

Appendix

Evaluating students' abilities to construct mathematical models from data using latent class analysis

Alexandra Brandriet, Charlie A. Rupp, Katherine Lazenby, Nicole M. Becker

Appendix 1: The Online Rate Law Assessment:

Activity Instructions (p. 1 of 3):

In the following activity, you are going to be asked questions related to some of the chemical kinetics topics that you learned about in your Principles II course. In this activity, you will receive a series of 3 different chemical reaction scenarios, and you will be asked a series of questions related to the rate law of these reactions. There are 17 questions on this survey. You can use paper, a pencil, and a calculator to solve these problems.

Please note that you may not be able to return to a previous screen once you have completed a step.

Click the next button to continue.

Activity Instructions (p. 2 of 3):

Please answer the following questions on your own. **DO NOT** use additional resources, such as your textbook or the internet. If you do not know how to answer a question, try to provide as much information as you can, to the best of your ability.

Your score will reflect the amount of *effort* that you put into taking this survey, rather than the correctness of your answers. So please respond with what you honestly think.

Click the next button to continue.

Activity Instructions (p. 3 of 3):

During this activity, you will be asked to explain your problem solving approach. It may be helpful to use the following shorthand notation:

Mathematical Expression	Shorthand Notation
5^2	5^2
2×3	$2*3$
$2 \div 3$	$2/3$

Click the next button to continue.

1. What is your understanding of the concept of a rate law?

2. What does the term 'order' of a reaction mean to you?

3. The generic form of a rate law is often expressed as: $\text{Rate} = k[\text{A}]^m[\text{B}]^n$

a. What do you think the m and n represent?

b. What do you think the k represents?

4. Why do you think rate laws are important to chemists? (*hint: what information can chemists get from a rate law and why might this information be useful*)

Directions: Use Chemical Scenario 1 to help you answer questions 5 through 9.

Chemical Scenario 1:

The reaction of A and B at 1280oC is



The following data was collected for this reaction:

Experiment	[A] (M)	[B] (M)	Initial Rate (M/s)
1	0.100	0.100	5.00
2	0.100	0.200	10.0
3	0.300	0.100	44.0
4	0.300	0.200	88.0

5. Which of the following rate laws makes sense based on the information provided in **Chemical Scenario 1**? If none of them makes sense, type your rate law in the space below.

- Rate = $k[\text{A}]^2[\text{B}]^2$
- Rate = $k[\text{A}]^3[\text{B}]$
- Rate = $[\text{C}]^2 / [\text{A}]^4[\text{B}]^3$
- Rate = $k[\text{A}]^2[\text{B}]$
- None of the above are correct. The rate law for this reaction is: *(type your rate law in the box below)*

6. A student in your class is confused about how to determine the rate law for this reaction. In the space below, explain to your classmate how you determined that the exponent for reactant A [*would be 2, 3, or 4*].

Be **specific** about your **procedure** and your **reasoning**.

7. Explain to your classmate how you determined that the exponent for reactant B [*would be 1, 2, or 3*].

Be **specific** about your **procedure** and your **reasoning**.

(the prompt for questions 6 and 7 changed depending on the multiple-choice response the student chose in question 5)

Directions: Use Chemical Scenario 1 to help you answer questions 5 through 9.

Chemical Scenario 1:

The reaction of A and B at 1280oC is



The following data was collected for this reaction:

Experiment	[A] (M)	[B] (M)	Initial Rate (M/s)
1	0.100	0.100	5.00
2	0.100	0.200	10.0
3	0.300	0.100	44.0
4	0.300	0.200	88.0

8. Megan, a student in your class, determined that the rate law is: $\text{Rate} = k[\text{A}][\text{B}]^2$.

Do you think that Megan's rate law is correct?

- Yes
- No

9. Using the information provided in Chemical Scenario 1, ***explain*** to Megan why you think that her rate law is [*correct or incorrect*].

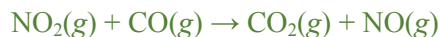
Be ***specific*** in your response.

(the prompt for question 9 changed depending on the multiple-choice response the student chose in question 8)

Directions: Use Chemical Scenario 2 to help you answer questions 10 through 14.

Chemical Scenario 2:

The reaction of NO₂ and CO at 400K is



The following data was collected for this reaction:

Experiment	[NO ₂] (M)	[CO] (M)	Initial Rate (M/s)
1	0.050	0.025	7.768x10 ⁻¹¹
2	0.050	0.075	7.768x10 ⁻¹¹
3	0.100	0.025	3.107x10 ⁻¹⁰
4	0.100	0.075	3.107x10 ⁻¹⁰

10. Which of the following rate laws makes sense based on the information provided in **Chemical Scenario 2**? If none of them makes sense, type your rate law in the space below.

- Rate = k[NO₂]²[CO]³
- Rate = k[NO₂]²
- Rate = [CO₂][NO] / [NO₂][CO]
- Rate = k[NO₂]
- None of the above are correct. The rate law for this reaction is: *(type your rate law in the box below)*

11. A student in your class is confused about how to determine the rate law for this reaction. In the space below, explain to your classmate how you determined that the exponent for NO₂ [*would be 1 or 2*].

Be **specific** about your **procedure** and your **reasoning**.

12. Explain to your classmate how you determined that the exponent for CO [*would be 0, 1, or 3*].

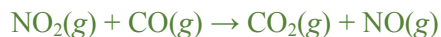
Be **specific** about your **procedure** and your **reasoning**.

(the prompt for questions 11 and 12 changed depending on the multiple-choice response the student chose in question 10)

Directions: Use Chemical Scenario 2 to help you answer questions 10 through 14.

Chemical Scenario 2:

The reaction of NO₂ and CO at 400K is



The following data was collected for this reaction:

Experiment	[NO ₂] (M)	[CO] (M)	Initial Rate (M/s)
1	0.050	0.025	7.768x10 ⁻¹¹
2	0.050	0.075	7.768x10 ⁻¹¹
3	0.100	0.025	3.107x10 ⁻¹⁰
4	0.100	0.075	3.107x10 ⁻¹⁰

13. Fred, another student in your class, determined that the rate law is: Rate = k[CO].

Do you think that Fred's rate law is correct?

- Yes
- No

14. Using the information provided in Chemical Scenario 2, ***explain*** to Fred why you think that his rate law is [*correct or incorrect*].

Be ***specific*** in your response.

(the prompt for question 14 changed depending on the multiple-choice response the student chose in question 13)

Directions: Use Chemical Scenario 3 to help you answer questions [15 through 17 or 15 & 16].

Chemical Scenario 3:

The reaction of O₂ and NO at 298K is



The following data was collected for this reaction:

Experiment	[O ₂] (M)	[NO] (M)	Initial Rate (M/s)
1	0.015	0.045	0.220
2	0.030	0.045	0.440
3	0.025	0.025	0.113
4	0.050	0.050	0.905

15. Which of the following rate laws makes sense based on the information provided in Chemical Scenario 3? If none of them makes sense, type your rate law in the space below.

- Rate = $k[\text{O}_2][\text{NO}]^2$
- Rate = $k[\text{O}_2]$
- Rate = $[\text{NO}_2]^2 / [\text{O}_2][\text{NO}]^2$
- It is impossible to determine the rate law with the information provided.
- None of the above are correct. The rate law for this reaction is: (type your rate law in the box below)

16. A student in your class is confused about how to determine the rate law for this reaction. In the space below, explain to your classmate how you determined that the exponent for O₂ [would be 1].

Be **specific** about your **procedure** and your **reasoning**.

17. Explain to your classmate how you determined that the exponent for NO [would be 0 or 2].

Be **specific** about your **procedure** and your **reasoning**.

(the directions and questions 16 & 17 prompt changed depending on the multiple-choice response the student chose in question 15)

Appendix 2: Summary of the changes made to the coding scheme from Becker et al. (2017)

Table A1. Comparison of the levels coding scheme across the qualitative and quantitative studies

Level	In-depth, qualitative study (Becker <i>et al.</i> , 2017)	Large scale, quantitative study (Brandriet <i>et al.</i> , submitted)
5	Students examine the experimental data by selecting two trials such that one variable is held constant; they correctly describe the pattern in the data verbally (e.g. concentration doubles, rate quadruples) and appropriately use exponents in the rate law expression to model observed changes. Responses may include some explicit reflection as to how the selected exponent models change in the reaction rate.	Students are able to interpret the changes in concentration and rate, while holding a variable constant (or accounting for another variable), and appropriately reason about how concentration exponentially influences the rate, depending on the order.
4	Students select two trials such that concentration of one variable is constant; they may correctly describe some aspect of the pattern in the data verbally (e.g. rate doubles, concentration doubles) However, they either 1) use incorrect intuitive reasoning and therefore select an incorrect exponent for the rate law or 2) make an error in translating their qualitative description into an exponent for the rate law. Students' reasoning at this level often does not include elaborating how the trend in the data relates to the selected exponent.	Students are able to interpret the changes in concentration and rate, while holding a variable constant (or accounting for another variable). However, students have difficulty reasoning about how the exponent relates to the interpretation of the data; this includes (1) no explicit reasoning, or (2) difficulties with exponential reasoning.
3	Students select two trials such that one variable is held constant; they discuss the impact of changing reactant concentration on the change in initial rate and attempt to interpret the trend in the data qualitatively using language such as "doubles" "triples", "increases by more than double". Students' interpretations of either rate or concentration may be incorrect or too vague to serve as an appropriate foundation for determining a reaction order. If the student selects an exponent, they typically do so without clear relation to the pattern in the data. That is, there is no explicit reasoning.	Students are able to recognize that they need to interpret how the concentration <u>and</u> rate vary, while holding another variable constant (or accounting for another variable), but their responses includes inappropriate or low-level interpretations of the changes in concentration or rate.
2	Students examine the experimental data and use an algorithmic approach to the find orders of the reactants. They describe what they did, but not why they did it or how the outcome of the calculation relates to the general pattern in the data. They do not discuss how the selected exponents models the trend in the data.	Students use the experimental data as evidence for determining the exponents; however, their arguments and/or use of the data is low level. Students in this level are (1) using procedures without interpreting the change in concentration and rate in terms of the exponent (<i>e.g.</i> $2^2 = 4$), (2) only focusing on concentration or rate, but not both, or (3) providing an argument that only the experimental trials used, or (4) failing to hold one variable constant when they interpret the data (or accounting for the change in O_2 in scenario 3).
1	Students use the stoichiometric coefficients as the exponents in the rate law; they do not use the rate data and do not discuss holding one variable constant to see the impact of concentration on rate.	Students reason using external features of the problem without attempting to use the data. Most of the students in this category used the coefficients in the chemical equation to determine the order or conflated a rate law with an equilibrium constant. When students used the data table, they used memorized procedures that did not require them to compare data values, such as focusing on the units of the initial rate.

Becker *et al.* (2017) conducted 15 in-depth interviews with students using a method of initial rate task like Task 1. The levels coding scheme emerged through an analysis of the interview data. We then used the coding scheme from Becker *et al.* (2017) as a closed-coding scheme for our current study; however, our current study included a larger sample of students that completed three method of initial rate tasks, rather than one. Therefore, we expanded the coding scheme to accommodate the variations in the students' responses. Table A1 provides a side-by-side comparison of the coding schemes used in both studies, and the following is a description of the differences between the coding schemes (the following student quotes are shown verbatim, which includes students' spelling and grammatical errors).

Level 1

Level 1 responses frequently derived the reaction order using the coefficients in the chemical equation. In Becker *et al.* (2017), the students exclusively used this approach; however, in the current study, a small number of students focused on other explicit features of the problems. As an example, for A in Task 1, one student described that they used the units of the rate to determine the order:

"First you determine which of all the molecules/compounds involved are the reactants. These are the concentrations that will be in the rate law. second you raise each reactant to it's order number, NOT its stoichiometric coefficient, although the two numbers can be the same. The overall order number is zero, seeing as the units are in M/s." (F226, Task 1, Rate law: "rate=k", Class 1 Level 1)

Level 2

Level 2 responses used the datasets to derive their reaction orders, but did so in a low-level manner. Of the five levels, Level 2 required the most expansion, because the students used a wider range of low-level approaches in the quantitative study. In Becker *et al.* (2017), two out of 15 students used the divide two trials approach to solve for the exponents in the rate law. While this approach can be used to solve for the answer, students who used this approach did not include meaningful interpretations of the exponent relative to the concentration and rate variables. Similarly, many students in our current study also used this approach, as shown in the following quote for A in Task 1:

"I determined the exponent for reactant A would be two by looking at the table of values and choosing two separate reactions. These reactions would have to have 2 different values for A while keeping the value for B constant. I then wrote a skeleton rate law for each reaction. I then set up a ratio of the two rate laws, cancelling constants and then solving for m." (S28, Task 1, Selected rate law: "Rate = k[A]²[B]", Class 2 Level 2)

In contrast to Becker *et al.* (2017), we also found some additional low-level data use, such as when students only focused on the concentration or rate variables. Further, several students did not select two trials for which one concentration variable was held constant or accounted for, while identifying how the other concentration variable influenced the rate. As an example, for A in Task 1, the following student chose two trials for which the [A] was constant, compared this to the change in the rate, and determined that since the [A] was not changing, A must be zeroth order:

"Looking at Experiments 1 and 2, for example, the exponent for [A] should be 0. The rate appears to be independent of the concentration." (F175, Task 1, Rate law: "Rate = k[A]⁰[B]²", Class 3 Level 2)

Finally, in the current study, we found a small subset of students that provided low-level arguments for how they analyzed the data. While these students may not necessarily have provided incorrect

reaction orders or have used incorrect procedures, their discussion of their analysis strategies were low-level:

“Compare Experiments 1 and 3 because the concentration of B is the same, allowing us to compare concentrations of A and the initial rates.” (S114, Task 1, Selected rate law: “Rate = $k[A]^3[B]$ ”, Class 2 Level 2)

Level 3

Level 3 was used very similarly in Becker *et al.* (2017) and the current study. In both studies, Level 3 responses inferred that the student was able to select two trials for which the additional concentration variable was held constant, but often struggled to analyze changes in the concentration and rate variables. Students often used terms like “doubling,” or “tripling” to describe the changes in the concentration and rate variables, but their responses implied that they struggled to identify those changes.

Level 4

In Becker *et al.* (2017), Level 4 student responses described the patterns in the data (*i.e.* the rate doubles as the [B] doubles), but used intuitive reasoning relative to the exponent, rather than mathematically-based reasoning. As an example, students tended to describe that the exponent with respect to B in Task 1 was two because of the “doubling” described in their analysis of the data.

While these types of responses would also be considered Level 4 in our current study, we expanded this category to also include other types of incorrect reasoning, such as the multiplication or division strategies shown in the body of the text (*i.e.* $2/2 = 1$).

Level 4+ was used in our current study to characterize students’ responses that neglected reasoning, but provided the correct reaction order. The Level 4+ category was unique to this study and did not emerge in Becker *et al.* (2017). This was likely because Becker *et al.* (2017) used semi-structured interviews to collect students’ responses. Therefore, the interviewer was able to further probe the reasoning that student used, which was not possible in our current study.

Level 5

Level 5 responses were very similar across Becker *et al.* (2017) and the current study. In both studies, the students’ responses indicated that the students could interpret the patterns in the data and appropriately reason about how those patterns were used to derive the correct reaction order.

References:

Becker N. M., Rupp C. A., Brandriet A., (2017), Engaging students in analyzing and interpreting data to construct mathematical models: an analysis of students’ reasoning in a method of initial rate task, *Chemistry Education Research and Practice*, **18**(4), 798-810.

Appendix 3: LCA model selection

We ran six LCA models that fit 2-7 latent classes to the data; the resulting fit indices for each model are shown in [Table A2](#). When choosing the best model, experts recommend considering a combination of statistical output, parsimony, and interpretability (Collins and Lanza, 2010, p. 82). We eliminated the two-class solution, because there was evidence to reject the null hypothesis (*i.e.* the model does not fit the data when $p < 0.05$).

We also chose to eliminate the six and seven class solutions, because there was not enough evidence to suggest that these models were well identified. The algorithm that estimates the LCA parameters requires an initial starting value to start the algorithm, but different starting values may produce different estimates for the same model (Collins and Lanza, 2010). As a result, we ran each model 1000 times using random starting values, and the final column in [Table A2](#) describes the percent of times the models converged to the highest log-likelihood solution. Dziak and Lanza (2015) suggest that a researcher may choose to consider a model well-identified if at least 25% of the models converge to the highest log-likelihood solution. Additionally, we explored some of the other unique, log-likelihood solutions and found a high degree of similarity between these solutions and the five-class model presented in [Table A2](#).

Table A2. Summary of model identification and fit information used to select LCA model

Classes	Log Likelihood	df	G ²	p-value	AIC	BIC	Percent of best fitted model
2	-2154.84	99	241.2	<0.001	291.2	407.3	100%
3	-2068.85	86	69.2	0.907	145.2	321.7	100%
4	-2054.21	73	40.0	0.999	142.0	378.8	100%
5	-2047.63	60	26.8	>0.999	154.8	452.0	40%
6	not well identified ^c						
7	not well identified ^c						

a. $n = 768$

b. Convergence criterion set at <0.000001000; all models converged to a solution

c. Poorly identified model with best model with less than 25% of seeds (Dziak and Lanza, 2015), based on 1000 random starting values

[Figure A1](#) shows the latent class prevalence ([Figure A1a](#)) and item response probability estimates for the three-, four-, and five-class solutions ([Figures A1b-d](#)). The gradual increase in the probability of higher level responses across each class in each solution suggests a potential ordering of the latent classes in each model. LCA does not make assumptions about the ordered nature of the observed variables (*i.e.* Levels 1-5) or the classes that emerge from the model (*i.e.* Classes 1-5). The numbers assigned to the latent classes are arbitrary, so the authors assigned values to the results in [Figure A1](#) for the sake of clarity.

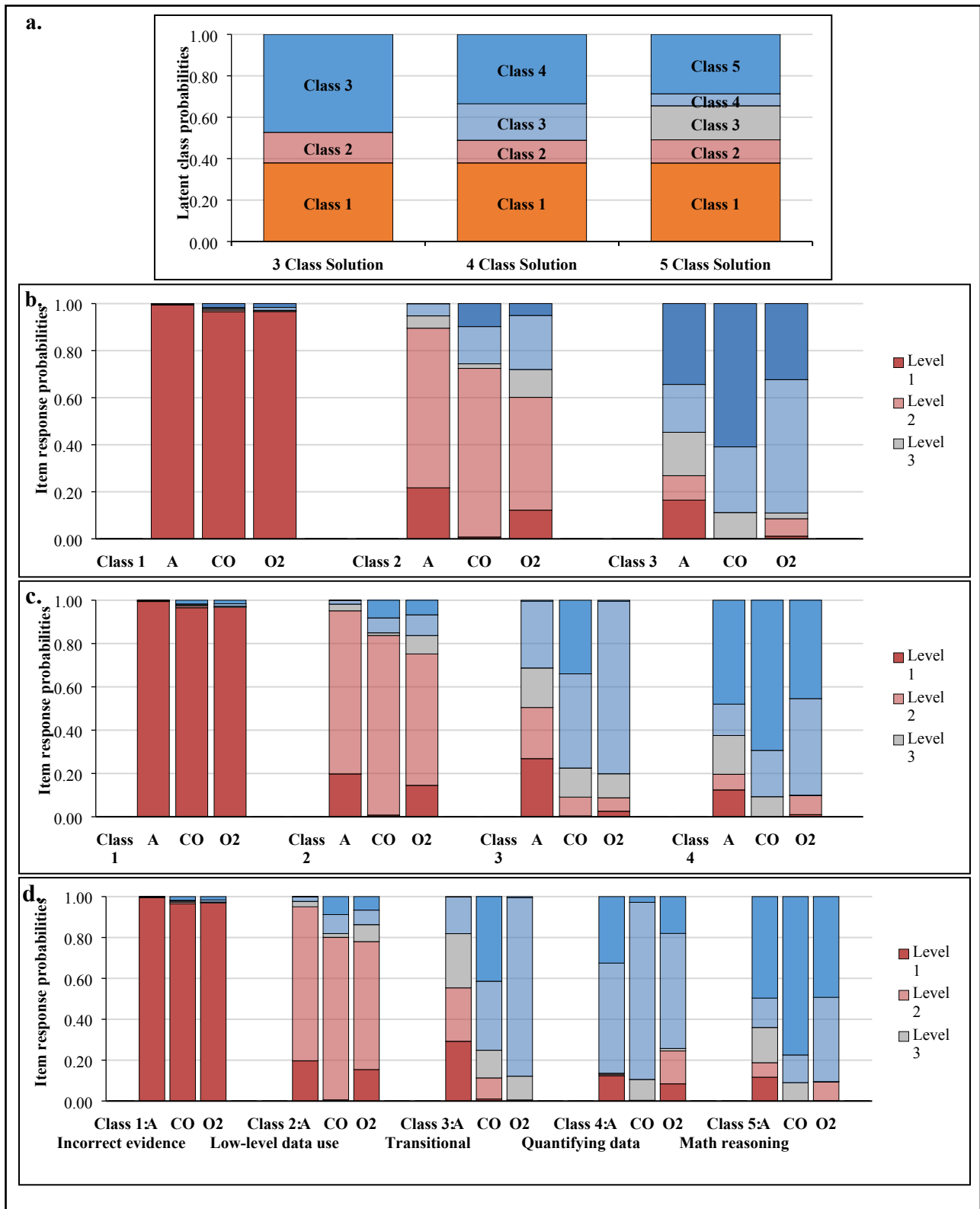


Fig. A1 Comparison of the latent class prevalence values (a) and the item response probabilities for the three- (b), four- (c), and five- (d) class model solutions.

While the low BIC value (Table A2) suggested the three-class model (Figure A1b), we believed that the response patterns across the three classes were still too heterogeneous to produce an interpretable solution. Class 3 in the three-class model held a wide range of Levels 3-5 responses, and as a result, was not sensitive enough to characterize these responses.

When comparing the four and five class solutions in Figures A1c and A1d, we noticed similarity in Classes 1, 2, and 3 across the four and five class solutions. Analogously, we noticed that Class 4 in the four-class solution and Class 5 in the five-class solution were quite similar. The distinguishing feature was the addition of Class 4 in the five-class model; this class had a high probability of Level 4 responses across all three questions, which suggests that the students did not provide correct reasoning that linked their synthesis of the data to their chosen order. The students classified in Class 4 in the five-class solution (based on posterior probability estimates) were primarily classified in Classes 3 and 4 (the highest two classes) in the four-class solution; this is shown in Figure A2.

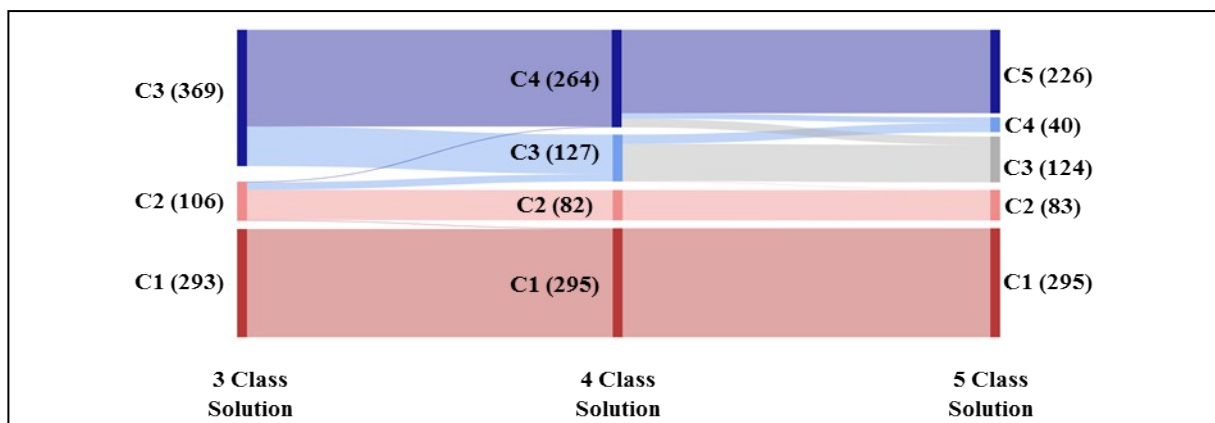


Fig. A2 Comparison of the students' class assignments for the three, four, and five latent class model solutions. Students were assigned to classes using the posterior probabilities estimates.

We decided to move forward with the five-class solution, because it included Class 4. Previous research has described that the reasoning component can be particularly difficult for students and that appropriate instruction can facilitate students' abilities to engage in appropriate reasoning (McNeill *et al.*, 2006). Therefore, we believe that this solution, which includes Class 4, provides additional insights into students' missing or incorrect reasoning (*i.e.* Level 4 responses).

We interpreted and labeled each latent class in the five-class model based on the item response probabilities associated with that class (*i.e.* Figure A1d). For each of the classes, except Class 3, there was one level that was clearly dominant across all three questions. Therefore, we labeled Classes 1, 2, 4, and 5 analogously with our coding scheme. Class 3 had more variation across the three questions, where students seemed to have the most difficulty with task A, which was likely because the [A] had a second order relationship with the rate. We labeled this class "Transitional" because the response patterns suggested that the students used both higher and lower-level arguments, potentially depending on the difficulty of the task.

References

- Collins L. M. and Lanza S. T., (2010), *Latent class and latent transition analysis: With applications in the social, behavioral, and health sciences*, Hoboken, N. J.: John Wiley and Sons, Inc.
- Dziak J. J. and Lanza S. T., (2015), *SAS graphics macros for latent class analysis users' guide (Version 2)*, University Park: The Methodology Center, Penn State, retrieved from <http://methodology.psu.edu>.
- McNeill K. L., Lizotte D. J., Krajcik J. and Marx R. W., (2006), Supporting students' construction of scientific explanations by fading scaffolds in instructional materials, *J. Learn. Sci.*, **15**(2), 153–191.