

Quantitative Analysis of Indicators on the RTOP and ITC Observation Instruments

Martha A. Henry, Keith S. Murray, Mark Hogrebe, & Marcia Daab

Abstract: Classroom observation is an important component of mathematics and science teacher professional development programs and other educational evaluation activities. In this paper, the authors extend their earlier qualitative examination and comparison of two popular classroom observation tools (Horizon Research, Inc.'s *Inside the Classroom Observation and Analytic Protocol* (ITC) and Sawada et al.'s *Reform Teacher Observation Protocol* (RTOP)) with a quantitative analysis. The instruments substantially are based on comparable assumptions, foundational philosophies and domains of interest, and appear to be used with an expectation that they yield similar results, not to mention that their domains adhere internally. A single well-trained educator observed 21 teachers from a Mathematics and Science Partnership over the course of two years. Pearson Correlation Coefficient analysis was applied to items across and within instruments. Because of the relatively small number of teachers and the use of one rater, the authors applied a relatively strict interpretation ($\geq .75$) of high correlations. While some items correlated as expected within and between instruments based on domain and item construction, many items lacked matches, including those with an apparently similar focus. RTOP items showed greater alignment than did ITC items. Among summative "synthesis" domain ratings in the ITC, none of the four aligned with any of the specific ITC item indicators within their own domain categories. Subtle differences in wording and implicit differences in overall focus especially appear to restrict both internal matches and matches between superficially comparable items across the two instruments. These findings complicate the interpretation of observational results and challenge assumptions that the instruments are interchangeable or necessarily internally consistent. Additional research and development of observation instruments is needed, and users of existing instruments must carefully assess their own needs and understandings before attempting to draw conclusions about classroom practice based on them.

Nationally funded science and mathematics teacher development projects typically incorporate a classroom component in their treatment activities. These classroom components serve as the logical culmination for these projects' activities, given that enhanced teacher classroom performance and increased student achievement represent their ultimate aims. Evaluation of project-related classroom activities therefore comprises a crucial element of understanding project impact.

The authors have served as evaluators, researchers or developers for several of these projects, including four National Science Foundation (NSF) Math and Science Partnerships (MSPs) and four Department of Education state MSPs. In trying to select the instrument for classroom observations that most accurately represents implementation of reform-based mathematics and science teaching, the authors have analyzed numerous observation instruments. Being the two most widely used or adapted instruments, the *Reform Teacher Observation Protocol* (RTOP) (Sawada et al., 2000) and the *Inside the Classroom Observation and Analytic Protocol* (ITC) (Horizon Research, Inc., 2002) were examined closely. Results of this first examination were released as a qualitative analysis of the individual items for similarities and differences across the two instruments (Henry, Murray & Phillips, 2007). The previous paper also provides a detailed description of the instruments, which will assist in following the analyses contained here.

The present paper represents a second look, comprising a quantitative analysis of the indicators on these two instruments. Its purpose is to statistically clarify the matches made during the qualitative analysis and, where appropriate, adjust results of the previous alignment. Ultimately, an increased understanding of what each instrument measures and does not measure will assist evaluators in selecting and using the most appropriate instrument for their data collection processes.

In preparing this paper, the authors do not suggest that the developers of the instruments intended any use beyond that specifically noted in their manuals and training materials. However, practically speaking, educational professionals seeking to implement an observation component appear likely to draw from instruments that are commonly used and that appear to correspond to their domains of interest. As this and the previous qualitative paper demonstrate, a closer look can reveal not just a lack of thorough validation and other testing, but subtle differences in construction that can affect scores in unexpected ways and reduce comparability across apparently comparable items.

Procedure and Results

Twenty-one teachers in a Mathematics and Science Partnership were observed twice by the same observer over the course of two years. A single observer was used to eliminate reliability concerns and ensure the greatest likelihood of consistency, in keeping with the aims of the research. The observer holds a doctorate in education and has considerable experience in science and math teaching, supervision, professional development, assessment, and classroom observation with the ITC and RTOP instruments.

Immediately following the observations, both the RTOP and the ITC were completed for each teacher, in no particular order. Pearson Correlation Coefficient analysis was applied to items across and within instruments. Because of the relatively small number of teachers and the use of one rater, a strict interpretation of correlation has been applied. Items with a correlation of .75 or above are considered to measure the same qualities of classroom activities. Items with correlations from .61 to .74, usually considered moderately high correlations, will need further study to establish confidence in their apparent relationships. Correlations below .60 were omitted in this analysis.

Expected Similarities

Both instruments, being based on the same foundational documents and concepts, contain similar language to represent constructs appearing in the instruments. Constructs such as “design,” “content,” and “culture” would be expected to reflect similar indicators. As such, one might expect to see the items loading in this manner (Table 1).

Table 1. Expected Item Loading Across Instruments

| RTOP Domains | ITC Domains | | | |
|--|-------------|----------------|---------------------------------|-------------------|
| | Design | Implementation | Mathematics/ Science Content | Classroom Culture |
| Design and Implementation | | | | |
| Content: Propositional Knowledge | | | | |
| Content: Procedural Knowledge | | | | |
| Classroom Culture: Communicative Interactions | | | | |
| Classroom Culture: Student Teacher Relationships | | | | |

A closer look at highly correlated items (Table 2) shows where the similarities in these instruments actually occur. (Items with relatively low or no correlations, within the limits already cited, have been omitted from Table 2.)

Table 2. RTOP and ITC Items Showing High ($\geq .75$) Correlation

| RTOP Domains/Items | ITC Domains/Items | | | | | | | | | | | | | | | |
|---|-------------------|-------|----------------|-------|-------|-------|-------|---------|-------|-------|-------|-------|---------|-------|-------|-------|
| | Design | | Implementation | | | | | Content | | | | | Culture | | | |
| | I.1 | I.S | II.1 | II.3 | II.6 | II.7 | II.S | III.4 | III.6 | III.7 | III.9 | III.S | IV.1 | IV.3 | IV.5 | IV.S |
| Design and Implementation | | | | | | | | | | | | | | | | |
| 1 | | | | | 0.783 | | | | | | | | | | | |
| 2 | | 0.820 | | | | | | | | | | | | | | |
| 3 | | | | | | | | | 0.778 | | | | | | | |
| Content: Propositional Knowledge | | | | | | | | | | | | | | | | |
| 7 | | | | | | | 0.780 | | | | | | 0.760 | | | |
| 9 | | | | | | | | | 0.833 | | | | | | | |
| Content: Procedural Knowledge | | | | | | | | | | | | | | | | |
| 13 | 0.755 | | 0.815 | | | | | | | | | | | | | |
| 14 | | | | | | 0.752 | 0.787 | | 0.809 | | 0.795 | | | | 0.794 | |
| 15 | | | | | | | | | | | | | 0.826 | | | |
| Classroom Culture: Communicative Interactions | | | | | | | | | | | | | | | | |
| 17 | | | | | | | | | 0.750 | | | | | | | |
| 18 | | | | | | | | | | | | | | 0.784 | | |
| 19 | | | | | | | | | | | | 0.754 | | | | |
| 20 | | | | | | | | | | | | | 0.827 | | | |
| Classroom Culture: Student Teacher Relationships | | | | | | | | | | | | | | | | |
| 21 | | | | | 0.781 | | 0.844 | 0.776 | 0.819 | | 0.760 | | | | 0.819 | 0.772 |
| 22 | | 0.799 | 0.804 | | | | 0.757 | | 0.889 | | | | | | 0.752 | |
| 23 | | | | 0.816 | | | 0.767 | 0.853 | | | | | | | | 0.782 |
| 24 | | 0.766 | 0.805 | | 0.853 | 0.834 | | | 0.770 | | | | | | 0.767 | |
| 25 | | | 0.896 | | 0.819 | | | | | | | | | | | |

The items in Table 2 do not fall entirely within the expected conceptual categories indicated in Table 1. In Table 2, the broader categories of Classroom Culture: Student Teacher Relationships (RTOP) load on many unexpected factors, including content, design and procedural knowledge. With the exception of item 25, the other Student/Teacher Relationships items load on items in at least three of the four ITC categories. Because of their lack of discrimination, RTOP items 21

through 25 will not be considered in the cross-instrument analysis. Care should be taken in the interpretation of this RTOP category because of its multiple loading.

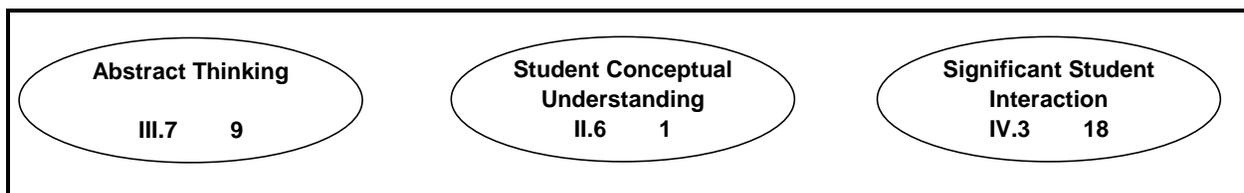
One problem arose with the Design indicators in both of the instruments. It is difficult to “see” Design during a classroom observation; it is often inferred from how the lesson progresses. Even then, one must determine whether what is seen is actually what was planned. Unless the lesson is provided in written form, the design indicators often are rated as the result of something observed and inferred as the lesson plays out.

Aligned Items

Single Item Correlation

A closer examination of those items that correlate reveals additional concerns. Indicators that only loaded on one item across instruments will be examined first. Figure 3 shows the three single-item correlations resulting from the analysis. (ITC items are indicated by a Roman numeral indicating the domain and the item number within that domain; RTOP items are represented as their item number.)

Figure 3. ITC and RTOP Single Item Correlations



ITC III.7 and RTOP 9 are similarly worded items within the Content domain in both instruments, providing evidence of the similar lineage of the items (Table 3). This pair can be labeled “Abstract Thinking.” A slight difference appears in these two indicators, with the ITC indicator using “included” and the RTOP using “encouraged” as the action to observe. Although these items have a strong correlation (0.833), they are slightly different indicators. Implicit in the RTOP indicator is some action on the part of the teacher that encourages students to use elements of abstraction. If abstraction was used with no indication that the teacher encouraged students to use elements of abstraction, the item can be rated high on the ITC. The RTOP indicator should be rated high only if encouragement was observed.

Table 3. ITC III.7 and RTOP 9 Alignment (Abstract Thinking)

| | |
|---|--|
| ITC III.7 (M/S Content) | Elements of mathematical/science abstraction (e.g., symbolic representations, theory building) were included when it was important to do so. |
| RTOP 9 (Content: Propositional Knowledge) | Elements of abstraction (i.e., symbolic representations, theory building) were encouraged when it was important to do so. |

The ITC II.6 and RTOP 1 items focus on the teacher’s pedagogical skills in designing and implementing instruction (Table 4). This grouping may be labeled “Student Conceptual

Understanding.” While they are not the only pedagogical skills assessed in the instruments, these two indicators focus specifically on what the student is bringing to the classroom and how that information can be used to enhance conceptual understanding. Without the parenthetical language in the ITC indicator, these two indicators do not seem to address the same behaviors.

Evaluators have found that observers have a tendency to focus on the parenthetical examples when rating an indicator with these types of concrete examples. Students’ prior conceptions/knowledge and preconceptions and misconceptions are mentioned in the two indicators, which may have provided the push needed for the rater to more closely rate these two indicators. Without the parenthetical descriptions, these two indicators may not have exhibited a close correlation.

Table 4. ITC II.6 and RTOP 1 Alignment (Student Conceptual Understanding)

| | |
|------------------------------------|--|
| ITC II.6 (Implementation) | The teacher’s questioning strategies were likely to enhance the development of student conceptual understanding/problem solving (e.g. emphasized higher order questions, appropriately used “wait time,” identified prior conceptions and misconceptions.) |
| RTOP 1 (Design and Implementation) | The instructional strategies and activities respected students’ prior knowledge and the preconceptions inherent therein. |

The only other matched pair is ITC IV.3 and RTOP 18 (Table 5), addressing “Significant Student Interactions.” Both of these items appear in the Culture categories of the two instruments and reflect students talking and working together within the context of the lesson.

Table 5. ITC IV.3 and RTOP 18 Alignment (Significant Student Interactions)

| | |
|---|--|
| ITC IV.3 (Classroom Culture) | Interactions reflected collegial working relationships among students (e.g., students worked together, talked with each other about the lesson). |
| RTOP 18 (Classroom Culture: Communicative Interactions) | There was a high proportion of student talk and a significant amount of it occurred between and among students. |

The ITC IV.3 and RTOP 18 match reveals differences in the two instruments pointed out in the authors’ previous article (Henry, Murray and Phillips, 2007). The ITC tends to word indicators so that they are more conceptual and inclusive of several behaviors that could be seen in a classroom.

ITC IV.3 qualifies the type of student interaction, such as “collegial working relationships.” Students should be working together collegially, on-task, and interacting with each other about the lesson. The RTOP 18 item focuses the observer on a significant amount of student talk (“high proportion”). The quality or focus of the talk does not necessarily have to be lesson-oriented or productive. The quality of the talk is not indicated. Any observer of classroom interactions would be expected to infer that the student talk should be productive, but this is not necessary for rating this indicator high. An observer may see group work occurring with a high level of student talk that was not focused on the lesson. Unless the rater went to the groups and listened to the talk, she or he would not know whether the talk was lesson-oriented, but could still rate the indicator

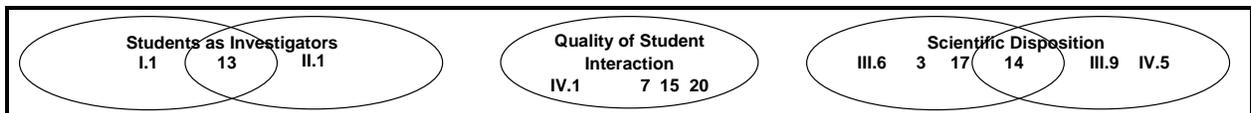
high. The fact that these two indicators correlated highly shows that this observer, who was trained on both instruments, interpreted the items similarly. This may not be the case for multiple observers or in situations where listening to student conversations is not possible.

These three pairs of indicators are the only unique matches across instruments. One could conclude that the instruments are not as similar as a cursory consideration may indicate. An examination of additional correlations across instruments provides a picture of other areas of congruence.

High Multiple Correlations

Items that correlated in multiple ways are show in Figure 2.

Figure 2. ITC and RTOP Multiple Item Correlations



The following analysis examines each set of items that correlates at .75 or above. From this analysis, new categories are identified. For example, in the first set, ITC I.1 and RTOP 13 correlated. RTOP 13 also correlated with ITC II.1. These indicators are considered as a set for analysis. The other sets are examined in the same way. Synthesis ratings will be examined in another section. (As previously explained, RTOP items 21-25 have been omitted from this analysis.)

The analysis of the items is organized around the ITC items.

Table 6. ITC 1.1, RTOP 13 and ITC II.1 Alignment (Students as Investigators)

| | |
|---|---|
| ITC I.1 (Design) | The design of the lesson incorporated tasks, roles, and interactions consistent with investigative mathematics/science. |
| RTOP 13 (Content: Procedural Knowledge) | Students were actively engaged in thought-provoking activity that often involved the critical assessment of procedures. |
| ITC II.1 (Implementation) | The instructional strategies were consistent with investigative mathematics/science. |

The set of items in Table 6 may be named “Students as Investigators.” This set incorporates indicators from the Design, Implementation and Procedural Knowledge categories from the two instruments. In a classroom where these ratings would be high, an observer would see students participating in investigative activities where they are actively involved in solving problems, either of their own or teacher design. This classroom is not one where labs are designed with one outcome in mind or where the procedures have been predetermined for the students. It is also not a classroom where teachers have a preconceived process and the class comes to consensus on the teacher’s design.

The next set of multiple correlations addresses classroom interactions (Table 8). These indicators focus on how students participate in the classroom. Do they listen to others and are they listened to? Are they discussing concepts in a rigorous way? Are they challenged to support and explain their thinking? This category could be labeled “Quality of Student Interactions.”

Table 8. ITC IV.1, RTOP 7, 15 and 20 Alignment (Quality of Student Interactions)

| | |
|---|---|
| ITC IV.1 (Classroom Culture) | Active participation of all was encouraged and valued. |
| RTOP 7 (Content: Propositional Knowledge) | The lesson promoted strongly coherent conceptual understanding. |
| RTOP 15 (Content: Procedural Knowledge) | Intellectual rigor, constructive criticism, and the challenging of ideas were valued. |
| RTOP 20 (Classroom Culture: Communicative Interactions) | There was a climate of respect for what others had to say. |

In the last grouping, ITC III.6 aligned with RTOP 3, 14, and 17. RTOP 14 also aligned with III.9, and IV.5 (Table 7).

Table 7. ITC III.6, RTOP 3, 14, 17, ITC III.9 and IV.5 Alignment (Scientific Disposition)

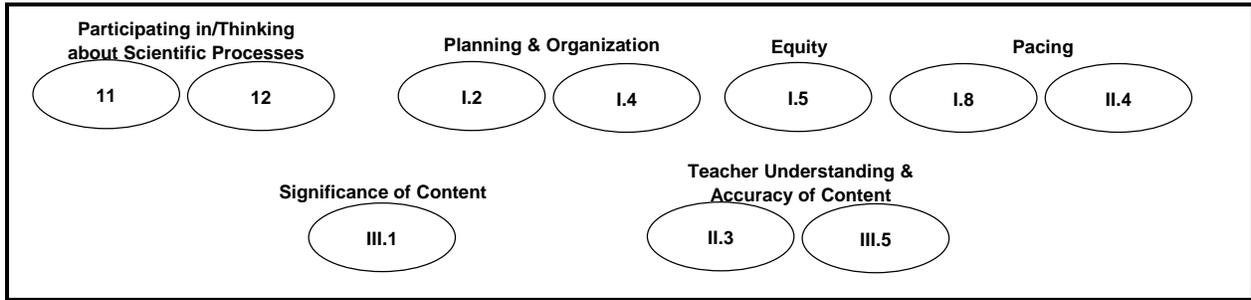
| | |
|---|---|
| ITC III.6 (M/S Content) | Mathematics/science was portrayed as a dynamic body of knowledge continually enriched by conjecture, investigation analysis, and/or proof/justification. |
| RTOP 3 (Design and Implementation) | In this lesson, student exploration preceded formal presentation. |
| RTOP 14 (Content: Procedural Knowledge) | Students were reflective about their learning. |
| RTOP 17 (Classroom Culture: Communicative Interactions) | The teacher’s questions triggered divergent modes of thinking. |
| ITC III.9 (M/S Content) | The degree of “sense-making” of mathematics/science content within this lesson was appropriate for the developmental levels/needs of the students and the purposes of the lesson. |
| ITC IV.5 (Classroom Culture) | The climate of the lesson encouraged students to generate ideas, questions, conjectures, and/or propositions. |

This category could be named “Scientific Disposition” because it incorporates various processes of understanding and scientific sense-making including reflective thinking, divergent thinking, conjectures, analysis, and justification of ideas. Elements of scientific processes are apparent in these six indicators.

Non-aligned Items

An analysis of the items on each instrument that do not align with any of the items on either instrument provides insight into what is unique about each instrument. Figure 3 shows the items that did not match any others.

Figure 3. ITC and RTOP Items with No Correlations



The RTOP contains only two items (11 and 12) that do not align in some statistically significant way with the items on the ITC COP or to other items on the RTOP (Table 9).

Table 9. RTOP Items with No Matches

| RTOP Items with No ITC Matches | |
|--------------------------------|--|
| RTOP 11 | Students used a variety of means (models, drawings, graphs, concrete materials, manipulatives, etc.) to represent phenomena. |
| RTOP 12 | Students made predictions, estimations and/or hypotheses and devised means for testing them. |

These two items address the student’s participation in thinking about and participating in the scientific processes (12) and in communicating her or his ideas to others (11). No ITC indicators specifically address student behavior or actions in this way. The RTOP generally tends to orient items toward student actions, as shown in these indicators.

There were eight ITC items that had no statistically significant matches. (See Table 10.)

Table 10. ITC Items with No RTOP Matches

| ITC Items With No RTOP Matches | |
|--------------------------------|---|
| I.2 | The design of the lesson reflected careful planning and organization. |
| I.4 | The resources available in this lesson contributed to accomplishing the purposes of the instruction. |
| I.5 | The instructional strategies and activities reflected attention to issues of access, equity, and diversity for students (e.g. cooperative learning, language appropriate strategies/materials). |
| I.8 | Adequate time and structures were provided for wrap-up |
| II.4 | The pace of the lesson was appropriate for the developmental levels/needs of the students and the purposes of the lesson. |
| III.1 | The mathematics/science content was significant and worthwhile. |
| III.3 | Teacher-provided content information was accurate. |
| III.5 | The teacher displayed an understanding of mathematics/science concepts (e.g. in his/her dialogue with students). |

The unmatched ITC items represent teacher planning and organization, including resources for classroom use (items I.2, I.4), issues of equity (item I.5), lesson pacing including time for wrap-up (items II.4 and I.8), issues of content including the teacher’s understanding and accuracy of presentation of content (III.3 and III.5) and the significance of the content being presented (III.1).

More attention is given in the ITC to the teacher’s planning process. It is assumed that this information could be ascertained in a pre-observation conference or inferred through the obvious lesson structure and implementation. Observers report that these are difficult items to rate with only a classroom observation because the care of planning and organization must be interpreted through how smoothly the lesson is implemented, which could be the result of other factors.

The ITC item related to resources (I.4) stands alone among the two instruments in assessing the materials and supplies that support the lesson. For projects focusing on use of manipulatives, specific equipment or technology, this would be an important item to consider during the observation.

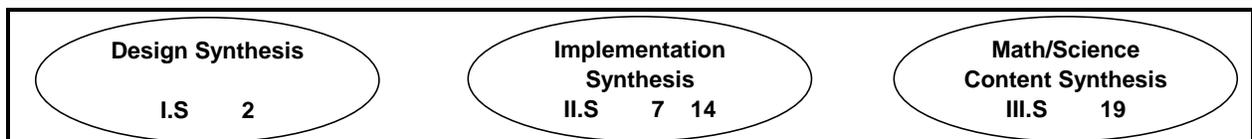
Likewise, equity is overtly stated in only item I.5 of the ITC. Equity is inferred in items in the RTOP, especially to Item 20 referring to a “climate of respect for what others had to say.” This RTOP item does not go nearly far enough to bring a teacher’s attention to the other issues involved in classroom equity and diversity.

Lesson implementation that focuses on pacing and the all-important time for wrap-up is found only in the ITC COP (I.8, II.4). The RTOP focuses lesson implementation on student behavioral factors such as students as learners in a community, rather than on teacher intentions and actions.

Synthesis Alignments

In the ITC, Synthesis ratings are given for each of the four domains, Design, Implementation, Mathematics/Science Content and Classroom Culture. According to the instructions, these ratings are not to be considered an average of the indicators within that domain, but a general rating of that domain, with some indicators keyed as important in the decision-making process. There was no comparable rating on the RTOP. One would expect the Synthesis ratings to closely correlate with the indicators within the domain, but that is not the case. No synthesis rating of any of the four aligned with any of the ITC indicators within its same domain. The Synthesis ratings did align with one or more of the RTOP ratings. Figure 4 shows correlations found among domain synthesis items.

Figure 4. ITC and RTOP Domain Synthesis Rating Correlations



Three Synthesis Ratings correlated with four RTOP indicators. There was no correlating item from either instrument for the Classroom Culture Synthesis Rating (IV.S).

Table 11. ITC Synthesis Correlations with RTOP Items

| ITC Synthesis Correlations with RTOP Items | | |
|--|---|---|
| ITC I.S (Design Synthesis) | RTOP 2 (Design and Implementation) | The lesson was designed to engage students as members of a learning community |
| ITC II.S (Implementation Synthesis) | RTOP 7 (Content: Propositional Knowledge) | The lesson promoted strongly coherent conceptual understanding |
| | RTOP 14 (Content: Procedural Knowledge) | Students were reflective about their learning |
| ITC III.S (Mathematics/ Science Content) | RTOP 19 (Classroom Culture: Communicative Interactions) | Student questions and comments often determine the focus and direction of classroom discourse |

The ITC Design Synthesis (I.S) rating correlated with RTOP 2 (Design and Implementation) indicator, “The lesson was designed to engage students as members of a learning community.” RTOP 2 is one of the more broadly stated RTOP indicators and captures one of the characteristics of reform-based classroom designs.

The ITC Implementation Synthesis (II.S) rating correlates with RTOP 7 (Content: Propositional Knowledge) “The lesson promoted strongly coherent conceptual understanding” and RTOP 14 (Content: Procedural Knowledge) “Students were reflective about their learning.” These two domain items reflect a classroom where students are reflective in their study of the concepts of mathematics or science, a characteristic of reform-based classes, where concepts are the central teaching goal rather than disparate items of knowledge and students are involved in meaningful thinking about the concepts.

The ITC Mathematics/Science Content Synthesis (III.S) rating correlates only with RTOP 19 (Classroom Culture: Communicative Interactions) “Student questions and comments often determine the focus and direction of classroom discourse.” Students’ directing their learning represents an important reform-based strategy. Student interest may lead the specific content within big ideas and direct the focus of the discussion in these classrooms.

The authors were curious about so few correlations being seen between the specific items and the Synthesis ratings. One would expect that if the indicators are reflective of best practice in mathematics and science and the Synthesis rating is a reflection of what is occurring in the classroom, the individual indicators would more closely align with the Synthesis ratings. One could question the value of the Synthesis rating if it does not represent the category. It ultimately may not be a helpful rating because of the subjective or otherwise hidden considerations being applied in arriving at it.

Authors have found during interrater reliability training that the synthesis ratings are often places where the raters tend to apply their individual preferences or biases. Statements like, “I knew that was what the teacher intended to do here so I just increased the score a little,” or “I figured that was what the teacher was planning to do tomorrow.” It takes concerted effort to avoid these types of tendencies in raters during training, and it may still have existed in this rater. More research on the efficacy of the Synthesis Ratings is called for.

Discussion

The gold standard of classroom observation instruments has yet to be developed. Without research to show they are rating what they claim, the ITC and the RTOP are still the most widely used instruments when classroom observations are undertaken for large-scale mathematics and science projects. This quantitative analysis of the ITC and the RTOP indicates even more than the previous qualitative analysis that these instruments, though similar in some aspects, are very different, both in design and in effect. It is recognized that this analysis is preliminary and based on a relatively small number of cases. With an increase in the number of subjects and raters who have undergone interrater reliability training, more commonalities may emerge.

The authors acknowledge that observer accuracy and consistency, being inherent concerns when considering observational data, may be questioned in interpreting these results. The choice of a single, experienced educator as an observer was made to alleviate interrater reliability considerations and enable a more determined focus on underlying instrumental concerns. The authors appreciate that, especially in the differentiated scoring scale for ITC synthesis ratings, observer summative judgment inevitably is involved. An observer who is generally consistent would generally score synthesis items in a particular way, which in the case of these results could affect the apparent matches or lack of matches that appear. However, the authors, in evaluating the synthesis scores and performing post-observation interviews with the observer have detected no patterns that would yield scoring incompatible with the stated directions of the instrument protocols.

Despite the same foundational philosophy, seminal documents and similar indicator categories, these instruments do not measure the same things. They have been developed by designers with their own underlying definition of commonly used pedagogical terms, such as inquiry and sense-making. Examples in each of the items illustrate those understandings and may cause a narrowing of interpretation of the item.

Meanwhile, the authors recommend that teacher professional development and other educational programs, when considering use of these instruments for classroom observations, examine the items closely to determine which instrument more closely represents what their project ideally would look like in the field. Raters should spend time to establish a common understanding of the items prior to using either of these instruments.

Likewise, raters bring their own interpretations and understandings built through individual experiences and biases. Researchers caution sending inexperienced education students into the classroom using either of these instruments. If it is necessary to use graduate students, time should be spent assuring that they understand the intent of the item and how it aligns with the

project being rated. Ongoing interrater reliability training will be necessary for all raters as observations occur across years.

Future Research

The need for a reliable observation instrument representing current understanding of reform teaching is critical. Some researchers and evaluators have stopped including classroom observations as part of their data collection because of the unreliability of data and the amount of time and resources needed to complete the observations. The authors feel that to eliminate classroom observations is closing the black box of implementation and any outcomes may not be reliably attributable to the project's intervention. The authors will continue to conduct additional field work to increase the sample sizes so additional correlations will emerge within and across these instruments, if present. We would encourage other researchers to join in this analysis.

This research was supported in part by the National Science Foundation EHR # 06-34423.

Martha A. Henry, Ed.D., is President and Lead Evaluator of M.A. Henry Consulting, LLC in St. Louis, Missouri. mahenry@mahenryconsulting.com

Keith S. Murray, CEP, is CEO and Lead Evaluator/Researcher of M.A. Henry Consulting, LLC. keithsmurray@mahenryconsulting.com

Mark Hoglebe, Ph.D., is Institutional Researcher at Washington University in St. Louis, mhoglebe@wustl.edu

Marcia Daab, Ed.D., is a former science supervisor and presently Educational Consultant. marciajd@sbcglobal.net

References

Henry, M. A., Murray, K. S., & Phillips, K. (2007). *Meeting the challenge of STEM classroom observation in evaluating teacher development projects: A comparison of two widely used instruments*. Retrieved from <http://hub.mspnet.org/index.cfm/14975>

Horizon Research, Inc. (2002). *Inside the classroom interview and analytic protocol*. Retrieved from <http://www.horizon-research.com/instruments/clas/cop.php>

Sawada, D., Piburn, M., Falconer, K., Turley, J., Benford, R., & Bloom, I. (2000). *Reformed Teaching Observation Protocol: Technical Report No. IN00-1*. Tempe, AZ: Arizona State University.