# Qeios

Peer Review

# Review of: "Enhancing Annotated Bibliography Generation with LLM Ensembles"

**Alexandru Tăbușcă[1]**

1. Romanian-American University, Romania

The submitted research article presents a novel approach to improving the generation of annotated bibliographies using ensembles of Large Language Models (LLMs). The paper is well-structured and provides a comprehensive overview of the methodology, experimental setup, and results. This review will highlight the strengths and weaknesses of the paper, providing a balanced critique of its contributions and limitations.

*Among the strengths of the article, I can list:*

**Innovative Approach**: The paper introduces an innovative three-tier LLM ensemble architecture for generating annotated bibliographies. This approach leverages the strengths of multiple LLMs in different roles—controllable text generation, evaluation, and summarization—to enhance the quality and reliability of the generated annotations. The use of LLM ensembles to address the limitations of individual models is a significant contribution to the field.

**Comprehensive Methodology**: The methodology is well-detailed, with clear explanations of the three-tier architecture. The paper describes the roles of each LLM in the ensemble, the hyperparameter settings used for generating diverse responses, and the evaluation criteria employed by the LLM acting as a judge. This thorough explanation allows readers to understand the process and potentially replicate the study.

**Experimental Validation**: The paper provides robust experimental validation of the proposed approach. The results demonstrate that the ensemble methods (Top M Responses and Top Temperature) outperform individual LLMs in terms of readability and conciseness. The inclusion of quantitative metrics, such as average sentence length and readability score, adds credibility to the findings.

**Clear Presentation of Results**: The results are presented clearly, with tables and figures that effectively illustrate the improvements achieved by the ensemble methods. The comparison between the baseline individual LLM and the ensemble methods is well-articulated, highlighting the advantages of the proposed approach.

**Potential for Practical Applications**: The study highlights the potential of LLM ensembles to automate complex scholarly tasks like annotated bibliography generation. This has significant implications for research productivity, as it can save time and improve the quality of annotations, making it a valuable tool for researchers.

*Nevertheless, there are also several weaknesses that are present:*

**Limited Scope of Evaluation**: While the paper provides a comprehensive evaluation of the proposed approach, the scope of the evaluation is somewhat limited. The experiments are conducted using specific LLM models (Gemini 1.5 ash and Gemini 1.5 pro), and it is unclear how the approach would perform with other models or in different domains. A broader evaluation would strengthen the generalizability of the findings.

**Potential Biases in LLMs**: The paper acknowledges the potential biases within LLMs but does not provide a detailed analysis of how these biases might affect the generated annotations. Addressing this issue in more depth would enhance the credibility of the study, as biases in training data and model architecture can significantly impact the quality and fairness of the generated content.

**Dependence on LLM-as-a-Judge**: The reliance on an LLM acting as a judge to evaluate the generated annotations introduces a potential weakness. While the LLM-as-a-judge approach aims to achieve greater objectivity, it is still subject to the limitations and biases of the underlying model. The paper could benefit from a discussion on the limitations of this approach and potential strategies to mitigate its impact.

**Lack of Qualitative Analysis**: The paper focuses primarily on quantitative metrics to evaluate the performance of the ensemble methods. Including a qualitative analysis of the generated annotations would provide additional insights into the strengths and weaknesses of the approach. For example, examining specific examples of annotations and discussing their relevance, coherence, and accuracy in detail would add depth to the evaluation.

**Future Research Directions**: While the paper concludes with recommendations for future research, these suggestions are somewhat general. Providing more specific directions for addressing the limitations

identified in the study, such as exploring different ensemble techniques or evaluating the approach in various domains, would be beneficial for guiding future work in this area.

## Conclusion

Overall, the paper is a well-written and insightful work that makes a significant contribution to the field of automated scholarly tasks. The innovative use of LLM ensembles to improve the quality of annotated bibliographies is a notable achievement. However, the study would benefit from a broader evaluation, a deeper analysis of potential biases, and a more detailed qualitative assessment. Despite these limitations, the paper provides a strong foundation for future research and has the potential to significantly enhance research productivity through automation.

## Declarations

**Potential competing interests:** No potential competing interests to declare.