

# UCSC QUD Empowered LLM Co-Writing Tool

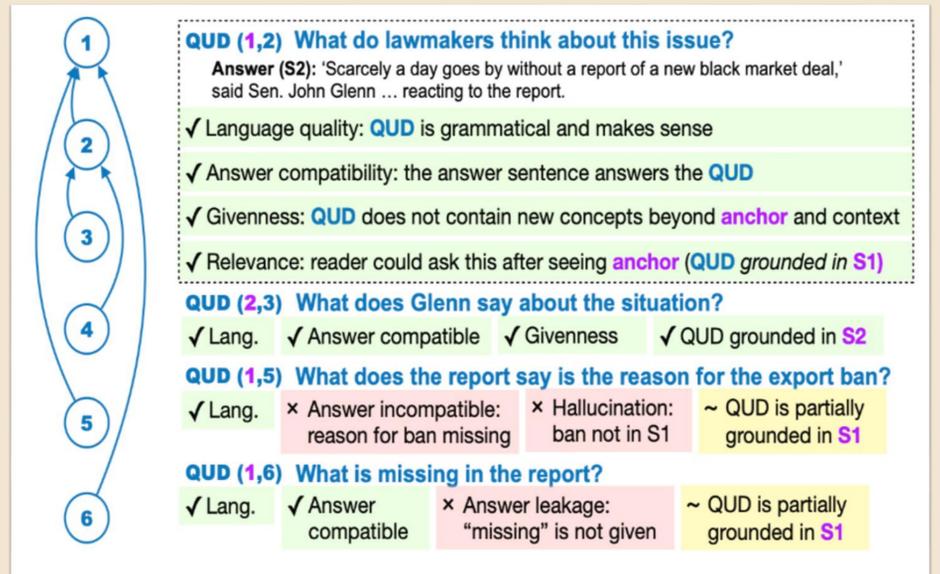
Ting-Yu Chou, Yousuf Golding, Kiara LaRocca, Neng Wan

## INTRO

Many students use large language models (LLMs) as part of their writing process, but the impact on learning and comprehension remains unclear. **We aim to identify where LLMs can provide support without diminishing the learner's active reasoning and skill development.**

We leveraged **Question Under Discussion (QUD)** to provide targeted, question-based feedback that guides students to reason through their summaries to find missing content, without relying on the model to generate answers.

## QUD FRAMEWORK



## STUDIES

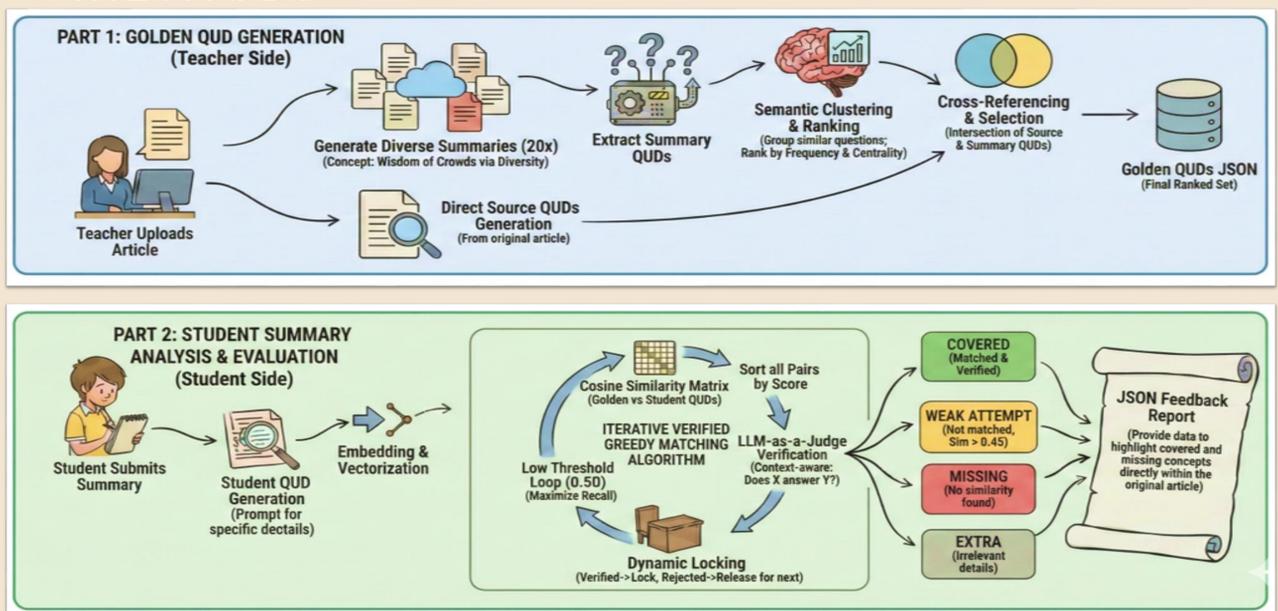
### Study 1: Question Extraction Study

- ❖ Validates if AI frequency ranking correlates with human judgements.
- ❖ Participants rate ~15 questions per article on relevance.
- ❖ 10 raters per article across 10 articles.

### Study 2: Human Preference Study

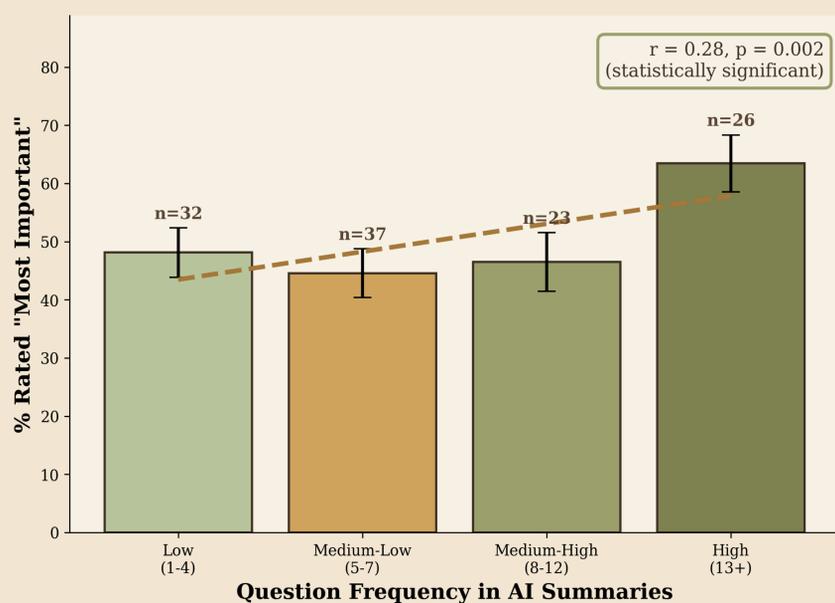
- ❖ Participants use both QUD system and Vanilla LLM feedback.
- ❖ Task: Read -> Summarize -> Receive Feedback -> Revise
- ❖ Metrics: Summary quality, user understanding, ownership and ease of use.

## METHODS

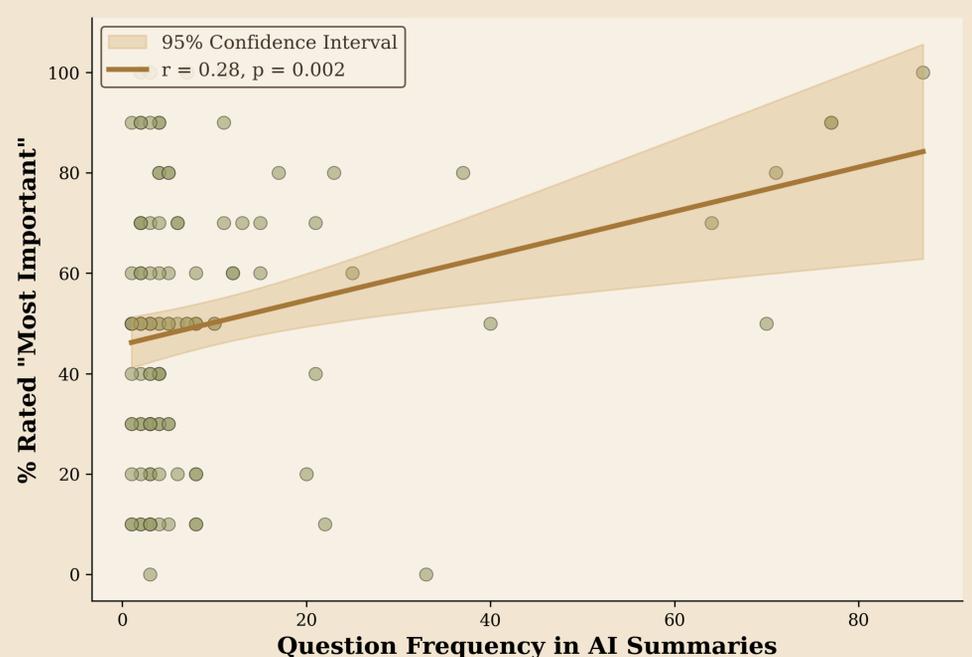


## RESULTS

Higher AI Frequency → Higher Human Importance Ratings



AI Prediction vs Human Rating (N=118 questions)



## CONCLUSION & ONGOING WORK

- ❖ Analysis of ~118 questions over 10 articles showcased a statistically significant correlation between AI-predicted significance and human ratings.
- ❖ Per-article correlation varied (-0.47 to 0.72) but showcased positive relationships in 80% of the articles.
- ❖ The variation suggests the proposed method is feasible but requires refinement to better capture meaningful aspects of human judgement.
- ❖ There are at least 7 more potential tool evaluations we can do
- ❖ We would like to train our own summarization model
- ❖ Testing the confidence threshold of the system might give us insights into comparing the summaries more accurately and efficiently