

Using large language models to summarize student feedback

Nicholas T. Young

Christopher Overton¹, Ania Majewska², Hina Shaikh^{1,3}, Nandana Weliweriya¹

¹ Department of Physics and Astronomy, University of Georgia

² Department of Physiology and Pharmacology, College of Veterinary Medicine,
University of Georgia

³ Institute for Astronomy and Astrophysics, Eberhard Karls University of
Tübingen



**Georgia Physics and Astronomy
Education Research**

Franklin College of Arts and Sciences

UNIVERSITY OF GEORGIA



Student feedback is useful but can be time-consuming to understand



Student feedback is useful but can be time-consuming to understand



The image displays six screenshots of student feedback forms, arranged in a 2x3 grid. Each form is a table with multiple columns and rows, containing handwritten student comments and numerical ratings. The forms are from the University of Georgia, as indicated by the logo in the top right corner of the slide. The comments are written in black ink on a white background. The ratings are numerical values, likely representing a score or grade. The forms are titled 'Student Feedback' and 'Feedback Form'.

Chatbots can summarize text and could be a useful tool to quickly understanding student feedback



What can I help with?

Ask anything



Search



Reason



Pu et al 2023; Parker et al 2024

What's the problem then?

What's the problem then?

- Unstructured text



What's the problem then?

- Unstructured text
- Sample size




What's the problem then?

- Unstructured text
- Sample size
- Hallucinations

Xu et al 2025



 Nicholas Young



Determine how well LLMs
can extract key insights from
student feedback

Data

- Student responses to two-questions on end of course feedback survey
 1. *What do you feel are your instructor's strengths and weaknesses?*
 2. *What do you feel are the strong and weak aspects of the course?*
- 9 courses taught by 3 unique instructors

Methodology

- 5 instructors read each set of student feedback and create a summary
- Same feedback files shared with 4 AI tools
 - LLMs have inherent randomness so did 5 trials with each model



Methodology

- Prompt: “For responses to open-ended questions, the goal is to focus on the useful information and identify trends or themes that appear. Note the frequency of themes, areas of agreement and disagreement among students, and suggestions students have for changes you might make. Please ignore the comments that are nonspecific. For the remaining comments, please sort them into three categories: positive, actionable suggestions, and nonactionable suggestions before identifying trends or themes”

<https://ctl.uga.edu/teaching-resources/feedback-and-evaluation-of-teaching/interpreting-responding-to-student-evaluations-of-teaching/>

Analysis

- Majority voting (3 of 5) to be included in final summary for each tool and 3 of 5 instructors

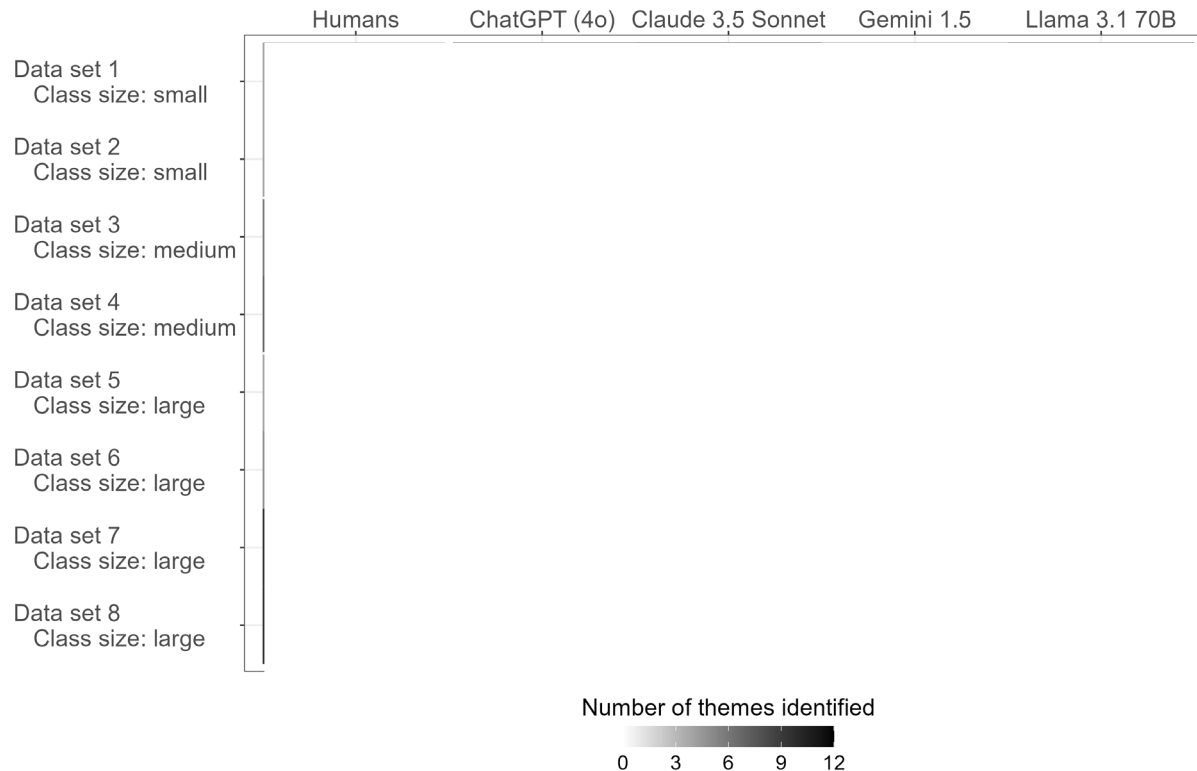
Analysis

- Majority voting (3 of 5) to be included in final summary for each tool and 3 of 5 instructors
- Take instructor summary as “true” summary and compare LLM summaries with human summaries

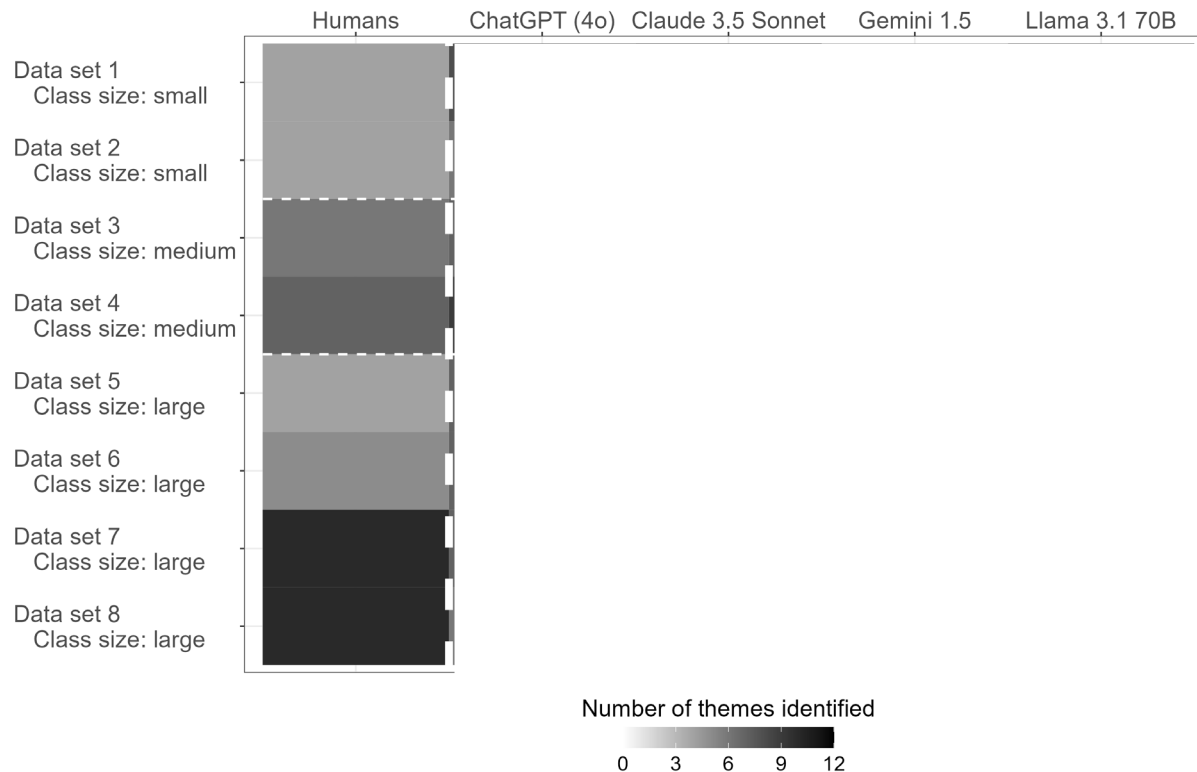
Analysis

- Majority voting (3 of 5) to be included in final summary for each tool and 3 of 5 instructors
- Take instructor summary as “true” summary and compare LLM summaries with human summaries
- Results preliminary, about halfway through data

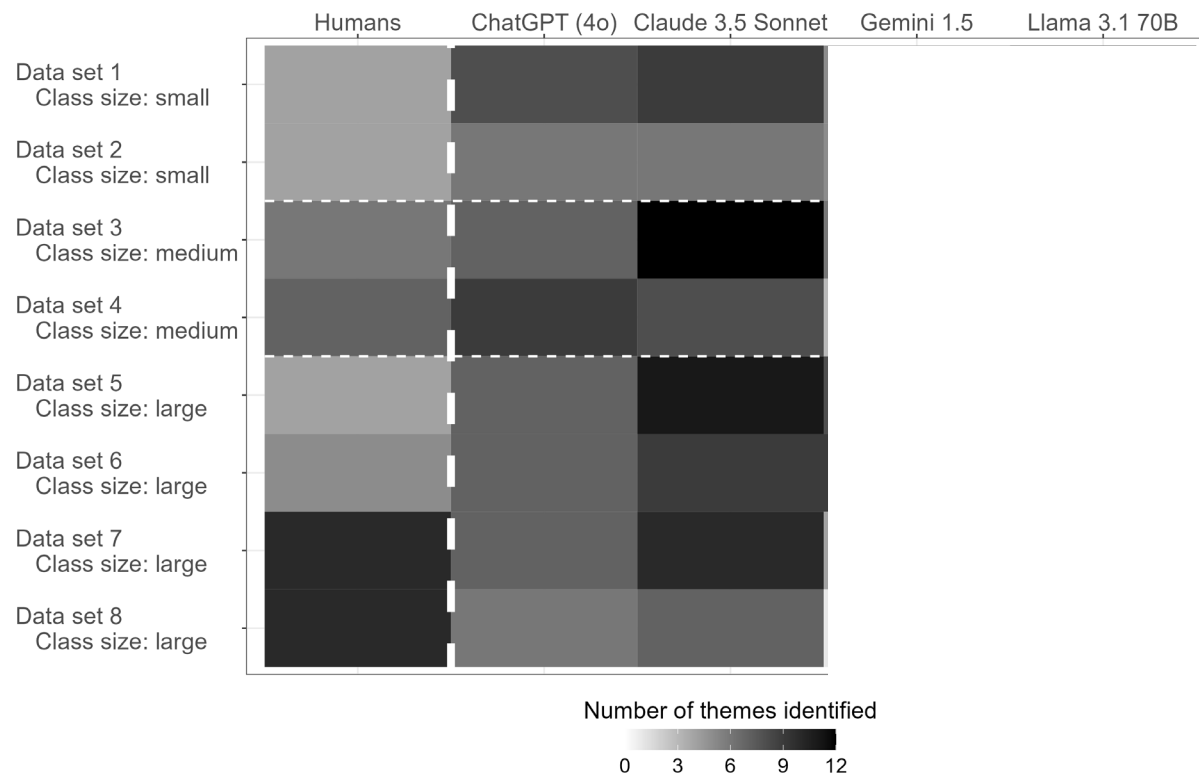
Results: themes identified by group



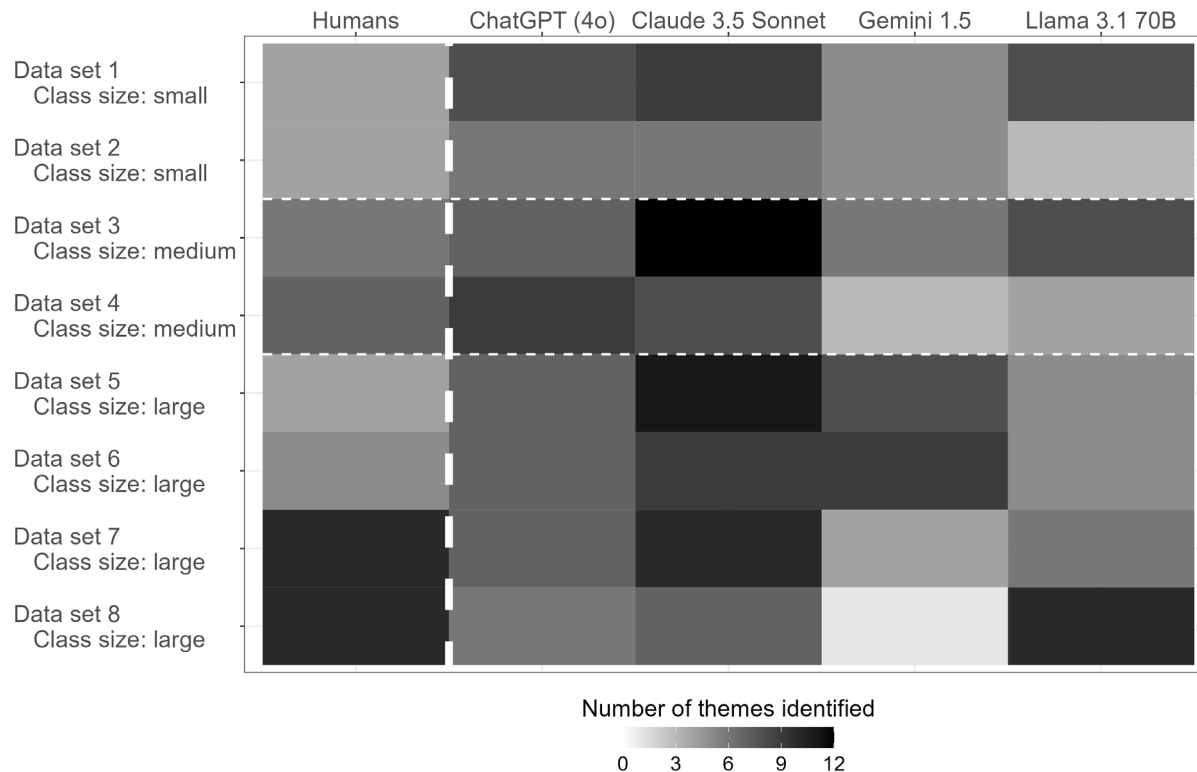
Results: themes identified by group



Results: ChatGPT and Claude identify more themes than instructors do



Results: Gemini and Llama identify fewer themes than instructors do

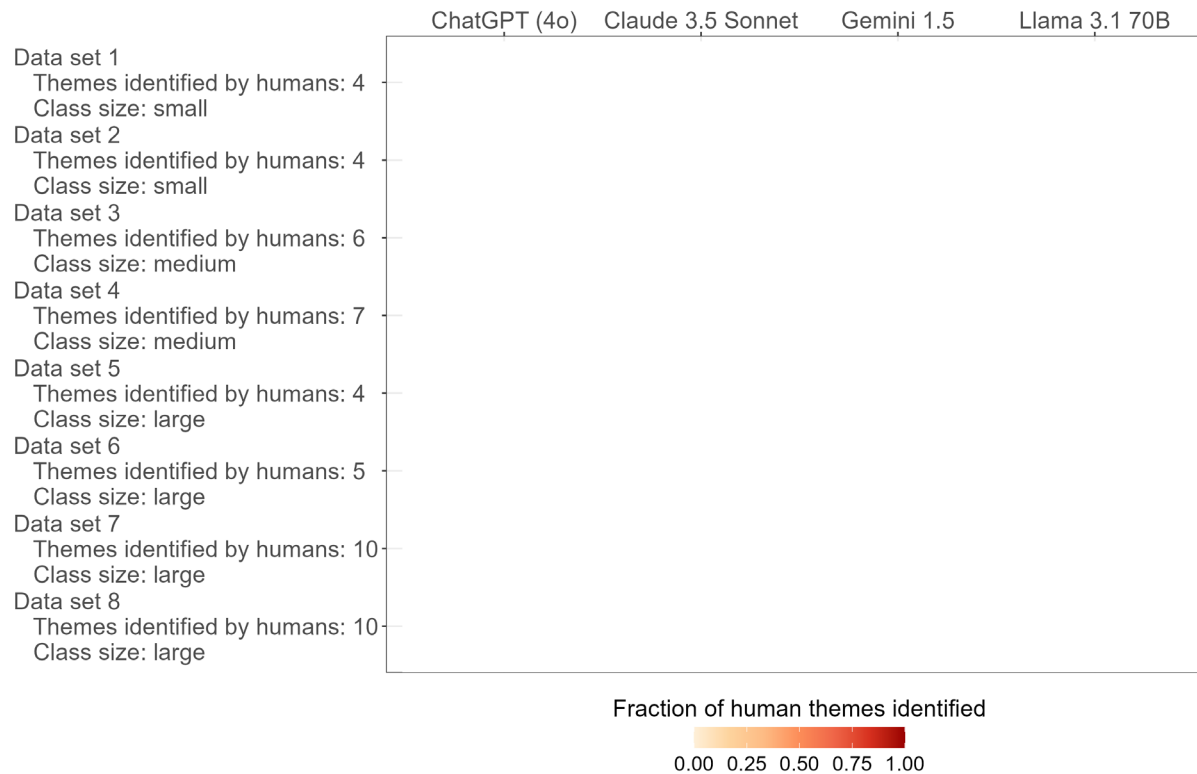




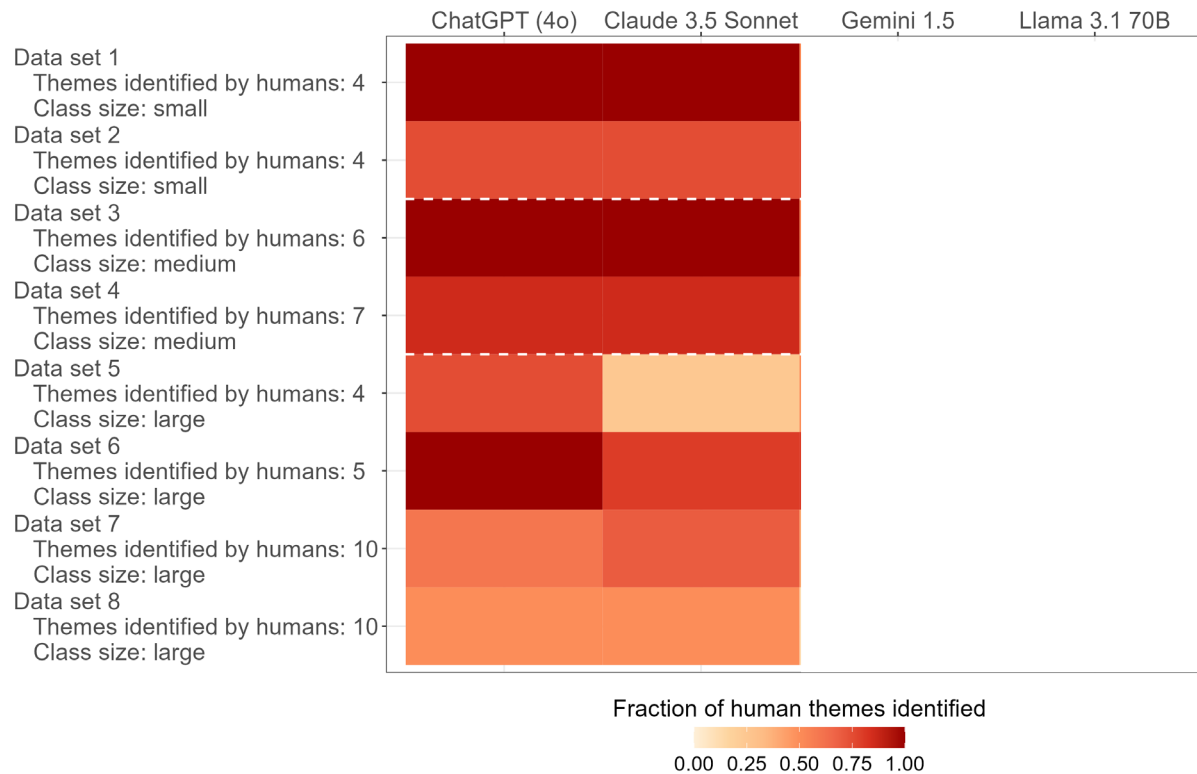
How do the LLMs do in comparison to the instructors?



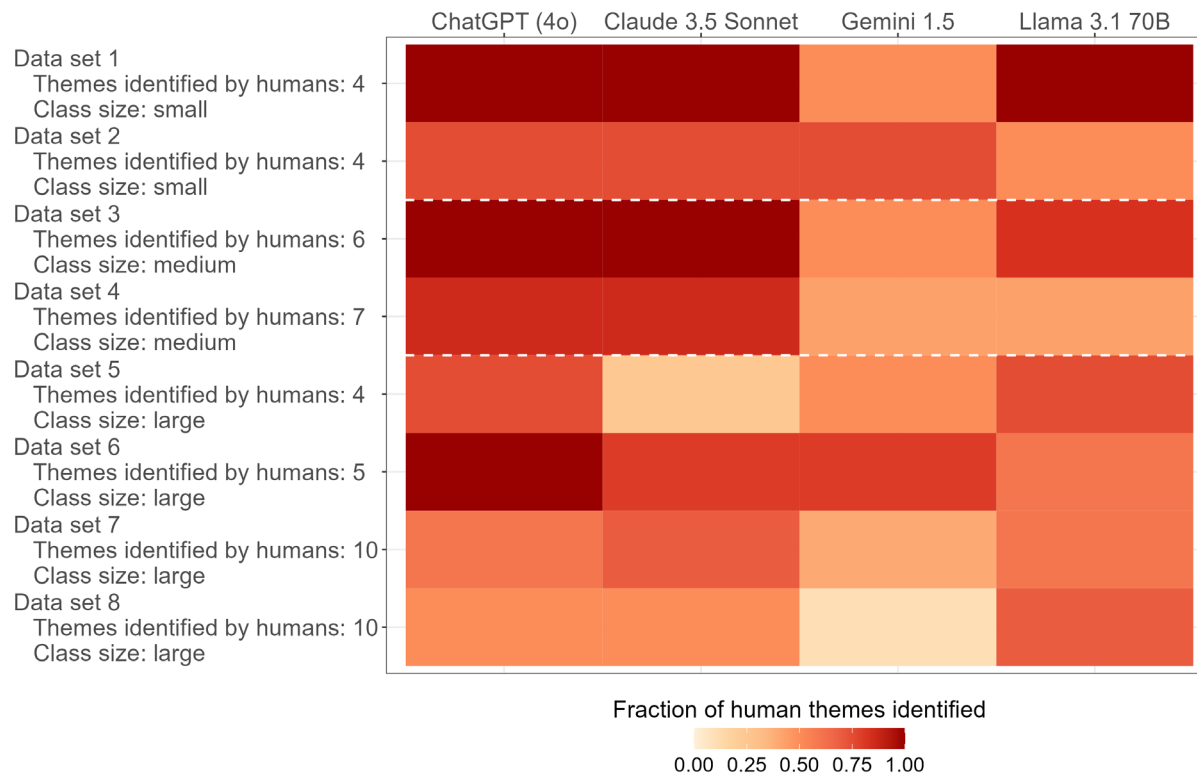
Results



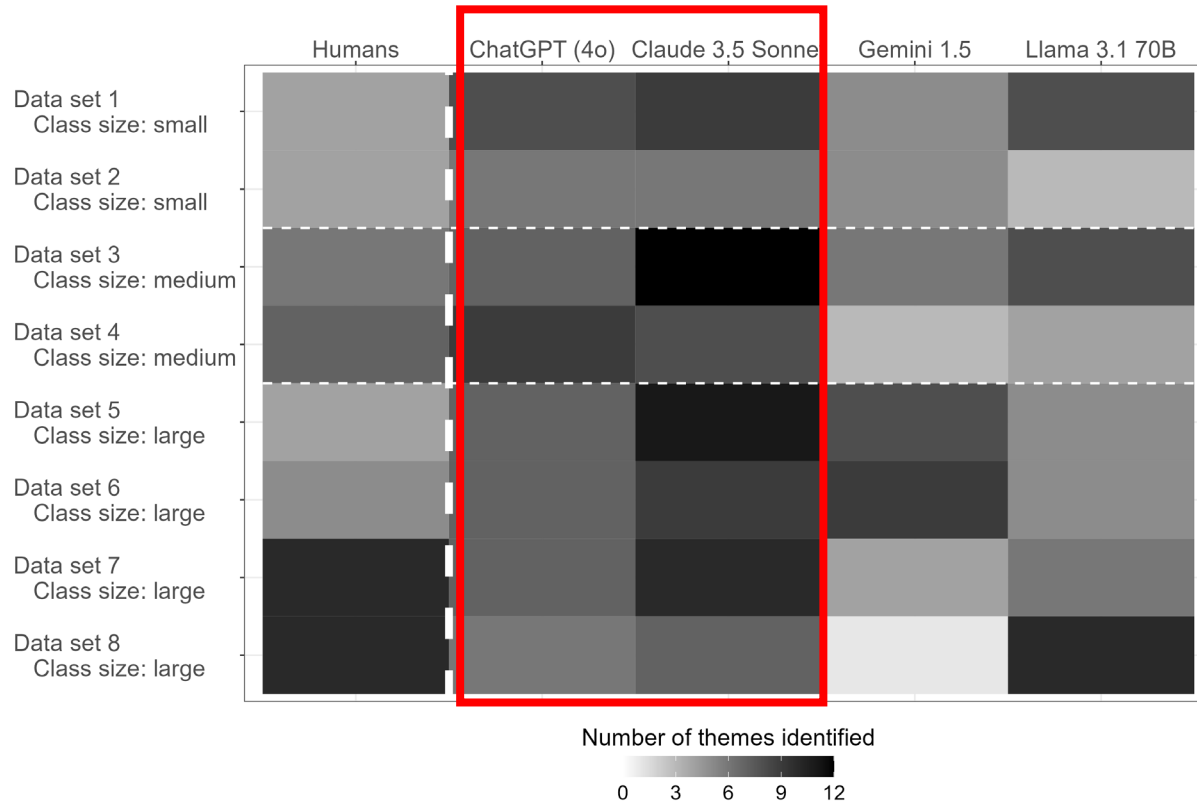
Results: ChatGPT and Claude generally detect most themes



Results: Gemini and Llama tend to miss many of the human themes



What about the themes not identified by humans?



Are the extra themes hallucinations?

- Analysis still in progress



Are the extra themes hallucinations?

- Analysis still in progress
- Probably not
 - Likely result of only 1 or 2 instructors identifying theme

Takeaways

- LLMs can be useful for summarizing student feedback

Takeaways

- LLMs can be useful for summarizing student feedback
- ChatGPT and Claude showed better performance than Gemini and Llama

Takeaways

- LLMs can be useful for summarizing student feedback
- ChatGPT and Claude showed better performance than Gemini and Llama
- Instructors likely want to ask for a summary multiple times or spot check results

Takeaways

- LLMs can be useful for summarizing student feedback
- ChatGPT and Claude showed better performance than Gemini and Llama
- Instructors likely want to ask for a summary multiple times or spot check results

Questions?

nicholas.young@uga.edu