developed a
24
25 research methodology that models epistemic levels of argument. This framework
26
27 includes six levels ranging from the lowest, data charts and representations, to the
28
29 highest, general geological (the specific context for this work) knowledge not specific to
30
31 the data presented. The levels represented students' ability to abstract from data to make
32
33 claims. With this framework, the authors analysed a subset of undergraduate
34
35 oceanography students' assignments by labelling each sentence with an epistemic level.
36
37 These epistemic criteria were then weighted in order to rank the 24 student arguments
38
39 (research reports in oceanography) from best to worst. While this framework was very
40
41 useful for characterizing the quality of students' arguments based on their use of
42
43 evidence, it did possess limitations. Namely, the assessment of quality determined by
44
45 the framework did not always align with content experts' evaluation of quality, the
46
47 framework did not consider inference logic (how the data led to theoretical claims), and
48
49 the authors had to make inferences in their application of the framework. These
50
51 limitations are difficult to overcome for everyone aiming to evaluate ill-defined
52
53 constructs like use of evidence. However, Kelly and Takao (2002) revealed that claims
54
55
56
57
58
59
60

1
2
3 can be made about students' reasoning from their written work. In order for the insight
4
5 into students' reasoning provided by writing to be useful for informing students'
6
7 development, tools for assessing writing must be efficient, systematic, and offer tailored
8
9 feedback.
10

11
12 Ongoing discussions about writing assessment distinguish between holistic
13
14 scoring, assigning a single score to a broad variable like writing proficiency, and
15
16 analytic scoring, breaking down variables like writing proficiency into components that
17
18 are individually scored (Hamp-Lyons, 2016a, b). High-stakes assessments, such as
19
20 college entrance examinations, motivated the use of both holistic and analytic scoring
21
22 approaches to assign students general writing scores, but researchers have begun to
23
24 identify shortcomings of both (Neill, 2002; Hamp-Lyons, 2016b, Chapman, 2016). With
25
26 more complex analytical tools (i.e., multivariate analyses), Hamp-Lyons (2016) calls for
27
28 a movement to multiple trait scoring of writing. In multiple trait scoring, there is no
29
30 single score given, whether composite or holistic. Rather, a set of scores is assigned
31
32 with multiple traits each warranting a score, thus lending to a richer description of
33
34 students' ability (Hamp-Lyons, 2016a, b).
35
36
37
38
39
40

41 ***Rationale and Research Objectives***

42
43 In any effort to measure a student's reasoning, choices must be made about what
44
45 characterizes quality. Specific to evaluating extensive writing, additional decisions must
46
47 be made about the grain size, level, and nature of rubric that will be used to determine
48
49 quality. In order to address these challenges, the cognitive operations used by Grimberg
50
51 and Hand (2009) were modified and applied to students' writing on ocean acidification.
52
53 Motivated to leverage the rich insight into student thinking that constructed responses
54
55 offer, the work presented herein aimed to test an approach for analysing extensive
56
57
58
59
60

1
2
3 scientific writing by applying it to a new context. We aimed to answer the following
4
5 questions regarding this approach:
6

- 7
8 (1) Can cognitive operations be used to make sense of general chemistry students'
9 argumentative writing? If so, how?
10
11
12 (2) What features of students' argumentative writing do cognitive operations serve
13 to explain?
14
15
16 (3) What is the relationship between framework estimates of quality and conceptual
17 correctness?
18
19
20
21
22

23 **Methods**

24 *Participants, setting, and data collection*

25
26 A writing prompt was designed and administered in a first-semester General Chemistry
27 course serving primarily students in the College of Engineering and undeclared students
28 in the College of Literature, Science, and Arts. This course had an enrolment of 1413
29 students, most of whom were freshman and sophomores. The content of the course
30 covered traditional general chemistry concepts, ranging from dimensional analysis,
31 quantum mechanical atomic models, bonding theories, to reactions, enthalpy,
32 intermolecular forces, chemical equilibrium, and acid-base theories.
33
34
35
36
37
38
39
40
41
42
43

44 This course is structured with three lectures led by an instructor and one
45 discussion session led by a teaching assistant per week. During each discussion section,
46 students complete a quiz. The writing assignment for this study was administered as a
47 substitute for a quiz. The writing assignment was uploaded as a .pdf file to the course
48 management site one week before the due date. This writing assignment directed
49 students to consider a set of concepts in their response and to keep their post between
50
51
52
53
54
55
56
57
58
59
60

1
2
3 350 and 500 words. Though the majority of students' responses were within this range,
4
5 some wrote less and some wrote more.
6

7
8 Students all submitted their writing assignment online the following week at the
9
10 start time of their specific discussion session. For this reason, some students with
11
12 discussion sessions later in the week had more time to write than those with earlier
13
14 recitations. During the discussion, students formed teams of two or three, switched
15
16 papers and gave feedback to each other. Six hundred seventy-three students gave
17
18 consent to have their writing analysed. Ethical review board approval was gained in
19
20 order to collect and analyse written assignments that students consented to have
21
22 analysed. Students did not receive any feedback on their written work beyond the
23
24 conversation that took place in their discussion session. We found that students did not
25
26 make meaningful revisions following the peer review discussion. Additionally, they
27
28 were not required to submit a revision. However, if they did, that revised draft was used
29
30 for analysis.
31
32
33
34
35

36 ***Writing activity development and design***

37
38 The writing prompt was developed iteratively through correspondence with authors who
39
40 had expertise in writing to learn and the development of meaningful writing prompts
41
42 (AG) and the faculty members teaching the course who collectively held more than two
43
44 decades of experience teaching chemical equilibrium in general chemistry. WTL
45
46 prompts are generally designed to provide students with an audience, an identity, and an
47
48 authentic context that require students to engage with a specific concept. This WTL
49
50 assignment was intentionally designed with elements empirically determined to
51
52 contribute to meaningful learning through writing (Gere et al., 2019). In this case, the
53
54 prompt showed a fake social media post, in which 'Ernie Clueless' shares a plot
55
56 illustrating the trend of concentration of atmospheric carbon dioxide and ocean pH over
57
58
59
60

1
2
3 time. Ernie claims that these things are unrelated. Students were tasked with explaining
4 the relationship to Ernie, given the relevant equilibria. The prompt targeted the concept
5 of chemical equilibrium, drawing on Le Châtelier's principle. It inherently supported
6 argumentation by requiring students to differentiate their perspective from Ernie's.
7
8
9
10
11
12

13
14 ***Data analysis: Development and application of analytical framework***

15 A list of cognitive operations was modified from a list used by Grimberg and Hand
16 (2009) to analyse reports from a Science Writing Heuristic (SWH) laboratory. In this
17 context, a cognitive operation is a written discursive move that serves some cognitive
18 objective. Cognitive operations, then, determined the grain size for breaking down an
19 essay into smaller analysable units (i.e. the number of sentences that served the
20 objective of a claim, for example, were coded as such). For this reason, a claim could be
21 one sentence in one essay and three sentences in another. The amount of text that was
22 assigned a code was determined by the function of that text. The list of cognitive
23 operations used by Grimberg and Hand (2009) was refined iteratively by testing it
24 against the data. This involved using Grimberg and Hand's original set to code the text,
25 identifying text that could not be coded with this set and operations from this set that
26 did not serve to explain any of the data, and refining the set of operations to a set that
27 were all used to describe virtually all of the text. Multiple initial iterations occurred with
28 the same subset of essays (N=25) and subsequent iterations incorporated more essays on
29 an as needed basis. Table 1 shows the final list of cognitive operations that was used
30 throughout the final analysis. Included in Table 1 is a characterization of the
31 dimensionality of each operation. As will be explained further in the theoretical
32 framework, the complexity of cognitive operations was determined by its
33 dimensionality—number of ideas being drawn upon. Operations with two domains were
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

more cognitively complex than operations with one dimension. That is, they drew upon and connected more idea units.

Table 1. Finalized list of cognitive operations and descriptors used to analyse all essays, listed in order of increasing cognitive complexity

Cognitive Operation	Description	Dimensionality
1. Definition	Canonical description of a term, concept, idea, or theory	single domain
2. Observation	Qualitative description of change, trend, or transformation of a variable	
3. Measurement	Quantitative description of change, trend, or transformation of a variable	
4. Comparison	Relationship of change, trend, or transformation for two or more variables	
5. Example	Illustration of a class of objects by singling out one object	two domains
6. Claim	Assertion supported with a tentative explanation	
7. Consequences	Cause and effect explanation with either cause or effect falling outside the scope of the writing prompt	
8. Cause and effect	Explanation providing a mechanism with a causal agent and observed effects	multiple domains
9. Deduction	Application of a theory or principle to a specific system or scenario	
10. Argumentation	Explicit differentiation between the author's perspective and the fictional character's perspective*	

*This conceptualization of argumentation was specific to the context of writing prompt used in this study. It is expected that it could be easily translated to other writing contexts that require argumentation.

Once a list of cognitive operations was finalized, the first two authors began coding assignments and built a detailed rubric that included definitions, linguistic markers, and examples for each operation. This rubric was further refined through multiple iterations of analysis by a team of chemistry education researchers in an effort to establish inter-rater reliability (IRR). This stage included a team consisting of the first two authors and another chemistry education researcher trained in qualitative coding of writing. The graduate student was trained on an existing rubric, the whole team coded ten assignments, and an IRR coefficient in the form of Krippendorff's alpha (KA) was calculated. This coefficient was quite low after the first round, so revisions were made to the rubric, another training session was conducted, and subsequent rounds of analysis

1
2
3 were conducted until a Krippendorff's alpha value of 0.69, the minimum acceptable
4 value, was achieved. Training involved discussing the rubric which included examples
5 from many essays and then illustrating the coding process by coding a few essays all
6 together.
7
8
9
10
11

12 The first two authors then coded approximately 200 assignments each, with
13 large overlap. Having two researchers code many of the same assignments lent to the
14 reliability of the coding. Throughout analysis, IRR 'checks' were performed to ensure
15 that the rubric was being applied consistently. This involved selecting overlapping
16 assignments and determining a KA. This stage resulted in a KA of 0.89, which was a
17 desirable value (Krippendorff, 2004).
18
19
20
21
22
23
24
25
26

27 *Estimate of quality*

28
29 Once all assignments were coded for cognitive operations, estimates of quality were
30 assigned according to the cognitive complexity of the essay. The cognitive operations
31 are ordered in Table 1 according to increasing cognitive complexity; that is, a definition
32 has the lowest cognitive complexity (1) whereas argumentation has the highest
33 cognitive complexity (10). Overall cognitive complexity for the essay was determined
34 by taking a weighted average of operations used. Because the magnitude of text
35 included in an operation varied from one essay to another, the average was weighted by
36 the number of sentences within that operation. An assignment, then, that was 50%
37 argumentation would likely have a higher average complexity than an assignment with
38 only 10% argumentation. This approach resulted in a single number characterizing the
39 quality of students' reasoning in an essay. This process of producing a single number is
40 illustrated with the example essay in Table 2.
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60

Table 2. Example of student essay and determination of cognitive complexity score

<p>Le Chatlier's principle is a way to show how if the equilibrium of a certain reaction is altered by changing certain aspects, then the reaction or equilibrium position will actually shift to fix the change. In Ernie's post there are multiple reactions that follow one after the other, ultimately showing how CO₂ transforms to H₂CO₃, which changes to H⁺ and HCO₃, which finally changes to 2H⁺. These reactions are connected in that as one breaks down, it spurs another reaction to then form. This shows how CO₂ in the atmosphere actually affects PH concentration and ocean acidification (or the number/concentration of H⁺ ions in the reaction). Because of this correlation, certain components like temperature, pressure, volume, and concentration can affect CO₂ levels and acidification. In our case, we will look at the concentration or amount of CO₂ in correlation to the amount of H⁺ that is formed. If the concentration of CO₂ (as a reactant) were to increase, the system would try to decrease it; this would mean that the concentration of the other reactants will increase to react with the CO₂, but then more product would be formed. Thus increasing the concentration of CO₂ (atmos) would cause the system to shift towards the products, producing more CO₂(aq). Then because there is more CO₂ (aq), more water would have to be use so more H₂CO₃ would form. The next reaction would then proceed as the others with an increase in the formation of the products H⁺ and HCO₃, which in turn would once again increase the next reaction producing more 2H⁺ product. This may seem confusing, but to summarize if CO₂ as a reactant increases, it would need to counteract this change by producing more product of 2H⁺ (acidification). Now, if CO₂ (atmos) reactant were to decrease, the equilibrium will actually try to increase it so as to set the system back at equilibrium. This would mean that a decrease in that reactant would cause more product to be formed so to produce more CO₂; thus 2H⁺ product would increase to make more CO₂.</p> <p>The relationship between CO₂ and pH is that as CO₂ increases, the pH will decrease and make the oceans more acidic. The pH scale runs from 1 to 14, with the acids being numbers 1 to 6. So the lower the pH, the more acidic and the more H⁺ ions that are formed (the H⁺ are indicators of acidic properties). Le Chatlier's principle shows that increasing one will have to increase the other so that the system is able to be equal again, thus there is correlation between CO₂ (atmos) increase and H⁺ ions increase(acidification)."</p>	<p>Definition (1)</p> <p>Observation (2)</p> <p>Claim (6)</p> <p>Deduction (9)</p> <p>Comparison (4) Definition (1) Comparison (4)</p> <p>Cause and Effect (8)</p>
<p>Cognitive complexity = $\frac{\sum(\text{cognitive operation score} * \# \text{ sentences used})}{\text{total number of sentences}}$</p> <p>Definition (1) + 3 * Observation (2) + Claim (6) + 9 * Deduction (9) + Comparison (4) + Definition (1) + Comparison (4) + Cause & Effect (8) = 111</p> <p>111/18 =6.2 (Average complexity for this essay)</p>	

Faculty ranking of essays

One motivation of this work is to use a framework to systematically characterize students' writing. In order to understand how this framework could accomplish that task, we compared framework estimates of quality with instructors' estimates of quality. The precedent for this approach is offered by Kelly and Takao's testing of the framework for epistemic levels (Kelly & Takao, 2002). Further, evidence from interviews with STEM faculty about writing suggested to us that the knowledge experts use to evaluate writing is tacit (Moon, Gere, & Shultz, 2018). By comparing expert rankings with essays ordered according to a framework, we can identify ways that the framework may capture this tacit knowledge. To compare framework estimates of quality and instructor estimates, we compared essays ranked according to the framework to instructor rankings. To select the essays, we split cognitive complexity into four roughly equivalent ranges (3.1-4.8, 4.8-6.5, 6.5-8.2, 8.2-10), where the highest range included the 'best' essays, according to the framework. Within each range, an essay was randomly selected. The output of this step was four essays ranging from most complex to least complex, as determined by the framework. These essays were then provided to instructors who were tasked with ranking the essays from best to worst according to the quality of scientific reasoning, which they were directed to evaluate as they saw fit. The rationale for not providing faculty with more extensive quality criteria was to target the kind of evaluative work that is inherent to grading this sort of task; that is, we wanted instructors to make the kinds of decisions that are required to define and evaluate scientific reasoning. Five faculty with a range of experience teaching general chemistry from different institutions ranked the writing tasks and one instructor volunteered the reasoning behind their ranking.

Chemistry Content Analysis

1
2
3 Essays were examined to understand how students employed ideas about
4
5 chemical equilibrium within their argument. Essays were coded for conceptual
6
7 correctness, which involved flagging all occurrences of inaccurate ideas. Every time an
8
9 inaccurate idea was identified, the cognitive operation containing that incorrect idea was
10
11 marked. The pairing of the incorrect idea with the cognitive operation was intentional,
12
13 based on the theoretical framework explained below, which posits that what students
14
15 articulate in written forms are representations of the ideas they hold. In this case those
16
17 ideas related to chemical equilibrium.
18
19
20
21
22

23 *Theoretical framework*

24
25 The primary assumption made in this work is that writing reveals students'
26
27 cognitive structures or understanding of the meaning of a concept (Emig, 1977; Novak,
28
29 2002). Meaning is defined in this case as 'the totality of propositions linked to any
30
31 given concept,' excluding the emotional association with the concept and the context in
32
33 which the concept was learned (Novak, 2002). This assumption is grounded in the
34
35 capacity of writing to 1) connect or relate propositions in the author's mind and 2) make
36
37 'evolutionary development of thought graphically visible and available (Emig, 1977).'

38
39 According to Novak's theory of meaningful learning, the complexity of the meanings
40
41 can be evaluated, the quantity and quality of which will determine meaningful learning
42
43
44
45
46 (Novak, 2002).
47

48
49 In this work, units of written text were coded according to their cognitive
50
51 function as a cognitive operation. Each cognitive operation was constituted by some
52
53 number of ontological domains, elements within those domains, and relationships
54
55 between each (Grimberg & Hand, 2009; Halford & McCredden, 1998). Cognitive
56
57 complexity, then, was defined in terms of dimensionality; in which cognitive operations
58
59 with higher numbers of domains, elements, and relationships were more complex than
60

1
2
3 cognitive operations with fewer (Halford & McCredden, 1998). These domains and
4
5 elements can be conceptualized similarly to the discourse analysis framework used by
6
7 Moreira et al. (2018) where they refer to ‘things’ in the system being considered. All
8
9 cognitive operations were organized on a spectrum from least complex to most complex
10
11 according to this criterion, illustrated in Table 1. A students’ scientific reasoning can
12
13 then be considered as the progression of operations used in a text.
14
15

16
17 Therefore, one way of evaluating the quality of meanings was through cognitive
18
19 complexity. It is possible, however, that a student can employ cognitively complex
20
21 reasoning without necessarily using correct content (Sandoval & Millwood, 2005; Kelly
22
23 & Takao, 2002). For this reason, a second way of evaluating the quality of meanings
24
25 was considering their conceptual correctness; that is, their agreement with scientifically
26
27 accepted knowledge. A conceptual change perspective suggests that problems of
28
29 incorrect conceptions arise from Limited or Inappropriate Propositional Hierarchies
30
31 (LIPs), the way that concepts are inappropriately organized in the learner’s mind,
32
33 which means that we as instructors must consider both the *content and structure* of
34
35 incorrect conceptions. Further, this implies a greater instructional effort is needed to
36
37 remediate stable LIPs (Novak, 2002). The cognitive operations framework used in this
38
39 study helps characterize the complexity or stability—as a function of the number of
40
41 domains, elements, and relationships—of conceptions.
42
43
44
45
46
47

48 **Results**

49 ***Research Question 1: Can cognitive operations be used to make sense of*** 50 ***general chemistry students’ argumentative writing? If so, how?***

51
52 Table 2 illustrates how these cognitive operations were interpreted to analyse the
53
54 writing with examples from students’ essays. The examples are useful for discussing the
55
56 difficulties of applying this framework. The lower complexity operations (*definition*
57
58
59
60

through *example*) were relatively easy to identify. The difference between *observation* and *measurement* was essentially a difference between qualitative and quantitative, with *measurements* requiring some numerical component. Because students could sufficiently respond to this writing prompt with qualitative reasoning, *measurement* occurred less frequently (Table 1). A *comparison* was distinct from *observation* and *measurement* in describing more than one variable relative to each other. A *claim* was similar to *comparison* in referencing a relationship between two variables but was distinct in that it required a tentative explanation of the relationship. For the *claim*, then, there were two domains (explanation and observation). *Cause and effect* and *consequences* were similar and easier to identify in text, with common linguistic markers being ‘caused,’ ‘leads to,’ or ‘drives.’ *Consequences* used cause and effect reasoning but relied on causes or effects that fell outside the scope of the prompt. In this case, the student referenced the effect of ocean acidification on coral. A primary marker of *deduction* was the invocation of a principle or theory that was then applied to a specific system. In this context, students frequently invoked Le Chatalier’s principle or equilibrium. Finally, *argumentation* was undoubtedly the most difficult to identify as it often drew on multiple operations. So, we used this feature as an identifier. *Argumentation*, then, required indistinguishable use of multiple operations and to explicitly differentiate between the author’s perspective and Ernie’s (or any opposing position in another context). Table 1, 2, and 3 together can be used to apply this framework to other contexts.

Table 3. Examples of cognitive operations, arranged in order of increasing complexity. Numbers describe cognitive complexity ordering.

Operation	Example from student essay
Definition (1)	Equilibrium is a state where a reaction is occurring forwards and backwards at equal rates with no overall change. When a change occurs to the system, the reaction will shift in a direction to counteract this change.

1 2 3 4 5	Observation (2)	The trend in the graphs shows an increase in atmospheric carbon dioxide over time.
6 7 8	Measurement (3)	The ocean pH has dropped from 8.2 to 8.1 since the Industrial Revolution.
9 10 11 12	Comparison (4)	The plot that you shared illustrates that as the concentration of atmospheric increases over time, the pH of the ocean seawater decreases.
13 14 15	Example (5)	An example is hydrochloric acid, which has hydrogen ions attached and is an acid with a rather low pH level.
16 17 18 19	Claim (6)	As a matter of fact, Ernie, the correlation between CO ₂ levels in the atmosphere and the pH of the oceans makes sense according to chemistry.
20 21 22 23 24 25 26 27	Consequences (7)	When this happens, calcifying organisms will become weaker, such as coral. They will be significantly affected as a result and may be unable to live in the current environment in which they live. In addition to ocean acidification wreaking havoc on the environment, other factors such as climate change can do the same thing and increase the amount of damage that is done to it.
28 29 30 31	Cause and Effect (8)	In each equation, the amount of reactants increases, which drives the reaction forward, meaning the amount of products will increase until the reaction reaches equilibrium.
32 33 34 35 36 37 38 39		In accordance with Le Châtelier's Principle, increasing the amount (or the concentration) of atmospheric CO ₂ will shift this equation towards the dissolved CO ₂ in the ocean to make up for the increase in gas on the reactants (left) side. This dissolved CO ₂ is indicated as 'CO ₂ (aq)' in the equation, which stands for 'aqueous CO ₂ .' Consequently, the dissolved CO ₂ relates to the bicarbonate formation equation:
40 41 42		$\text{CO}_2 (aq) + \text{H}_2\text{O} (l) \rightleftharpoons \text{H}_2\text{CO}_3 (aq) \rightleftharpoons \text{H}^+ + \text{HCO}_3 (aq)$ (Doney et al., 2009 [<i>reference provided to student</i>])
43 44 45 46 47 48 49 50 51 52 53 54 55 56	Deduction (9)	Just as increasing the concentration of the atmospheric CO ₂ caused the equilibrium between dissolved CO ₂ to favor the formation of dissolved CO ₂ , a similar phenomenon will occur to favor production of bicarbonate (HCO ₃) and hydrogen, two byproducts of carbonic acid (H ₂ CO ₃). Increasing dissolved CO ₂ concentration—a result of increased atmospheric CO ₂ —will 'push' the equilibrium towards the formation of H ₂ CO ₃ . Likewise, a shift towards the formation of H ₂ CO ₃ will also shift the equation forward towards the formation of protons, H ⁺ , and HCO ₃ . Finally, these products of proton and bicarbonate will be in equilibrium with two protons and carbon trioxide (CO ₃ ²⁻):
57 58 59 60		$\text{H}^+ + \text{HCO}_3 (aq) \rightleftharpoons 2\text{H}^+ + \text{CO}_3^{2-}$ (Doney et al., 2009)

	Once more, an increase in the H^+ and HCO_3^- concentration will push the equilibrium forward towards the formation of two $2H^+$ and CO_3^{2-} .
Argumentation (10)	So using this information, you can now look at the graph you posted and understand the relationship between CO_2 levels and the pH of the water. As CO_2 is absorbed into the water, it produces H^+ ions, which then cause the pH of the water to decrease. You can see this trend on the graph. Even though they do not seem like they should be related in any way, a change in one would cause a change in the other. Something that is making the line representing the CO_2 in the atmosphere on the graph to increase so much is the amount of CO_2 humans emit every day. Whenever you drive a car you are releasing CO_2 into the atmosphere. Since there is so much more CO_2 being released, the oceans are absorbing more CO_2 . In fact, the oceans have absorbed almost 30% of the CO_2 humans have emitted since the Industrial Revolution. As more CO_2 is absorbed by the oceans because of human activity, the more H^+ ions are formed, and the more the pH of the ocean is decreased.

A total of 296 assignments have been coded (average weighted complexity: 6.3; average number of operations per essay: 9). We determined that saturation had been reached when no new codes or patterns were observed in the writing. Table 4 shows the frequency of operations used in the 296 essays. *Observation* was the most frequently used operation by students, with many using more than one observation in a single essay. Students heavily relied on making statements about how a variable was changing to counter 'Ernie's' claim. In this context, this means that students were able to understand how variables were changing from the graph provided. There was very little use of *Example* or *Measurement*. *Argumentation* was also used relatively infrequently. As mentioned above, *argumentation* was distinct as its own operation given indistinguishable use of multiple lower complexity operations. For this reason, *argumentation* required students to combine multiple operations (and domains). This difficulty combined with the infrequent use suggests that *argumentation* was indeed the most complex operation. The high frequency of *claim*, *cause and effect*, and *deduction* is likely tied to this writing context. Students were trying to convince Ernie (*claim* and *cause and effect*) by invoking chemical principles (*deduction*). It is expected that the

distribution of use will vary with different writing contexts, depending upon what the prompt elicits.

Table 4. Descriptive information from application of the framework to the data set presented herein.

Operations	Frequency
Observation (2)	411
Claim (6)	348
Cause and effect (8)	331
Definition (1)	312
Comparison (4)	302
Deduction (9)	235
Consequences (7)	146
Argumentation (10)	63
Example (5)	51
Measurement (3)	34

Research Question 2: What features of students' argumentative writing do cognitive operations serve to explain?

To determine what exactly this framework served to characterize, two comparisons were made. The first comparison was between framework estimates of complexity and instructor estimates of quality. This comparison was intended to demonstrate that this framework was telling us something that faculty would normally have to make a judgment about. Further this comparison was intended to reveal that this framework could make *similar* judgments to an instructor. Table 5 shows how five faculty ranked four assignments, and this is compared to our framework ranking. Instructors were tasked with ranking the assignments according to the quality of the scientific reasoning (as they saw fit to evaluate it). This approach was taken so as to elicit instructors' "gut reaction"—the kind of evaluation they would make if they were grading this sort of task in their class. Instructor rankings reveal a few trends. First and potentially the most

important, there is almost complete consensus on the ‘best’ essay (J) with the exception of Instructor 3. This finding speaks to the framework’s capacity to identify the best. Essay J received a high cognitive complexity score because of the presence of extended argumentation, which aligns with what instructors are valuing when evaluating scientific reasoning. Four of the five instructors ranked Essay G as second best, with the exception of Instructor 3, even though our framework estimates it as second from the worst. Finally, all five instructors consider Essay H and K to be the worst, where as our framework estimates Essay H to be the second best. The difference between instructor rankings for Essays H, G, and K and framework estimates illustrate an important limitation of the framework. Essay H contained a misconception in which the student claimed a relationship between atmospheric temperature and ocean acidification though there was no data regarding heat for any of the chemical reactions provided. However, this student’s reasoning was rather sophisticated, with multiple high complexity cognitive operations. Our framework does not account for scientific accuracy of students’ essays. We chose to include this essay as it authentically represents what an instructor might encounter with grading writing. One possible explanation for the ranking difference is that for instructors scientific content and reasoning are inextricably linked, which is consistent with feedback from one instructor who explained that content accuracy factored into their ranking. It is for this reason that we also analysed the scientific accuracy of students’ essays (see Table 5 below).

Table 5. Instructor ranking of four assignments of varying quality compared to framework estimates of ranking

Essay name	Cognitive complexity ranking (4 worst, 1 best)	Expert rankings (4 worst, 1 best)				
		Instructor 1	Instructor 2	Instructor 3	Instructor 4	Instructor 5
J	1	1	1	2	1	1
H	2	3	3	4	4	4
G	3	2	2	1	2	2

K 4 4 4 3 3 3

The second comparison made was between cognitive complexity—framework estimates of quality—and common student characteristics that are frequently used as measures or predictors of success (Hein & Smerdon, 2013). The purpose of this comparison was to determine if characterizing the quality of students' reasoning in this way was revealing something about students that could have been predicted by a metric that was already collected (e.g., ACT math score). In other words, this framework is useful only in so far as it tells us something interesting about students that other metrics do not. Table 6 shows the correlations between common student characteristics and cognitive complexity. There were no significant correlations between cognitive complexity and any common characteristics, which would not be expected for measures of constructs distinct from that captured by this framework (i.e. math). These findings may mean that the framework captures something distinct from what is measured by other standardized tests. A strong negative correlation exists between the number of operations used and the cognitive complexity. This means that students with higher cognitive complexity essays used fewer moves, which could indicate a synthesis of ideas in order to produce higher complexity operations.

Table 6. Pearson correlations between student characteristics and cognitive complexity (p-values reported for t-tests used for categorical variables: gender and ethnicity [white and non-white students compared])

Variables	Cognitive complexity
Number of operations	-0.649*
Final exam grade	-0.018
Final course grade	-0.025
CHEM placement	-0.081
MATH placement	-0.020
ACT math [†]	0.003
Current GPA	-0.062
Cumulative GPA	-0.060
Gender	0.401
Ethnicity	0.071

* indicates p (two-tailed) < 0.01

† For students with only SAT math scores, their scores were converted to ACT math scores using contingency tables

Research Question 3: What is the relationship between framework characterizations of complexity and conceptual correctness?

The data above reveal that the cognitive operations framework is characterizing students' reasoning in a way that other measures do not. However, the instructor rankings reveal that there exists a relationship between reasoning and accuracy. The motivation for considering this relationship partially sources from the concern that any information this framework provides is irrelevant if students are largely scientifically inaccurate. In order to explore this relationship, we coded all data that had already been coded according to operations for 'correctness.' That is, when scientifically inaccurate information was identified in an essay, the cognitive operation containing that information was marked as incorrect. In this way, all student writing was coded for both correctness and cognitive function (i.e., content and structure as highlighted in theoretical framework). Table 7 shows the number of incorrect operations per the total number of cognitive operations. Further, there were no correlations between the cognitive complexity and number of incorrect operations or between the number of cognitive operations and the number of incorrect operations. This finding suggests that overall, producing a more complex essay does not make it more likely that a student will use more incorrect ideas, but as Table 7 shows, there may be specific operations that elicit more incorrect ideas. Further, writing more operations, or introducing more separate idea units, does not make a student more likely to put forth incorrect ideas.

Table 7. Number of incorrect cognitive operations relative to total number of operations [def.=definitions, obs.=observation, meas.=measurement, comp.=comparison, ex.=example, claim=claim, cons.=consequences, C&E=cause and effect, Ded.=deduction, Arg.=argumentation]

	Def.	Obs.	Meas.	Comp.	Ex.	Claim	Cons.	C&E	Ded.	Arg.
# incorrect	9	11	0	13	1	8	6	40	18	4
# operations	312	411	34	302	51	348	146	331	235	63

% incorrect per total ops.	3	3	0	4	2	2	4	10	8	6
----------------------------------	---	---	---	---	---	---	---	----	---	---

Evident in Table 7 is a relatively infrequent use of scientifically inaccurate information. That is, given that our unit of analysis is ideas, students are largely generating scientific ideas employing correct scientific information. This further justifies the move beyond simply considering scientific accuracy of students' conceptions towards considering the sophistication of their reasoning about scientific ideas. In this case, only considering the accuracy would have provided a very limited picture of what these students were doing in their writing. Because of the relative infrequency, it became important to consider the nature of the inaccuracies. For this inquiry, categorizing the inaccuracies by operation led to an interesting finding. The highest percentages of inaccuracy, though still relatively small, occurred with *cause and effect*, *deduction*, and *argumentation*. It is possible that higher complexity operations surface alternative conceptions more effectively. Further, the alternative conceptions elicited are potentially more deeply held, keeping in mind the Limited or Inappropriate Propositional Hierarchies (LIPs). That is, higher complexity operations draw on multiple domains and elements and may have the potential to reveal more of students' mental structures, and thus, expose LIPs.

Limitations

Though this framework provides a useful way to evaluate students' written work, it has a number of limitations. First, as noted above, this framework does not capture the scientific accuracy of students' written ideas. The utility of this tool, then, is limited to a narrower research goal—characterizing students' reasoning. When combined with an analysis of the content accuracy, however, this framework can provide unique insights about students' understanding. Further, this framework was conceptualized, tested,

1
2
3 refined, and ultimately applied to a corpus of writing in a very specific context—general
4
5 chemistry argumentative writing about ocean acidification. It is possible that some of
6
7 the ways that cognitive operations have been conceptualized in this study are specific to
8
9 this context. For this reason, applications to other contexts are needed to ensure the
10
11 domain-general nature of this framework. Finally, due to the relatively low occurrence
12
13 of certain operations in this context, we have a weaker understanding of some of the
14
15 operations (i.e., measurement). Because of the complete absence of the inductive
16
17 reasoning operation from Grimberg and Hand’s original framework in this set of student
18
19 writing, it was not included in this application, even though it is likely to be employed
20
21 in other contexts. Finally, this data was collected at a selective institution and it is likely
22
23 that different incorrect ideas or reasoning patterns would emerge from other student
24
25 populations. Again, this can be addressed by applying this framework to student writing
26
27 in other contexts.
28
29
30
31
32
33

34 **Discussion and Implications**

36 The first research question posed in this work considered how a cognitive operations
37
38 framework can be used to characterize students’ reasoning evident in their
39
40 argumentative writing. In this article, we show what this framework is like and how it
41
42 can be applied to students’ writing. We refined a list of cognitive operations generated
43
44 by Grimberg and Hand (2009) and organized them according to complexity, and then
45
46 used these operations to code general chemistry students’ writing on ocean
47
48 acidification. This framework has some key affordances that make it useful to both
49
50 research and practice. It is domain general, which means that it can be applied to writing
51
52 in a variety of contexts. We recommend, then, that others apply this to writing in a
53
54 variety of contexts across STEM and across levels (introductory to advanced student
55
56 populations). The domain-general nature of this potentially enables the identification of
57
58
59
60

1
2
3 differences in students' reasoning across disciplines and levels. For example, do
4
5 advanced students employ more complex reasoning than introductory students?
6

7
8 Another affordance of this framework is the 'score' that is a product of
9
10 application—the cognitive complexity. The single score output provides an estimate of
11
12 construct that is rather difficult to measure —student reasoning. This framework, then,
13
14 can potentially overcome some of the difficulties with evaluating writing reported in the
15
16 literature (Hamp-Lyon, 2016; Neill, 2002). This framework provides a novel approach
17
18 to assigning a holistic score to writing. Further, the use of cognitive operations enables
19
20 the identification of patterns in students' writing. That is, it can be used to characterize
21
22 the movement between cognitive operations and the likelihood of moving towards high
23
24 complexity operations, as shown in Grimberg and Hand's original application (2009).
25
26 This framework's capacity to capture temporal patterns makes it very useful for
27
28 understanding how students reason in extensive writing (Grimberg & Hand, 2009; Kelly
29
30 & Takao, 2002; Moreira et al., 2018).
31
32
33
34

35
36 The second research question aimed to elucidate what features of student
37
38 thinking were understandable with this framework. That is, what does this framework
39
40 evaluate the quality of? This was achieved in two ways. The first was to compare
41
42 framework estimates to instructor estimates of quality. This approach was intended to
43
44 determine if the framework estimates were similar to the instructor estimates and if both
45
46 were evaluating a similar construct. This revealed that perhaps for the upper bound of
47
48 the construct—argumentation—there was agreement between instructors and
49
50 framework estimates. There was less agreement for the other-than-best essays. Kelly
51
52 and Takao (2002) identified similar disparities between their framework estimates and
53
54 expert rankings and explained them as common occurrences when evaluating writing
55
56 (Wolcott & Legg, 1998). In our case, we argue that the variety was an artefact of the
57
58
59
60

1
2
3 presence of inaccurate scientific information in one of the essays. Instructors may not
4 separate content and reasoning as this framework does. However, we argue, similar to
5
6 Kelly and Takao (2002), that this framework may provide a tool for evaluating the
7
8 validity of instructor's estimates of quality. More research is necessary to establish
9
10 interrater reliability amongst instructor ratings and identify ways in which the
11
12 framework can serve as a tool for supporting instructors in systematically assessing
13
14 students' writing.
15
16
17
18

19 To determine if this framework was providing unique information about
20
21 students' ability, we compared cognitive complexity to other common performance
22
23 measures. There were no correlations. We posit two potential explanations for this. The
24
25 first is that this metric of cognitive complexity is indeed measuring something unique
26
27 from what typical performance metrics measure (National Research Council, 2001).
28
29 The second is that students who perform well on typical performance metrics do not
30
31 necessarily perform equally well on more extensive writing tasks (National Research
32
33 Council, 2001). Both of these explanations warrant further investigation because of the
34
35 implications for assessment. Specifically, this framework could serve to equip the
36
37 evaluation of more interesting competencies in students than that measured by typical
38
39 performance measures *or* assignments of this nature could serve to minimize advantages
40
41 certain groups bring with them to typical performance measures. However, we also
42
43 recognize that there may be other performance measures that correlate with the
44
45 framework estimate. Particularly, we would expect that more generative or authentic
46
47 assessments might correlate more strongly with cognitive complexity (National
48
49 Research Council, 2001). Finally, we aimed to characterize the relationship between
50
51 framework estimates of quality and scientific accuracy. In order to do this, we analysed
52
53 writing for the presence of scientific inaccuracies and coded the respective operation in
54
55
56
57
58
59
60

1
2
3 which they appeared. This revealed that scientific inaccuracies occurred relatively
4
5 infrequently with about 10 percent of *cause and effect* operations including something
6
7 that did not agree with scientifically accepted knowledge. The percentage among *cause*
8
9 *and effect* operations was the highest. However, there appears to be a trend in which
10
11 higher complexity operations (i.e., cause and effect, deduction, and argumentation) had
12
13 higher frequencies of incorrect information than low complexity operations. We argue,
14
15 in light of Novak's work on LIPs, that higher complexity operations as a
16
17 representation of students' mental models may better reveal LIPs (Novak, 2002). That
18
19 is, employing more complex reasoning may surface more deeply held LIPs. Students
20
21 who do not use higher complexity operations may be more limited in both their and
22
23 their instructors' capacity to address potential alternative conceptions. The relationship
24
25 between complexity and conceptual correctness warrants further investigation.
26
27 Understanding this relationship is important for designing formative assessments that
28
29 better elicit high complexity operations.
30
31
32
33

34
35 This framework also offers some unique implications for instructors who assign
36
37 similar tasks to their students. Scoring assignments in this way could permit an
38
39 instructor to draw conclusions about their students' collective access to complex
40
41 reasoning operations. For example, a low average score of cognitive complexity in their
42
43 course may motivate instructors to explicitly address and model complex reasoning
44
45 types for their students. However, we argue that the most important implication of this
46
47 framework for practice is providing a vocabulary to instructors for giving tailored
48
49 feedback to students. That is, applying this sort of framework would support an
50
51 instructor to give specific examples of when a student could have employed more
52
53 complex reasoning appropriately and instead used a less complex operation.
54
55
56
57
58
59
60

1
2
3
4
5
6
7 **Disclosure statement**
8

9
10 No potential conflict of interest was reported by the authors.
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 **References**

- 55 Emig, J. (1977). Writing as a Mode of Learning. *College Composition and*
56 *Communication*, 28(2), 122–128.
57
58 Gere, A. R., Limlamai, N., Wilson, E., MacDougall Saylor, K., & Pugh, R. (2019).
59 Writing and Conceptual Learning in Science: An Analysis of Assignments. *Written*
60 *Communication*, 36(1), 99–135. <https://doi.org/10.1177/0741088318804820>

- 1
2
3 Grimberg, B. I., & Hand, B. (2009). Cognitive Pathways : Analysis of students ' written
4 texts for science understanding. *International Journal of Science Education*, 31(4),
5 503–521. <https://doi.org/10.1080/09500690701704805>
6
7 Gunel, M., Hand, B., & Prain, V. (2007). Writing for learning in science: A secondary
8 analysis of six studies. *International Journal of Science and Mathematics*
9 *Education*, 5, 615–637. <https://doi.org/10.1007/s10763-007-9082-y>
10
11 Ha, M., Nehm, R. H., Urban-Lurain, M., & Merrill, J. E. (2011). Applying
12 computerized-scoring models of written biological explanations across courses and
13 colleges: Prospects and limitations. *CBE Life Sciences Education*, 10(4), 379–393.
14 <https://doi.org/10.1187/cbe.11-08-0081>
15
16 Hein, V., & Smerdon, B. (2013). *Predictors of Postsecondary Success. College and*
17 *Career Readiness and Success Center at American Institutes for Research.*
18
19 Kelly, G. J., & Bazerman, C. (2003). How Students Argue Scientific Claims : A
20 Rhetorical-Semantic Analysis. *Applied Linguistics*, 24(1), 28–55.
21
22 Kelly, G. J., Chen, C., & Prothero, W. (2000). The Epistemological Framing of a
23 Discipline : Writing Science in University Oceanography, 37(7).
24
25 Kelly, G. J., Regev, J., & Prothero, W. (2007). Analysis of Lines of Reasoning in
26 Written Argumentation. In *Argumentation in Science Education* (pp. 137–157).
27
28 Kelly, G. J., & Takao, A. (2002). Epistemic levels in argument: An analysis of
29 university oceanography students' use of evidence in writing. *Science Education*,
30 86(3), 314–342. <https://doi.org/10.1002/sce.10024>
31
32 Keys, C. W. (1994). The development of scientific reasoning skills in conjunction with
33 collaborative writing assignments: An interpretive study of six ninth-grade
34 students. *Journal of Research in Science Teaching*, 31(9), 1003–1022.
35 <https://doi.org/10.1002/tea.3660310912>
36
37 Klein, P. D. (1999). Reopening Inquiry into Cognitive Processes in Writing-To-Learn.
38 *Educational Psychology Review*, 11(3), 203–270.
39 <https://doi.org/10.1023/A:1021913217147>
40
41 Klein, P. D. (2015). Mediators and Moderators in Individual and Collaborative Writing
42 to Learn. *Journal of Writing Research*, 7(1), 201–214.
43
44 Krippendorff, K. (2004). Reliability in Content Analysis : Some Common
45 Misconceptions and Recommendations Reliability in Content Analysis : Some
46 Common Misconceptions and. *Human Communication Research*, 30(3), 411–433.
47
48 Liu, O. L., Rios, J. A., Heilman, M., Gerard, L., & Linn, M. C. (2016). Validation of
49 automated scoring of science assessments. *Journal of Research in Science*
50 *Teaching*, 53(2), 215–233. <https://doi.org/10.1002/tea.21299>
51
52 Moon, A., Gere, A. R., & Shultz, G. V. (2018). Writing in the STEM classroom:
53 Faculty conceptions of writing and its role in the undergraduate classroom. *Science*
54 *Education*, 0(0). <https://doi.org/10.1002/sce.21454>
55
56 Moreira, P., Marzabal, A., & Talanquer, V. (2018). Using a mechanistic framework to
57 characterise chemistry students' reasoning in written explanations. *Chemistry*
58 *Education Research and Practice*. <https://doi.org/10.1039/C8RP00159F>
59
60 National Research Council. (2001). *Knowing what students know: The science and*
design of educational assessment. National Academies Press. Washington DC.
<https://doi.org/10.17226/10019>
Novak, J. D. (2002). Meaningful Learning: The Essential Factor for Conceptual Change
in Limited or Inappropriate Propositional Hierarchies Leading to Empowerment of
Learners. *Science Education*, 86(4), 548–571. <https://doi.org/10.1002/sce.10032>
Prain, V., & Hand, B. (2016). Coming to Know More Through and From Writing.
Educational Researcher, 45(7), 430–434.

- 1
2
3 <https://doi.org/10.3102/0013189X16672642>
4 Reynolds, J. A., Thaiss, C., Katkin, W., & Thompson, R. J. (2012). Writing-to-learn in
5 undergraduate science education: A community-based, conceptually driven
6 approach. *CBE-Life Sciences Education*, 11(1), 17–25.
7
8 Sandoval, W. A. (2003). Conceptual and Epistemic Aspects of Students ' Scientific
9 Explanations. *Journal of the Learning Sciences*, 12(1), 5–51.
10 <https://doi.org/10.1207/S15327809JLS1201>
11 Sandoval, W. A., & Millwood, K. A. (2005). The Quality of Students ' Use of Evidence
12 in Written Scientific Explanations. *Cognition and Instruction*, 23(1), 23–55.
13 <https://doi.org/10.1207/s1532690xci2301>
14 Sandoval, W. A., & Reiser, B. J. (2004). Explanation-driven inquiry: Integrating
15 conceptual and epistemic scaffolds for scientific inquiry. *Science Education*, 88(3),
16 345–372. <https://doi.org/10.1002/sce.10130>
17
18 Sevian, H., & Talanquer, V. (2014). Rethinking chemistry: a learning progression on
19 chemical thinking. *Chemistry Education Research and Practice*, 15(1), 10–23.
20 <https://doi.org/10.1039/C3RP00111C>
21
22 Takao, A. Y., & Kelly, G. J. (2003a). Assessment of Evidence in University Students '
23 Scientific Writing. *Science & Education*, 12, 341–363. <https://doi.org/10.1023/A>
24 Takao, A. Y., & Kelly, G. J. (2003b). Assessment of Evidence in University Students '
25 Scientific Writing. *Science & Education*, 12, 341–363.
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

1
2
3 **Application and testing of a framework for characterizing the quality**
4 **of scientific reasoning in chemistry students' writing on ocean**
5 **acidification**
6
7
8
9

10 Alena Moon, Robert Moeller, Anne Ruggles Gere†, and Ginger V. Shultz

11 *Department of Chemistry, University of Michigan, Ann Arbor, MI, USA*

12 *†Sweetland Center for Writing, University of Michigan, Ann Arbor, MI, USA*

13
14
15 Corresponding author: Ginger V. Shultz, gshultz@umich.edu, 2521 Chemistry,
16 University of Michigan, Ann Arbor, MI, 48109.
17
18
19

20 Provide short biographical notes on all contributors here if the journal requires them.
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

Development and testing of a framework for characterizing the quality of scientific reasoning in students' writing on ocean acidification

Science educators recognize the need to teach scientific ways of knowing and reasoning in addition to scientific knowledge. However, characterizing and assessing scientific ways of knowing and reasoning is challenging. Writing-to-learn offers one way of eliciting and supporting students' reasoning; further, writing serves to externalize and make traceable students' reasoning. For this reason, it is a useful formative assessment of scientific reasoning. The utility hinges on researchers' ability to understand what students can do and think from their writing. Given the challenges in assessing students' writing, this research offers an adapted framework for assessing students' scientific reasoning evident in writing. This work will introduce the adapted framework and show an application to general chemistry students' argumentative writing about ocean acidification. We provide evidence that this framework can be used to validly estimate the quality of students' reasoning. We argue that this framework offers some affordances that overcome challenges reported in the literature. It serves to define scientific reasoning in a domain-general way by breaking it down into its components, but in a way that can produce a composite score that tells us about how students reason using chemistry content. Further, the framework provides a way to characterize the scientific accuracy of students' reasoning that can inform instructors' treatment of alternative conceptions.

Keywords: Writing-to-learn; assessing writing; scientific reasoning

Background

Science educators recognize that it is insufficient to only teach students' scientific knowledge as a collection of concepts and topics. Rather, to enable students to use scientific knowledge, we must support the development of reasoning and thinking skills that scientists use (NRC, 2012; Sevian & Talanquer, 2014; Bulte, Westbroek, De Jong, & Pilot, 2006). Writing-to-learn (WTL) is one way of supporting the development of this skill by activating deep thinking and reasoning in students (Keys, 1999) and, more importantly, making that reasoning visible and traceable (Emig, 1977; Kelly & Takao,

1
2
3 2002; Kelly, Regev, & Prothero, 2007). From an assessment perspective, this evidence
4 of student reasoning is valuable in so far as researchers and practitioners can use it to
5 make an argument about students' abilities to reason scientifically (Lavery et al., 2016;
6 NRC, 2001). However, there are challenges that currently limit the utility of this
7 evidence. There are few widely agreed upon epistemic criteria for characterizing the
8 quality of students' reasoning (i.e., what makes one students' reasoning better than
9 another's). Further, actually applying these criteria to understand and evaluate students'
10 writing is difficult as writing requires the researcher to make choices about grain size,
11 whether to evaluate structure or content or both, and what the presence or absence of a
12 quality criterion actually looks like in students' writing (Kelly & Takao, 2002; Takao &
13 Kelly, 2003a). To address these challenges, we have modified and applied a framework
14 for characterizing and evaluating reasoning in students' argumentative writing. This
15 framework contributes meaningfully to efforts to conceptualize and evaluate scientific
16 reasoning, as well as to efforts to analyse writing, which poses unique challenges.

35 ***Writing to Learn***

37 Writing-to-learn refers to the kind of informal writing about science that facilitates
38 learning and ownership of scientific ideas. This informal writing is distinct in that its
39 primary aim is not to communicate or display mastery to an instructor, but to actually
40 facilitate sense-making by activating deep thinking and interaction with the concepts
41 (Keys, 1994). A secondary benefit of writing-to-learn, then, is promoting engagement
42 with disciplinary norms of writing and thinking (Prain & Hand, 2016). There is quite a
43 bit of variation around this primary aim, however; WTL assignments take a variety of
44 forms, lengths, methods of text production, audiences, and genres (Keys, 1994).

56 A secondary analysis of six writing-to-learn studies revealed some promising
57 gains as a result of writing-to-learn—the treatment condition outperformed comparison
58
59
60

1
2
3 groups on a total test scores and conceptual question scores and this effect was largely
4
5 due to the treatment (Gunel, Hand, & Prain, 2007). All six studies followed a similar
6
7 design including a pre-test/post-test design with the test having multiple-choice and
8
9 conceptual extended response questions. More importantly, all writing interventions
10
11 were grounded in the same theoretical considerations that have been identified as key
12
13 for successful learning from writing: 1) opportunities for brainstorming, 2) provision of
14
15 authentic audiences, 3) drafting and redrafting with feedback, 4) explicit instruction of
16
17 genre specifications, 5) focus on big ideas, 6) use of rubrics, and 7) diverse
18
19 opportunities to plan and draft writing (Gere, Limlamai, Wilson, MacDougall Saylor, &
20
21 Pugh, 2019; Gunel et al., 2007; Klein, 1999, 2015). The theoretical grounding afforded
22
23 comparisons across domains and writing assignment types and served to reveal the
24
25 benefits of WTL more broadly (Gunel et al., 2007; Prain & Hand, 2016). However, at
26
27 the undergraduate STEM level specifically, more work is needed to understand the
28
29 mechanism of effect for WTL assignments (Reynolds, Thaiss, Katkin, & Thompson,
30
31 2012) and we argue that to undertake investigations into the mechanism of effect, we
32
33 need a reliable and meaningful framework for interpreting and evaluating students'
34
35 written work.
36
37
38
39
40
41
42
43

Characterizing Students' Reasoning in Written Products

44
45 Constructed responses reveal rich insight into students' ideas and the coherence of those
46
47 ideas, but evaluating open responses remains a barrier to implementing such rich
48
49 assessments (Liu, Rios, Heilman, Gerard, & Linn, 2016). This barrier consists of two
50
51 distinct but interdependent challenges: characterizing the quality usually in some sort of
52
53 rubric (Kelly & Bazerman, 2003; Sandoval, 2003; Sandoval & Millwood, 2005) and
54
55 consistently and reliably applying the quality criteria (Ha, Nehm, Urban-Lurain, &
56
57 Merrill, 2011; Liu et al., 2016). Additionally, the nature and difficulty of these
58
59
60

1
2
3 challenges vary with the length of constructed response; for example, extended
4
5 arguments in the form of research reports in oceanography tended to have long and
6
7 complex chains of reasoning that are difficult to characterize with a single rubric (Kelly,
8
9 Regev, & Prothero, 2007; Kelly & Takao, 2002). Researchers have sought to overcome
10
11 these challenges with a variety of approaches for a variety of written products, ranging
12
13 from short written explanations (Ha et al., 2011; Liu et al., 2016; Moreira, Marzabal, &
14
15 Talanquer, 2018; Sandoval, 2003; Sandoval & Millwood, 2005; Sandoval & Reiser,
16
17 2004) to more extensive writing like laboratory reports or research reports (Grimberg &
18
19 Hand, 2009; Kelly, Chen, & Prothero, 2000; Kelly et al., 2007; Kelly & Takao, 2002;
20
21 Takao & Kelly, 2003b). A few of these approaches are distinct in that they break down
22
23 students' responses into smaller units to then identify patterns in students' reasoning, as
24
25 opposed to a rubric that considers the quality of the response as a whole (Grimberg &
26
27 Hand, 2009; Kelly et al., 2007; Kelly & Takao, 2002; Moreira et al., 2018). These
28
29 approaches will be the focus of this literature review.
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

Moreira et al. (2018) specifically sought to characterize the causal reasoning of 10th grade chemistry students' explanations of freezing point depression. To do so, the authors modified and applied a discourse analysis framework developed by Russ et al. (2008) to students' written explanations and drawings. The final form of the analysis scheme included four components—entities, properties, activities, and organisation—and the relationships between the components and the students' drawings that were identified in students' responses. Entities are the 'things' in the system that are being considered. Properties are characteristics of those entities and activities are actions of those entities. Organisation refers to the spatial-temporal relationship between the entities and activities or properties of the system. By coding the explanations for these components and relationships, they were able to elucidate patterns in students'

1
2
3 explanations and organize these patterns into four levels according to the quality of
4 causal reasoning. These levels increased in sophistication from descriptive to relational
5 to simple causal, culminating in emerging mechanistic. Unsurprisingly given the
6 authors' previous findings (Sevian & Talanquer, 2014), the majority of students (45%)
7 used relational causal reasoning (Moreira et al., 2018). Explanations at this level could
8 be modelled to show that students generally identified two entities and the properties of
9 one or both entities, and then related either the properties to each other or related
10 entities to properties. Such complex modelling of students' explanations can hopefully
11 equip teachers with more sophisticated approaches to understanding, interpreting, and
12 developing students' reasoning abilities (Moreira et al., 2018).

13
14
15
16
17
18
19
20
21
22
23
24
25
26 Grimberg and Hand (2009) similarly identified the presence or absence of
27 dimensions of reasoning and determined *patterns* in students' reasoning evident in their
28 laboratory reports. In this study, Grimberg and Hand (2009) identify cognitive
29 operations used in writing laboratory reports and then construct what they term
30 'cognitive pathways'—the sequence of cognitive operations used by author(s) of a lab
31 report. The authors argued that because writing a laboratory report was a meaning-
32 making activity, considering the sequence of cognitive operations revealed how students
33 were constructing meaning. Using a list of 11 cognitive operations derived partially
34 from the literature and from the students' data, authors coded students' writing for use
35 of cognitive operations. The cognitive operations included observation, measurement,
36 comparison, analogy, clarifications, claim, cause/effect, induction/generalization,
37 deduction, investigation design, and argumentation (Grimberg & Hand, 2009). When
38 comparing cognitive pathways of low achievers to high achievers, as determined by a
39 standardized skills test, both high and low achievers used the same range of operations,
40 but with a different structure. Though the cognitive structure was partially determined
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60

1
2
3 by the structure of the Science Writing Heuristic (SWH) activity (i.e., clarification
4 questions were posed during the research question portion of SWH activity), high
5
6 achievers began using complex operations earlier in the text than low achievers. This
7
8 ultimately demonstrated that SWH scaffolding supported all students in using high-
9
10 complexity operations, albeit at different rates (Grimberg & Hand, 2009).
11
12
13

14
15 Grimberg and Hand's (2009) work demonstrates the capacity of writing to make
16
17 students' thinking visible and traceable. Emig (1977) argues that writing is unique in its
18
19 capacity to do this. Kelly and Takao (2002) demonstrate the utility of students' writing
20
21 for understanding their reasoning by characterizing how undergraduate students use
22
23 evidence to construct arguments. To make this characterization, they developed a
24
25 research methodology that models epistemic levels of argument. This framework
26
27 includes six levels ranging from the lowest, data charts and representations, to the
28
29 highest, general geological (the specific context for this work) knowledge not specific to
30
31 the data presented. The levels represented students' ability to abstract from data to make
32
33 claims. With this framework, the authors analysed a subset of undergraduate
34
35 oceanography students' assignments by labelling each sentence with an epistemic level.
36
37 These epistemic criteria were then weighted in order to rank the 24 student arguments
38
39 (research reports in oceanography) from best to worst. While this framework was very
40
41 useful for characterizing the quality of students' arguments based on their use of
42
43 evidence, it did possess limitations. Namely, the assessment of quality determined by
44
45 the framework did not always align with content experts' evaluation of quality, the
46
47 framework did not consider inference logic (how the data led to theoretical claims), and
48
49 the authors had to make inferences in their application of the framework. These
50
51 limitations are difficult to overcome for everyone aiming to evaluate ill-defined
52
53 constructs like use of evidence. However, Kelly and Takao (2002) revealed that claims
54
55
56
57
58
59
60

1
2
3 can be made about students' reasoning from their written work. In order for the insight
4
5 into students' reasoning provided by writing to be useful for informing students'
6
7 development, tools for assessing writing must be efficient, systematic, and offer tailored
8
9 feedback.
10

11
12 Ongoing discussions about writing assessment distinguish between holistic
13
14 scoring, assigning a single score to a broad variable like writing proficiency, and
15
16 analytic scoring, breaking down variables like writing proficiency into components that
17
18 are individually scored (Hamp-Lyons, 2016a, b). High-stakes assessments, such as
19
20 college entrance examinations, motivated the use of both holistic and analytic scoring
21
22 approaches to assign students general writing scores, but researchers have begun to
23
24 identify shortcomings of both (Neill, 2002; Hamp-Lyons, 2016b, Chapman, 2016).
25
26 With more complex analytical tools (i.e., multivariate analyses), Hamp-Lyons (2016)
27
28 calls for a movement to multiple trait scoring of writing. In multiple trait scoring, there
29
30 is no single score given, whether composite or holistic. Rather, a set of scores is
31
32 assigned with multiple traits each warranting a score, thus lending to a richer
33
34 description of students' ability (Hamp-Lyons, 2016a, b).
35
36
37
38
39
40

41 ***Rationale and Research Objectives***

42
43 In any effort to measure a student's reasoning, choices must be made about what
44
45 characterizes quality. Specific to evaluating extensive writing, additional decisions must
46
47 be made about the grain size, level, and nature of rubric that will be used to determine
48
49 quality. In order to address these challenges, the cognitive operations used by Grimberg
50
51 and Hand (2009) were modified and applied to students' writing on ocean acidification.
52
53 Motivated to leverage the rich insight into student thinking that constructed responses
54
55 offer, the work presented herein aimed to test an approach for analysing extensive
56
57
58
59
60

1
2
3 scientific writing by applying it to a new context. We aimed to answer the following
4
5 questions regarding this approach:
6

- 7
8 (1) Can cognitive operations be used to make sense of general chemistry students'
9
10 argumentative writing? If so, how?
11
12 (2) What features of students' argumentative writing do cognitive operations serve
13
14 to explain?
15
16 (3) What is the relationship between framework estimates of quality and conceptual
17
18 correctness?
19
20
21
22

23 **Methods**

24 *Participants, setting, and data collection*

25
26 A writing prompt was designed and administered in a first-semester General Chemistry
27
28 course serving primarily students in the College of Engineering and undeclared students
29
30 in the College of Literature, Science, and Arts. This course had an enrolment of 1413
31
32 students, most of whom were freshman and sophomores. The content of the course
33
34 covered traditional general chemistry concepts, ranging from dimensional analysis,
35
36 quantum mechanical atomic models, bonding theories, to reactions, enthalpy,
37
38 intermolecular forces, chemical equilibrium, and acid-base theories.
39
40
41
42
43

44 This course is structured with three lectures led by an instructor and one
45
46 discussion session led by a teaching assistant per week. During each discussion section,
47
48 students complete a quiz. The writing assignment for this study was administered as a
49
50 substitute for a quiz. The writing assignment was uploaded as a .pdf file to the course
51
52 management site one week before the due date. This writing assignment directed
53
54 students to consider a set of concepts in their response and to keep their post between
55
56
57
58
59
60

1
2
3 350 and 500 words. Though the majority of students' responses were within this range,
4
5 some wrote less and some wrote more.
6

7
8 Students all submitted their writing assignment online the following week at the
9
10 start time of their specific discussion session. For this reason, some students with
11
12 discussion sessions later in the week had more time to write than those with earlier
13
14 recitations. During the discussion, students formed teams of two or three, switched
15
16 papers and gave feedback to each other. Six hundred seventy-three students gave
17
18 consent to have their writing analysed. Ethical review board approval was gained in
19
20 order to collect and analyse written assignments that students consented to have
21
22 analysed. Students did not receive any feedback on their written work beyond the
23
24 conversation that took place in their discussion session. We found that students did not
25
26 make meaningful revisions following the peer review discussion. Additionally, they
27
28 were not required to submit a revision. However, if they did, that revised draft was used
29
30 for analysis.
31
32
33
34
35

36 ***Writing activity development and design***

37
38 The writing prompt was developed iteratively through correspondence with authors who
39
40 had expertise in writing to learn and the development of meaningful writing prompts
41
42 (AG) and the faculty members teaching the course who collectively held more than two
43
44 decades of experience teaching chemical equilibrium in general chemistry. WTL
45
46 prompts are generally designed to provide students with an audience, an identity, and an
47
48 authentic context that require students to engage with a specific concept. This WTL
49
50 assignment was intentionally designed with elements empirically determined to
51
52 contribute to meaningful learning through writing (Gere et al., 2019). In this case, the
53
54 prompt showed a fake social media post, in which 'Ernie Clueless' shares a plot
55
56 illustrating the trend of concentration of atmospheric carbon dioxide and ocean pH over
57
58
59
60

1
2
3 time. Ernie claims that these things are unrelated. Students were tasked with explaining
4 the relationship to Ernie, given the relevant equilibria. The prompt targeted the concept
5 of chemical equilibrium, drawing on Le Châtelier's principle. It inherently supported
6 argumentation by requiring students to differentiate their perspective from Ernie's.
7
8
9
10
11

12 13 14 ***Data analysis: Development and application of analytical framework***

15 A list of cognitive operations was modified from a list used by Grimberg and Hand
16 (2009) to analyse reports from a Science Writing Heuristic (SWH) laboratory. In this
17 context, a cognitive operation is a written discursive move that serves some cognitive
18 objective. Cognitive operations, then, determined the grain size for breaking down an
19 essay into smaller analysable units (i.e. the number of sentences that served the
20 objective of a claim, for example, were coded as such). For this reason, a claim could be
21 one sentence in one essay and three sentences in another. The amount of text that was
22 assigned a code was determined by the function of that text. The list of cognitive
23 operations used by Grimberg and Hand (2009) was refined iteratively by testing it
24 against the data. This involved using Grimberg and Hand's original set to code the text,
25 identifying text that could not be coded with this set and operations from this set that
26 did not serve to explain any of the data, and refining the set of operations to a set that
27 were all used to describe virtually all of the text. Multiple initial iterations occurred with
28 the same subset of essays (N=25) and subsequent iterations incorporated more essays on
29 an as needed basis. Table 1 shows the final list of cognitive operations that was used
30 throughout the final analysis. Included in Table 1 is a characterization of the
31 dimensionality of each operation. As will be explained further in the theoretical
32 framework, the complexity of cognitive operations was determined by its
33 dimensionality—number of ideas being drawn upon. Operations with two domains were
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

more cognitively complex than operations with one dimension. That is, they drew upon and connected more idea units.

Table 1. Finalized list of cognitive operations and descriptors used to analyse all essays, listed in order of increasing cognitive complexity

Cognitive Operation	Description	Dimensionality
1. Definition	Canonical description of a term, concept, idea, or theory	single domain
2. Observation	Qualitative description of change, trend, or transformation of a variable	
3. Measurement	Quantitative description of change, trend, or transformation of a variable	
4. Comparison	Relationship of change, trend, or transformation for two or more variables	
5. Example	Illustration of a class of objects by singling out one object	two domains
6. Claim	Assertion supported with a tentative explanation	
7. Consequences	Cause and effect explanation with either cause or effect falling outside the scope of the writing prompt	
8. Cause and effect	Explanation providing a mechanism with a causal agent and observed effects	multiple domains
9. Deduction	Application of a theory or principle to a specific system or scenario	
10. Argumentation	Explicit differentiation between the author's perspective and the fictional character's perspective*	

*This conceptualization of argumentation was specific to the context of writing prompt used in this study. It is expected that it could be easily translated to other writing contexts that require argumentation.

Once a list of cognitive operations was finalized, the first two authors began coding assignments and built a detailed rubric that included definitions, linguistic markers, and examples for each operation. This rubric was further refined through multiple iterations of analysis by a team of chemistry education researchers in an effort to establish inter-rater reliability (IRR). This stage included a team consisting of the first two authors and another chemistry education researcher trained in qualitative coding of writing. The graduate student was trained on an existing rubric, the whole team coded ten assignments, and an IRR coefficient in the form of Krippendorff's alpha (KA) was calculated. This coefficient was quite low after the first round, so revisions were made to the rubric, another training session was conducted, and subsequent rounds of analysis

1
2
3 were conducted until a Krippendorff's alpha value of 0.69, the minimum acceptable
4 value, was achieved. Training involved discussing the rubric which included examples
5 from many essays and then illustrating the coding process by coding a few essays all
6 together.
7
8
9
10
11

12 The first two authors then coded approximately 200 assignments each, with
13 large overlap. Having two researchers code many of the same assignments lent to the
14 reliability of the coding. Throughout analysis, IRR 'checks' were performed to ensure
15 that the rubric was being applied consistently. This involved selecting overlapping
16 assignments and determining a KA. This stage resulted in a KA of 0.89, which was a
17 desirable value (Krippendorff, 2004).
18
19
20
21
22
23
24
25
26

27 *Estimate of quality*

28
29 Once all assignments were coded for cognitive operations, estimates of quality were
30 assigned according to the cognitive complexity of the essay. The cognitive operations
31 are ordered in Table 1 according to increasing cognitive complexity; that is, a definition
32 has the lowest cognitive complexity (1) whereas argumentation has the highest
33 cognitive complexity (10). Overall cognitive complexity for the essay was determined
34 by taking a weighted average of operations used. Because the magnitude of text
35 included in an operation varied from one essay to another, the average was weighted by
36 the number of sentences within that operation. An assignment, then, that was 50%
37 argumentation would likely have a higher average complexity than an assignment with
38 only 10% argumentation. This approach resulted in a single number characterizing the
39 quality of students' reasoning in an essay. This process of producing a single number is
40 illustrated with the example essay in Table 2.
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60

Table 2. Example of student essay and determination of cognitive complexity score

<p>Le Chatlier's principle is a way to show how if the equilibrium of a certain reaction is altered by changing certain aspects, then the reaction or equilibrium position will actually shift to fix the change. In Ernie's post there are multiple reactions that follow one after the other, ultimately showing how CO₂ transforms to H₂CO₃, which changes to H⁺ and HCO₃, which finally changes to 2H⁺. These reactions are connected in that as one breaks down, it spurs another reaction to then form. This shows how CO₂ in the atmosphere actually affects PH concentration and ocean acidification (or the number/concentration of H⁺ ions in the reaction). Because of this correlation, certain components like temperature, pressure, volume, and concentration can affect CO₂ levels and acidification. In our case, we will look at the concentration or amount of CO₂ in correlation to the amount of H⁺ that is formed. If the concentration of CO₂ (as a reactant) were to increase, the system would try to decrease it; this would mean that the concentration of the other reactants will increase to react with the CO₂, but then more product would be formed. Thus increasing the concentration of CO₂ (atmos) would cause the system to shift towards the products, producing more CO₂(aq). Then because there is more CO₂ (aq), more water would have to be use so more H₂CO₃ would form. The next reaction would then proceed as the others with an increase in the formation of the products H⁺ and HCO₃, which in turn would once again increase the next reaction producing more 2H⁺ product. This may seem confusing, but to summarize if CO₂ as a reactant increases, it would need to counteract this change by producing more product of 2H⁺ (acidification). Now, if CO₂ (atmos) reactant were to decrease, the equilibrium will actually try to increase it so as to set the system back at equilibrium. This would mean that a decrease in that reactant would cause more product to be formed so to produce more CO₂; thus 2H⁺ product would increase to make more CO₂.</p> <p>The relationship between CO₂ and pH is that as CO₂ increases, the pH will decrease and make the oceans more acidic. The pH scale runs from 1 to 14, with the acids being numbers 1 to 6. So the lower the pH, the more acidic and the more H⁺ ions that are formed (the H⁺ are indicators of acidic properties). Le Chatlier's principle shows that increasing one will have to increase the other so that the system is able to be equal again, thus there is correlation between CO₂ (atmos) increase and H⁺ ions increase(acidification)."</p>	<p>Definition (1)</p> <p>Observation (2)</p> <p>Claim (6)</p> <p>Deduction (9)</p> <p>Comparison (4) Definition (1) Comparison (4)</p> <p>Cause and Effect (8)</p>
<p>Cognitive complexity = $\frac{\sum(\text{cognitive operation score} * \# \text{ sentences used})}{\text{total number of sentences}}$</p> <p>Definition (1) + 3 * Observation (2) + Claim (6) + 9 * Deduction (9) + Comparison (4) + Definition (1) + Comparison (4) + Cause & Effect (8) = 111</p> <p>111/18 =6.2 (Average complexity for this essay)</p>	

Faculty ranking of essays

One motivation of this work is to use a framework to systematically characterize students' writing. In order to understand how this framework could accomplish that task, we compared framework estimates of quality with instructors' estimates of quality. The precedent for this approach is offered by Kelly and Takao's testing of the framework for epistemic levels (Kelly & Takao, 2002). Further, evidence from interviews with STEM faculty about writing suggested to us that the knowledge experts use to evaluate writing is tacit (Moon, Gere, & Shultz, 2018). By comparing expert rankings with essays ordered according to a framework, we can identify ways that the framework may capture this tacit knowledge. To compare framework estimates of quality and instructor estimates, we compared essays ranked according to the framework to instructor rankings. To select the essays, we split cognitive complexity into four roughly equivalent ranges (3.1-4.8, 4.8-6.5, 6.5-8.2, 8.2-10), where the highest range included the 'best' essays, according to the framework. Within each range, an essay was randomly selected. The output of this step was four essays ranging from most complex to least complex, as determined by the framework. These essays were then provided to instructors who were tasked with ranking the essays from best to worst according to the quality of scientific reasoning, which they were directed to evaluate as they saw fit. The rationale for not providing faculty with more extensive quality criteria was to target the kind of evaluative work that is inherent to grading this sort of task; that is, we wanted instructors to make the kinds of decisions that are required to define and evaluate scientific reasoning. Five faculty with a range of experience teaching general chemistry from different institutions ranked the writing tasks and one instructor volunteered the reasoning behind their ranking.

Chemistry Content Analysis

1
2
3 Essays were examined to understand how students employed ideas about
4
5 chemical equilibrium within their argument. Essays were coded for conceptual
6
7 correctness, which involved flagging all occurrences of inaccurate ideas. Every time an
8
9 inaccurate idea was identified, the cognitive operation containing that incorrect idea was
10
11 marked. The pairing of the incorrect idea with the cognitive operation was intentional,
12
13 based on the theoretical framework explained below, which posits that what students
14
15 articulate in written forms are representations of the ideas they hold. In this case those
16
17 ideas related to chemical equilibrium.
18
19
20
21
22

23 *Theoretical framework*

24
25 The primary assumption made in this work is that writing reveals students'
26
27 cognitive structures or understanding of the meaning of a concept (Emig, 1977; Novak,
28
29 2002). Meaning is defined in this case as 'the totality of propositions linked to any
30
31 given concept,' excluding the emotional association with the concept and the context in
32
33 which the concept was learned (Novak, 2002). This assumption is grounded in the
34
35 capacity of writing to 1) connect or relate propositions in the author's mind and 2) make
36
37 'evolutionary development of thought graphically visible and available (Emig, 1977).'

38
39 According to Novak's theory of meaningful learning, the complexity of the meanings
40
41 can be evaluated, the quantity and quality of which will determine meaningful learning
42
43
44
45
46 (Novak, 2002).
47

48 In this work, units of written text were coded according to their cognitive
49
50 function as a cognitive operation. Each cognitive operation was constituted by some
51
52 number of ontological domains, elements within those domains, and relationships
53
54 between each (Grimberg & Hand, 2009; Halford & McCredde, 1998). Cognitive
55
56 complexity, then, was defined in terms of dimensionality; in which cognitive operations
57
58 with higher numbers of domains, elements, and relationships were more complex than
59
60

1
2
3 cognitive operations with fewer (Halford & McCredden, 1998). These domains and
4
5 elements can be conceptualized similarly to the discourse analysis framework used by
6
7 Moreira et al. (2018) where they refer to ‘things’ in the system being considered. All
8
9 cognitive operations were organized on a spectrum from least complex to most complex
10
11 according to this criterion, illustrated in Table 1. A students’ scientific reasoning can
12
13 then be considered as the progression of operations used in a text.
14
15

16
17 Therefore, one way of evaluating the quality of meanings was through cognitive
18
19 complexity. It is possible, however, that a student can employ cognitively complex
20
21 reasoning without necessarily using correct content (Sandoval & Millwood, 2005; Kelly
22
23 & Takao, 2002). For this reason, a second way of evaluating the quality of meanings
24
25 was considering their conceptual correctness; that is, their agreement with scientifically
26
27 accepted knowledge. A conceptual change perspective suggests that problems of
28
29 incorrect conceptions arise from Limited or Inappropriate Propositional Hierarchies
30
31 (LIPs), the way that concepts are inappropriately organized in the learner’s mind,
32
33 which means that we as instructors must consider both the *content and structure* of
34
35 incorrect conceptions. Further, this implies a greater instructional effort is needed to
36
37 remediate stable LIPs (Novak, 2002). The cognitive operations framework used in this
38
39 study helps characterize the complexity or stability—as a function of the number of
40
41 domains, elements, and relationships—of conceptions.
42
43
44
45
46
47

48 **Results**

49 50 51 ***Research Question 1: Can cognitive operations be used to make sense of*** 52 53 ***general chemistry students’ argumentative writing? If so, how?***

54
55 Table 2 illustrates how these cognitive operations were interpreted to analyse the
56
57 writing with examples from students’ essays. The examples are useful for discussing the
58
59 difficulties of applying this framework. The lower complexity operations (*definition*
60

through *example*) were relatively easy to identify. The difference between *observation* and *measurement* was essentially a difference between qualitative and quantitative, with *measurements* requiring some numerical component. Because students could sufficiently respond to this writing prompt with qualitative reasoning, *measurement* occurred less frequently (Table 1). A *comparison* was distinct from *observation* and *measurement* in describing more than one variable relative to each other. A *claim* was similar to *comparison* in referencing a relationship between two variables but was distinct in that it required a tentative explanation of the relationship. For the *claim*, then, there were two domains (explanation and observation). *Cause and effect* and *consequences* were similar and easier to identify in text, with common linguistic markers being ‘caused,’ ‘leads to,’ or ‘drives.’ *Consequences* used cause and effect reasoning but relied on causes or effects that fell outside the scope of the prompt. In this case, the student referenced the effect of ocean acidification on coral. A primary marker of *deduction* was the invocation of a principle or theory that was then applied to a specific system. In this context, students frequently invoked Le Chatalier’s principle or equilibrium. Finally, *argumentation* was undoubtedly the most difficult to identify as it often drew on multiple operations. So, we used this feature as an identifier. *Argumentation*, then, required indistinguishable use of multiple operations and to explicitly differentiate between the author’s perspective and Ernie’s (or any opposing position in another context). Table 1, 2, and 3 together can be used to apply this framework to other contexts.

Table 3. Examples of cognitive operations, arranged in order of increasing complexity. Numbers describe cognitive complexity ordering.

Operation	Example from student essay
Definition (1)	Equilibrium is a state where a reaction is occurring forwards and backwards at equal rates with no overall change. When a change occurs to the system, the reaction will shift in a direction to counteract this change.

1 2 3 4 5	Observation (2)	The trend in the graphs shows an increase in atmospheric carbon dioxide over time.
6 7 8	Measurement (3)	The ocean pH has dropped from 8.2 to 8.1 since the Industrial Revolution.
9 10 11 12	Comparison (4)	The plot that you shared illustrates that as the concentration of atmospheric increases over time, the pH of the ocean seawater decreases.
13 14 15	Example (5)	An example is hydrochloric acid, which has hydrogen ions attached and is an acid with a rather low pH level.
16 17 18 19	Claim (6)	As a matter of fact, Ernie, the correlation between CO ₂ levels in the atmosphere and the pH of the oceans makes sense according to chemistry.
20 21 22 23 24 25 26 27	Consequences (7)	When this happens, calcifying organisms will become weaker, such as coral. They will be significantly affected as a result and may be unable to live in the current environment in which they live. In addition to ocean acidification wreaking havoc on the environment, other factors such as climate change can do the same thing and increase the amount of damage that is done to it.
28 29 30 31	Cause and Effect (8)	In each equation, the amount of reactants increases, which drives the reaction forward, meaning the amount of products will increase until the reaction reaches equilibrium.
32 33 34 35 36 37 38 39		In accordance with Le Châtelier's Principle, increasing the amount (or the concentration) of atmospheric CO ₂ will shift this equation towards the dissolved CO ₂ in the ocean to make up for the increase in gas on the reactants (left) side. This dissolved CO ₂ is indicated as 'CO ₂ (aq)' in the equation, which stands for 'aqueous CO ₂ .' Consequently, the dissolved CO ₂ relates to the bicarbonate formation equation:
40 41 42 43		$\text{CO}_2 (aq) + \text{H}_2\text{O} (l) \rightleftharpoons \text{H}_2\text{CO}_3 (aq) \rightleftharpoons \text{H}^+ + \text{HCO}_3 (aq)$ (Doney et al., 2009 [<i>reference provided to student</i>])
44 45 46 47 48 49 50 51 52 53 54 55 56	Deduction (9)	Just as increasing the concentration of the atmospheric CO ₂ caused the equilibrium between dissolved CO ₂ to favor the formation of bicarbonate (HCO ₃) and hydrogen, two byproducts of carbonic acid (H ₂ CO ₃). Increasing dissolved CO ₂ concentration—a result of increased atmospheric CO ₂ —will 'push' the equilibrium towards the formation of H ₂ CO ₃ . Likewise, a shift towards the formation of H ₂ CO ₃ will also shift the equation forward towards the formation of protons, H ⁺ , and HCO ₃ . Finally, these products of proton and bicarbonate will be in equilibrium with two protons and carbon trioxide (CO ₃ ²⁻):
57 58 59 60		$\text{H}^+ + \text{HCO}_3 (aq) \rightleftharpoons 2\text{H}^+ + \text{CO}_3^{2-}$ (Doney et al., 2009)

Once more, an increase in the H^+ and HCO_3^- concentration will push the equilibrium forward towards the formation of two $2H^+$ and CO_3^{2-} .

Argumentation
(10)

So using this information, you can now look at the graph you posted and understand the relationship between CO_2 levels and the pH of the water. As CO_2 is absorbed into the water, it produces H^+ ions, which then cause the pH of the water to decrease. You can see this trend on the graph. Even though they do not seem like they should be related in any way, a change in one would cause a change in the other. Something that is making the line representing the CO_2 in the atmosphere on the graph to increase so much is the amount of CO_2 humans emit every day. Whenever you drive a car you are releasing CO_2 into the atmosphere. Since there is so much more CO_2 being released, the oceans are absorbing more CO_2 . In fact, the oceans have absorbed almost 30% of the CO_2 humans have emitted since the Industrial Revolution. As more CO_2 is absorbed by the oceans because of human activity, the more H^+ ions are formed, and the more the pH of the ocean is decreased.

A total of 296 assignments have been coded (average weighted complexity: 6.3; average number of operations per essay: 9). We determined that saturation had been reached when no new codes or patterns were observed in the writing. Table 4 shows the frequency of operations used in the 296 essays. *Observation* was the most frequently used operation by students, with many using more than one observation in a single essay. Students heavily relied on making statements about how a variable was changing to counter 'Ernie's' claim. In this context, this means that students were able to understand how variables were changing from the graph provided. There was very little use of *Example* or *Measurement*. *Argumentation* was also used relatively infrequently. As mentioned above, *argumentation* was distinct as its own operation given indistinguishable use of multiple lower complexity operations. For this reason, *argumentation* required students to combine multiple operations (and domains). This difficulty combined with the infrequent use suggests that *argumentation* was indeed the most complex operation. The high frequency of *claim*, *cause and effect*, and *deduction* is likely tied to this writing context. Students were trying to convince Ernie (*claim* and *cause and effect*) by invoking chemical principles (*deduction*). It is expected that the

distribution of use will vary with different writing contexts, depending upon what the prompt elicits.

Table 4. Descriptive information from application of the framework to the data set presented herein.

Operations	Frequency
Observation (2)	411
Claim (6)	348
Cause and effect (8)	331
Definition (1)	312
Comparison (4)	302
Deduction (9)	235
Consequences (7)	146
Argumentation (10)	63
Example (5)	51
Measurement (3)	34

Research Question 2: What features of students' argumentative writing do cognitive operations serve to explain?

To determine what exactly this framework served to characterize, two comparisons were made. The first comparison was between framework estimates of complexity and instructor estimates of quality. This comparison was intended to demonstrate that this framework was telling us something that faculty would normally have to make a judgment about. Further this comparison was intended to reveal that this framework could make *similar* judgments to an instructor. Table 5 shows how five faculty ranked four assignments, and this is compared to our framework ranking. Instructors were tasked with ranking the assignments according to the quality of the scientific reasoning (as they saw fit to evaluate it). This approach was taken so as to elicit instructors' "gut reaction"—the kind of evaluation they would make if they were grading this sort of task in their class. Instructor rankings reveal a few trends. First and potentially the most

important, there is almost complete consensus on the ‘best’ essay (J) with the exception of Instructor 3. This finding speaks to the framework’s capacity to identify the best. Essay J received a high cognitive complexity score because of the presence of extended argumentation, which aligns with what instructors are valuing when evaluating scientific reasoning. Four of the five instructors ranked Essay G as second best, with the exception of Instructor 3, even though our framework estimates it as second from the worst. Finally, all five instructors consider Essay H and K to be the worst, where as our framework estimates Essay H to be the second best. The difference between instructor rankings for Essays H, G, and K and framework estimates illustrate an important limitation of the framework. Essay H contained a misconception in which the student claimed a relationship between atmospheric temperature and ocean acidification though there was no data regarding heat for any of the chemical reactions provided. However, this student’s reasoning was rather sophisticated, with multiple high complexity cognitive operations. Our framework does not account for scientific accuracy of students’ essays. We chose to include this essay as it authentically represents what an instructor might encounter with grading writing. One possible explanation for the ranking difference is that for instructors scientific content and reasoning are inextricably linked, which is consistent with feedback from one instructor who explained that content accuracy factored into their ranking. It is for this reason that we also analysed the scientific accuracy of students’ essays (see Table 5 below).

Table 5. Instructor ranking of four assignments of varying quality compared to framework estimates of ranking

Essay name	Cognitive complexity ranking (4 worst, 1 best)	Expert rankings (4 worst, 1 best)				
		Instructor 1	Instructor 2	Instructor 3	Instructor 4	Instructor 5
J	1	1	1	2	1	1
H	2	3	3	4	4	4
G	3	2	2	1	2	2

K 4 4 4 3 3 3

The second comparison made was between cognitive complexity—framework estimates of quality—and common student characteristics that are frequently used as measures or predictors of success (Hein & Smerdon, 2013). The purpose of this comparison was to determine if characterizing the quality of students' reasoning in this way was revealing something about students that could have been predicted by a metric that was already collected (e.g., ACT math score). In other words, this framework is useful only in so far as it tells us something interesting about students that other metrics do not. Table 6 shows the correlations between common student characteristics and cognitive complexity. There were no significant correlations between cognitive complexity and any common characteristics, which would not be expected for measures of constructs distinct from that captured by this framework (i.e. math). These findings may mean that the framework captures something distinct from what is measured by other standardized tests. A strong negative correlation exists between the number of operations used and the cognitive complexity. This means that students with higher cognitive complexity essays used fewer moves, which could indicate a synthesis of ideas in order to produce higher complexity operations.

Table 6. Pearson correlations between student characteristics and cognitive complexity (p-values reported for t-tests used for categorical variables: gender and ethnicity [white and non-white students compared])

Variables	Cognitive complexity
Number of operations	-0.649*
Final exam grade	-0.018
Final course grade	-0.025
CHEM placement	-0.081
MATH placement	-0.020
ACT math [†]	0.003
Current GPA	-0.062
Cumulative GPA	-0.060
Gender	0.401
Ethnicity	0.071

* indicates p (two-tailed) < 0.01

† For students with only SAT math scores, their scores were converted to ACT math scores using contingency tables

Research Question 3: What is the relationship between framework characterizations of complexity and conceptual correctness?

The data above reveal that the cognitive operations framework is characterizing students' reasoning in a way that other measures do not. However, the instructor rankings reveal that there exists a relationship between reasoning and accuracy. The motivation for considering this relationship partially sources from the concern that any information this framework provides is irrelevant if students are largely scientifically inaccurate. In order to explore this relationship, we coded all data that had already been coded according to operations for 'correctness.' That is, when scientifically inaccurate information was identified in an essay, the cognitive operation containing that information was marked as incorrect. In this way, all student writing was coded for both correctness and cognitive function (i.e., content and structure as highlighted in theoretical framework). Table 7 shows the number of incorrect operations per the total number of cognitive operations. Further, there were no correlations between the cognitive complexity and number of incorrect operations or between the number of cognitive operations and the number of incorrect operations. This finding suggests that overall, producing a more complex essay does not make it more likely that a student will use more incorrect ideas, but as Table 7 shows, there may be specific operations that elicit more incorrect ideas. Further, writing more operations, or introducing more separate idea units, does not make a student more likely to put forth incorrect ideas.

Table 7. Number of incorrect cognitive operations relative to total number of operations [def.=definitions, obs.=observation, meas.=measurement, comp.=comparison, ex.=example, claim=claim, cons.=consequences, C&E=cause and effect, Ded.=deduction, Arg.=argumentation]

	Def.	Obs.	Meas.	Comp.	Ex.	Claim	Cons.	C&E	Ded.	Arg.
# incorrect	9	11	0	13	1	8	6	40	18	4
# operations	312	411	34	302	51	348	146	331	235	63

% incorrect per total ops.	3	3	0	4	2	2	4	10	8	6
----------------------------------	---	---	---	---	---	---	---	----	---	---

Evident in Table 7 is a relatively infrequent use of scientifically inaccurate information. That is, given that our unit of analysis is ideas, students are largely generating scientific ideas employing correct scientific information. This further justifies the move beyond simply considering scientific accuracy of students' conceptions towards considering the sophistication of their reasoning about scientific ideas. In this case, only considering the accuracy would have provided a very limited picture of what these students were doing in their writing. Because of the relative infrequency, it became important to consider the nature of the inaccuracies. For this inquiry, categorizing the inaccuracies by operation led to an interesting finding. The highest percentages of inaccuracy, though still relatively small, occurred with *cause and effect*, *deduction*, and *argumentation*. It is possible that higher complexity operations surface alternative conceptions more effectively. Further, the alternative conceptions elicited are potentially more deeply held, keeping in mind the Limited or Inappropriate Propositional Hierarchies (LIPs). That is, higher complexity operations draw on multiple domains and elements and may have the potential to reveal more of students' mental structures, and thus, expose LIPs.

Limitations

Though this framework provides a useful way to evaluate students' written work, it has a number of limitations. First, as noted above, this framework does not capture the scientific accuracy of students' written ideas. The utility of this tool, then, is limited to a narrower research goal—characterizing students' reasoning. When combined with an analysis of the content accuracy, however, this framework can provide unique insights about students' understanding. Further, this framework was conceptualized, tested,

1
2
3 refined, and ultimately applied to a corpus of writing in a very specific context—general
4 chemistry argumentative writing about ocean acidification. It is possible that some of
5 the ways that cognitive operations have been conceptualized in this study are specific to
6 this context. For this reason, applications to other contexts are needed to ensure the
7 domain-general nature of this framework. Finally, due to the relatively low occurrence
8 of certain operations in this context, we have a weaker understanding of some of the
9 operations (i.e., measurement). Because of the complete absence of the inductive
10 reasoning operation from Grimberg and Hand’s original framework in this set of student
11 writing, it was not included in this application, even though it is likely to be employed
12 in other contexts. Finally, this data was collected at a selective institution and it is likely
13 that different incorrect ideas or reasoning patterns would emerge from other student
14 populations. Again, this can be addressed by applying this framework to student writing
15 in other contexts.

34 **Discussion and Implications**

36 The first research question posed in this work considered how a cognitive operations
37 framework can be used to characterize students’ reasoning evident in their
38 argumentative writing. In this article, we show what this framework is like and how it
39 can be applied to students’ writing. We refined a list of cognitive operations generated
40 by Grimberg and Hand (2009) and organized them according to complexity, and then
41 used these operations to code general chemistry students’ writing on ocean
42 acidification. This framework has some key affordances that make it useful to both
43 research and practice. It is domain general, which means that it can be applied to writing
44 in a variety of contexts. We recommend, then, that others apply this to writing in a
45 variety of contexts across STEM and across levels (introductory to advanced student
46 populations). The domain-general nature of this potentially enables the identification of

1
2
3 differences in students' reasoning across disciplines and levels. For example, do
4
5 advanced students employ more complex reasoning than introductory students?
6

7
8 Another affordance of this framework is the 'score' that is a product of
9
10 application—the cognitive complexity. The single score output provides an estimate of
11
12 construct that is rather difficult to measure —student reasoning. This framework, then,
13
14 can potentially overcome some of the difficulties with evaluating writing reported in the
15
16 literature (Hamp-Lyon, 2016; Neill, 2002). This framework provides a novel approach
17
18 to assigning a holistic score to writing. Further, the use of cognitive operations enables
19
20 the identification of patterns in students' writing. That is, it can be used to characterize
21
22 the movement between cognitive operations and the likelihood of moving towards high
23
24 complexity operations, as shown in Grimberg and Hand's original application (2009).
25
26 This framework's capacity to capture temporal patterns makes it very useful for
27
28 understanding how students reason in extensive writing (Grimberg & Hand, 2009; Kelly
29
30 & Takao, 2002; Moreira et al., 2018).
31
32
33
34

35
36 The second research question aimed to elucidate what features of student
37
38 thinking were understandable with this framework. That is, what does this framework
39
40 evaluate the quality of? This was achieved in two ways. The first was to compare
41
42 framework estimates to instructor estimates of quality. This approach was intended to
43
44 determine if the framework estimates were similar to the instructor estimates and if both
45
46 were evaluating a similar construct. This revealed that perhaps for the upper bound of
47
48 the construct—argumentation—there was agreement between instructors and
49
50 framework estimates. There was less agreement for the other-than-best essays. Kelly
51
52 and Takao (2002) identified similar disparities between their framework estimates and
53
54 expert rankings and explained them as common occurrences when evaluating writing
55
56 (Wolcott & Legg, 1998). In our case, we argue that the variety was an artefact of the
57
58
59
60

1
2
3 presence of inaccurate scientific information in one of the essays. Instructors may not
4 separate content and reasoning as this framework does. However, we argue, similar to
5
6 Kelly and Takao (2002), that this framework may provide a tool for evaluating the
7
8 validity of instructor's estimates of quality. More research is necessary to establish
9
10 interrater reliability amongst instructor ratings and identify ways in which the
11
12 framework can serve as a tool for supporting instructors in systematically assessing
13
14 students' writing.
15
16
17
18

19 To determine if this framework was providing unique information about
20
21 students' ability, we compared cognitive complexity to other common performance
22
23 measures. There were no correlations. We posit two potential explanations for this. The
24
25 first is that this metric of cognitive complexity is indeed measuring something unique
26
27 from what typical performance metrics measure (National Research Council, 2001).
28
29 The second is that students who perform well on typical performance metrics do not
30
31 necessarily perform equally well on more extensive writing tasks (National Research
32
33 Council, 2001). Both of these explanations warrant further investigation because of the
34
35 implications for assessment. Specifically, this framework could serve to equip the
36
37 evaluation of more interesting competencies in students than that measured by typical
38
39 performance measures *or* assignments of this nature could serve to minimize advantages
40
41 certain groups bring with them to typical performance measures. However, we also
42
43 recognize that there may be other performance measures that correlate with the
44
45 framework estimate. Particularly, we would expect that more generative or authentic
46
47 assessments might correlate more strongly with cognitive complexity (National
48
49 Research Council, 2001). Finally, we aimed to characterize the relationship between
50
51 framework estimates of quality and scientific accuracy. In order to do this, we analysed
52
53 writing for the presence of scientific inaccuracies and coded the respective operation in
54
55
56
57
58
59
60

1
2
3 which they appeared. This revealed that scientific inaccuracies occurred relatively
4 infrequently with about 10 percent of *cause and effect* operations including something
5 that did not agree with scientifically accepted knowledge. The percentage among *cause*
6 *and effect* operations was the highest. However, there appears to be a trend in which
7 higher complexity operations (i.e., cause and effect, deduction, and argumentation) had
8 higher frequencies of incorrect information than low complexity operations. We argue,
9 in light of Novak's work on LIPHS, that higher complexity operations as a
10 representation of students' mental models may better reveal LIPHS (Novak, 2002). That
11 is, employing more complex reasoning may surface more deeply held LIPHS. Students
12 who do not use higher complexity operations may be more limited in both their and
13 their instructors' capacity to address potential alternative conceptions. The relationship
14 between complexity and conceptual correctness warrants further investigation.
15 Understanding this relationship is important for designing formative assessments that
16 better elicit high complexity operations.

17
18
19 This framework also offers some unique implications for instructors who assign
20 similar tasks to their students. Scoring assignments in this way could permit an
21 instructor to draw conclusions about their students' collective access to complex
22 reasoning operations. For example, a low average score of cognitive complexity in their
23 course may motivate instructors to explicitly address and model complex reasoning
24 types for their students. However, we argue that the most important implication of this
25 framework for practice is providing a vocabulary to instructors for giving tailored
26 feedback to students. That is, applying this sort of framework would support an
27 instructor to give specific examples of when a student could have employed more
28 complex reasoning appropriately and instead used a less complex operation.

1
2
3
4
5
6
7 **Disclosure statement**
8

9
10 No potential conflict of interest was reported by the authors.
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 **References**

- 55 Emig, J. (1977). Writing as a Mode of Learning. *College Composition and*
56 *Communication*, 28(2), 122–128.
57
58 Gere, A. R., Limlamai, N., Wilson, E., MacDougall Saylor, K., & Pugh, R. (2019).
59 Writing and Conceptual Learning in Science: An Analysis of Assignments. *Written*
60 *Communication*, 36(1), 99–135. <https://doi.org/10.1177/0741088318804820>

- 1
2
3 Grimberg, B. I., & Hand, B. (2009). Cognitive Pathways : Analysis of students ' written
4 texts for science understanding. *International Journal of Science Education*, 31(4),
5 503–521. <https://doi.org/10.1080/09500690701704805>
6
7 Gunel, M., Hand, B., & Prain, V. (2007). Writing for learning in science: A secondary
8 analysis of six studies. *International Journal of Science and Mathematics*
9 *Education*, 5, 615–637. <https://doi.org/10.1007/s10763-007-9082-y>
10
11 Ha, M., Nehm, R. H., Urban-Lurain, M., & Merrill, J. E. (2011). Applying
12 computerized-scoring models of written biological explanations across courses and
13 colleges: Prospects and limitations. *CBE Life Sciences Education*, 10(4), 379–393.
14 <https://doi.org/10.1187/cbe.11-08-0081>
15
16 Hein, V., & Smerdon, B. (2013). *Predictors of Postsecondary Success. College and*
17 *Career Readiness and Success Center at American Institutes for Research.*
18
19 Kelly, G. J., & Bazerman, C. (2003). How Students Argue Scientific Claims : A
20 Rhetorical-Semantic Analysis. *Applied Linguistics*, 24(1), 28–55.
21
22 Kelly, G. J., Chen, C., & Prothero, W. (2000). The Epistemological Framing of a
23 Discipline : Writing Science in University Oceanography, 37(7).
24
25 Kelly, G. J., Regev, J., & Prothero, W. (2007). Analysis of Lines of Reasoning in
26 Written Argumentation. In *Argumentation in Science Education* (pp. 137–157).
27
28 Kelly, G. J., & Takao, A. (2002). Epistemic levels in argument: An analysis of
29 university oceanography students' use of evidence in writing. *Science Education*,
30 86(3), 314–342. <https://doi.org/10.1002/sce.10024>
31
32 Keys, C. W. (1994). The development of scientific reasoning skills in conjunction with
33 collaborative writing assignments: An interpretive study of six ninth grade
34 students. *Journal of Research in Science Teaching*, 31(9), 1003–1022.
35 <https://doi.org/10.1002/tea.3660310912>
36
37 Klein, P. D. (1999). Reopening Inquiry into Cognitive Processes in Writing-To-Learn.
38 *Educational Psychology Review*, 11(3), 203–270.
39 <https://doi.org/10.1023/A:1021913217147>
40
41 Klein, P. D. (2015). Mediators and Moderators in Individual and Collaborative Writing
42 to Learn. *Journal of Writing Research*, 7(1), 201–214.
43
44 Krippendorff, K. (2004). Reliability in Content Analysis : Some Common
45 Misconceptions and Recommendations Reliability in Content Analysis : Some
46 Common Misconceptions and. *Human Communication Research*, 30(3), 411–433.
47
48 Liu, O. L., Rios, J. A., Heilman, M., Gerard, L., & Linn, M. C. (2016). Validation of
49 automated scoring of science assessments. *Journal of Research in Science*
50 *Teaching*, 53(2), 215–233. <https://doi.org/10.1002/tea.21299>
51
52 Moon, A., Gere, A. R., & Shultz, G. V. (2018). Writing in the STEM classroom:
53 Faculty conceptions of writing and its role in the undergraduate classroom. *Science*
54 *Education*, 0(0). <https://doi.org/10.1002/sce.21454>
55
56 Moreira, P., Marzabal, A., & Talanquer, V. (2018). Using a mechanistic framework to
57 characterise chemistry students' reasoning in written explanations. *Chemistry*
58 *Education Research and Practice*. <https://doi.org/10.1039/C8RP00159F>
59
60 National Research Council. (2001). *Knowing what students know: The science and*
design of educational assessment. National Academies Press. Washington DC.
<https://doi.org/10.17226/10019>
Novak, J. D. (2002). Meaningful Learning: The Essential Factor for Conceptual Change
in Limited or Inappropriate Propositional Hierarchies Leading to Empowerment of
Learners. *Science Education*, 86(4), 548–571. <https://doi.org/10.1002/sce.10032>
Prain, V., & Hand, B. (2016). Coming to Know More Through and From Writing.
Educational Researcher, 45(7), 430–434.

- 1
2
3 <https://doi.org/10.3102/0013189X16672642>
4 Reynolds, J. A., Thaiss, C., Katkin, W., & Thompson, R. J. (2012). Writing-to-learn in
5 undergraduate science education: A community-based, conceptually driven
6 approach. *CBE-Life Sciences Education*, 11(1), 17–25.
7
8 Sandoval, W. A. (2003). Conceptual and Epistemic Aspects of Students ' Scientific
9 Explanations. *Journal of the Learning Sciences*, 12(1), 5–51.
10 <https://doi.org/10.1207/S15327809JLS1201>
11 Sandoval, W. A., & Millwood, K. A. (2005). The Quality of Students ' Use of Evidence
12 in Written Scientific Explanations. *Cognition and Instruction*, 23(1), 23–55.
13 <https://doi.org/10.1207/s1532690xci2301>
14 Sandoval, W. A., & Reiser, B. J. (2004). Explanation-driven inquiry: Integrating
15 conceptual and epistemic scaffolds for scientific inquiry. *Science Education*, 88(3),
16 345–372. <https://doi.org/10.1002/sce.10130>
17
18 Sevian, H., & Talanquer, V. (2014). Rethinking chemistry: a learning progression on
19 chemical thinking. *Chemistry Education Research and Practice*, 15(1), 10–23.
20 <https://doi.org/10.1039/C3RP00111C>
21 Takao, A. Y., & Kelly, G. J. (2003a). Assessment of Evidence in University Students '
22 Scientific Writing. *Science & Education*, 12, 341–363. <https://doi.org/10.1023/A>
23 Takao, A. Y., & Kelly, G. J. (2003b). Assessment of Evidence in University Students '
24 Scientific Writing. *Science & Education*, 12, 341–363.
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