# Finding Patterns in and Refining Characterizations of Students' Epistemic Cognition: A Computational Approach

*Christina (Stina) Krist, Northwestern University*
*Joshua Rosenberg, Michigan State University*

**Abstract:** Epistemic cognition has become an area of interest among education researchers. However, characterizing epistemic elements of cognition is difficult. Using computational approaches can contribute to conceptualizing and operationalizing what we mean by epistemic cognition. We present how we used two computational methods to see new dimensions of epistemic elements in students' written work. Our work demonstrates how computational tools refine our understanding of, and tools to examine, dimensions of epistemic cognition.

## Introduction

Epistemic cognition has increasingly become an area of interest among education researchers, particularly in understanding and supporting student's use of epistemic criteria for building knowledge (e.g., Chinn, Buckland, & Samarapungavan, 2011). However, characterizing epistemic elements of cognition is difficult. Computational approaches to data analysis have been applied for educational data mining and learning analytics, but are not yet commonly used as a part of Learning Sciences research. To explore the use of such methods for characterizing elements of epistemic cognition, we present a combination of approaches for identifying students' epistemic ideas. To do so, we build on Sherin's (2013) and others' (i.e., Beggrow, Ha, Nehm, Pearl, & Boone, 2013) arguments for the use of computational methods as a compliment to traditional qualitative methods. First, we present our use of two natural language processing approaches to analyze students' written work about their modeling and explanation-building tasks embedded within the curriculum ("embedded assessments"). We demonstrate how using these methods provides conceptual benefits: they help us identify gaps in our coding scheme and refine our conceptual models of students' epistemic cognition. We present this work-in-progress as a contribution to the ongoing conversation around learning analytics in Learning Sciences research.

## Method

In order to analyze the epistemic criteria that students used to build ideas, we collected embedded assessments designed to elicit students' consideration of the *Generality* of their diagrammatic models or written explanations constructed as a regular part of their science classroom activities. We asked whether their model or explanation should explain phenomena in general (e.g., a general way that chemical reactions occur) or just the specific situation on which they focused (e.g., how and why aluminum and copper chloride react), and why. We selected this criterion for our initial computational analyses because we had a relatively complete set of hand-coded student responses from a $7^{th}$ grade chemistry unit ($n = 178$). We adapted our initial coding scheme (Table 1) to differentiate between responses that discussed the rationale for a general or specific account without (1) and with (2) a rationale, and those that addressed both general and specific elements and/or trade-offs (4a and 5). Over time, increasing scores would indicate that students are developing an understanding that scientific work involves building general principles from specific phenomena (Godfrey-Smith, 2003).

Table 1: Original Coding Scheme for Generality

| Code | Description | Sample Response(s) |
|---|---|---|
| 0 | Unclear/not codeable | Because the cemig rauted (101466) |
| 1 | One level: No rationale | "Because that's the whole point of the model" (101309) |
| 2 | One level: Rationale | "That can explain the atoms [rearrangement] better" (101536) |
| 4a | Level-crossing | "This should help with all open systems in general because we know that if this happens with other reactants, the atoms would still leave in an open system" (104148) |
| 5 | Level-crossing: Boundary conditions of Gen and/or Spec | "Because the question is asking about only 1 specific thing and if I talked about all the chemical reactions in general it would not make sense because in different chemical reactions different things happen, like bubble, smell" (101323) |

We then used two natural language processing approaches that take advantage of patterns in relative frequencies and arrangements of words in a set of documents. First, we used a supervised (and therefore requiring hand-

coded data) approach to "train" a Naive Bayes classifier on the initial coding scheme. This was moderately successful but not convincing enough to use more widely. Next, we used an unsupervised or inductive approach that converts the text into vectors and groups them using centroid clustering (see Sherin, 2013). We reasoned that this approach would help identify patterns similar to those that a trained classifier would be likely to identify, which could then give us some leverage in "training" an automated classifier in the future.

## Findings and Discussion

Here, we compare the clusters of student responses to our assignment of our original codes to illustrate how this analysis is leading us to revise our coding scheme and eventually improve on our earlier supervised approach. For the 178 responses we analyzed, we found that a 9-cluster solution balanced interpretability and concerns of parsimony. We graphed the clustered responses by their manual code according to their new cluster (Figure 1). We found that some of these clusters corresponded to levels of the coding scheme. For example, most of the responses that clustered into the "Similarities and Comparisons Across Processes and Classes" topic, which had been hand-coded as either a 4a or a 5, which is conceptually similar.
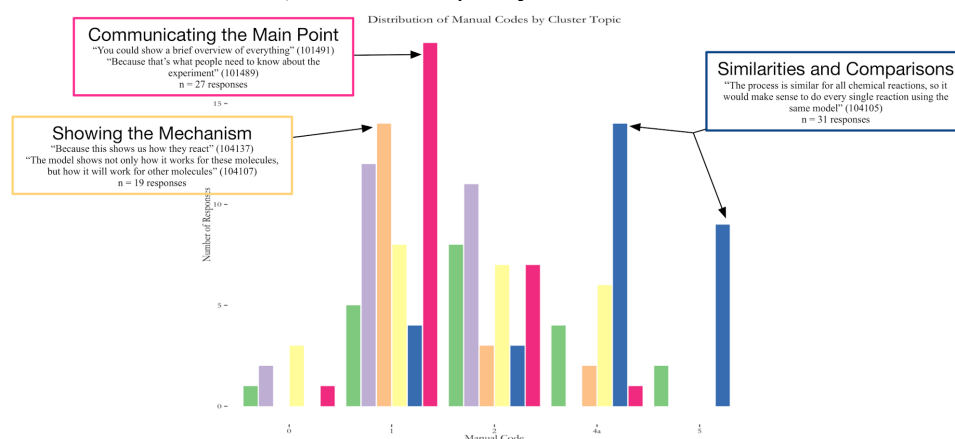


Figure 1: Distribution of Manual Codes by Cluster Topic.

In contrast, other responses did not correspond between manual code and cluster. For example, responses that we coded manually as 1 are the main responses in the clusters with the topics "Communicating the Main Point" and "Showing the Mechanism." The responses in these clusters differ as well: while both groups provide a rationale related to the intent of communicating some information, the rationales from students who communicated the main point are agnostic towards the substance of the main point, while rationales from students who show the mechanism focus on the importance of communicating mechanistic processes. This subtle difference may have important implications for how students' thinking about *Generality* develops.

These examples illustrate how computational methods may serve as tools to refine our coding scheme (and underlying conceptual models) as well as enhance the reliability of our analysis. Our moderately successful work with automated classification indicated that our initial coding scheme was not yet tailored to capture the patterns in our data. In addition, considering how the computer inductively clustered responses helps us to improve the already-coded data to train and use to a classifier. As computational tools become more widely used, we see great potential for them to support researchers' efforts to understand and develop ill-defined or difficult to capture constructs.

## References

Beggrow, E. P., Ha, M., Nehm, R. H., Pearl, D., & Boone, W. J. (2014). Assessing scientific practices using machine-learning methods: How closely do they match clinical interview performance? *Journal of Science Education and Technology*, 23(1), 160-182.

Chinn, C. A., Buckland, L. A., & Samarapungavan, A. L. A. (2011). Expanding the dimensions of epistemic cognition: Arguments from philosophy and psychology. *Educational Psychologist*, *46*(3), 141-167.

Godfrey-Smith, P. (2009). *Theory and reality: An introduction to the philosophy of science*. University of Chicago Press: Chicago, IL.

Sherin, B. (2013). A computational study of commonsense science: An exploration in the automated analysis of clinical interview data. *Journal of the Learning Sciences*, *22*(4), 600-638.