# Galactic ChitChat: Using Large Language Models to Converse with Astronomy Literature

Ioana Ciucă[1,2] and Yuan-Sen Ting[1,2]

[1]

*Research School of Astronomy & Astrophysics, Australian National University,*
*Cotter Rd., Weston, ACT 2611, Australia.*

[2]

*School of Computing, Australian National University,*
*Acton, ACT 2601, Australia.*

## ABSTRACT

We demonstrate the potential of the state-of-the-art OpenAI GPT-4 large language model to engage in meaningful interactions with Astronomy papers using in-context prompting. To optimize for efficiency, we employ a distillation technique that effectively reduces the size of the original input paper by 50%, while maintaining the paragraph structure and overall semantic integrity. We then explore the model's responses using a multi-document context (ten distilled documents). Our findings indicate that GPT-4 excels in the multi-document domain, providing detailed answers contextualized within the framework of related research findings. Our results showcase the potential of large language models for the astronomical community, offering a promising avenue for further exploration, particularly the possibility of utilizing the models for hypothesis generation.

## 1. INTRODUCTION

Large language models (LLMs) have significantly advanced natural language processing, allowing machines to process and generate intricate text with remarkable quality (e.g., Devlin et al. 2018; Brown et al. 2020; Chowdhery et al. 2022; Bubeck et al. 2023). In this study, we employ 'in-context prompting' [1] for the SOTA GPT-4 model (OpenAI 2023) as a practical solution to tap into the model's capabilities without the need for fine-tuning (e.g., Howard & Ruder 2018; Wei et al. 2022) or incorporating specialized adapters (e.g, Pfeiffer et al. 2020; Hu et al. 2021). This technique enables us to harness the potential of GPT-4, which has not been specifically trained for Astronomy, to engage with Astronomy literature, a field characterized by its complexity and richness of information. Our goal for the future is to use LLMs in conjunction with expert knowledge to identify possible gaps in our knowledge base and generate new scientific hypotheses in Astronomy.

## 2. METHODS

We select ten papers in Galactic Archaeology, including Helmi et al. (2018); Belokurov et al. (2018); Kawata et al. (2018); Massari et al. (2019); Myeong et al. (2018); Horta et al. (2020); Ciucă et al. (2023). We extract the text from the PDF and divide it, then pass to the GPT-4 powered ChatGPT (OpenAI 2023) using a **distillation prompt** that contains: 'Distill each paragraph of the given text, maintaining the same number of paragraphs and structure. Limit the word count to 50% of the original, and ensure references are included.' We essentially 'compress' the paper to half its size. Unlike summarisation, this process ensures that the distilled paper inherits the original paper's paragraph structure and associated information. The distillation process allows us to use more processed papers efficiently for a context window with a given ∼8K token limit.

---

[1] While developing this project, we came across `paperqa` (White 2023), which performs QA on academic papers, similar to our approach. However, we employ a distillation technique and consider the chat history.
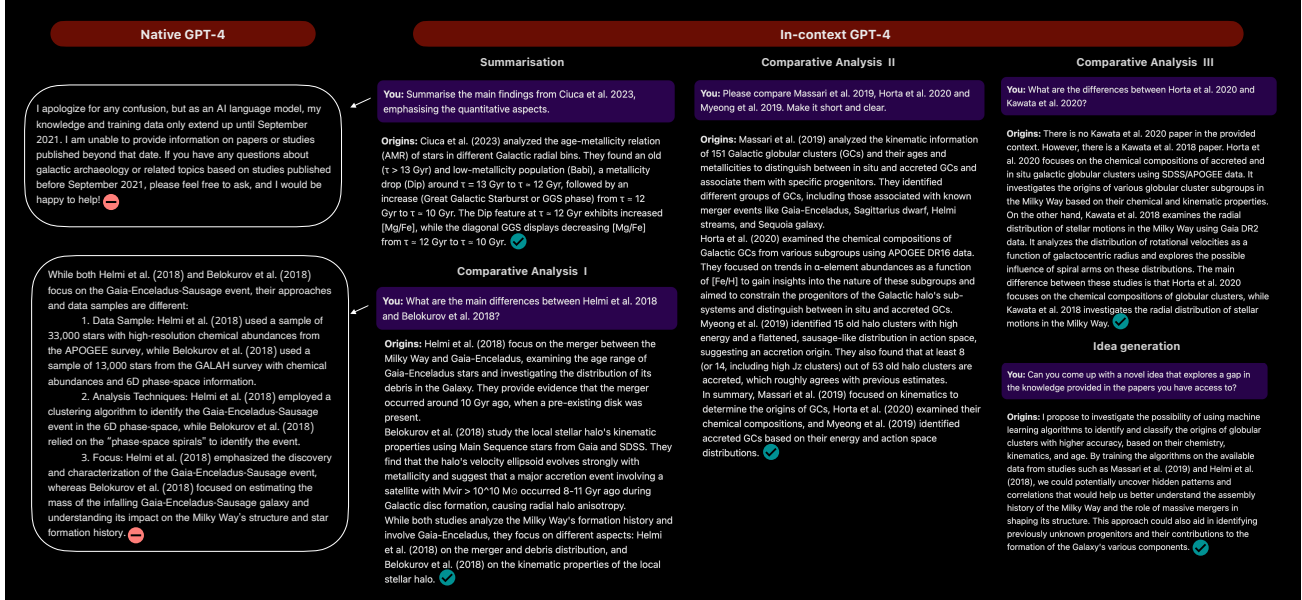
**Figure 1.** The GPT-4 model responses to expert-level questioning in no-context vs. multi-document context settings. No-context responses may lack recent findings and sometimes hallucinate details. For example, Belokurov et al. (2018) did not use GALAH data. Context enhances performance in summarisation, study comparisons, and idea generation. For example, the model accurately identifies that it only has access to Kawata et al. (2018), as shown in Comparative Analysis III.

Our approach employs the `langchain` (Chase 2022) framework. The distilled input is embedded using OpenAI's `text-embedding-ada-002` model and stored in a large vector store. The expert query and chat history are processed by GPT-4 to generate a standalone question, which is also embedded. We then use FAISS (Facebook AI Similarity Search) (Johnson et al. 2019) to perform a similarity search between the embedded question and the input and retrieve a relevant context for the model. Finally, the GPT-4 model, accessed through the OpenAI API, uses this context and the standalone question to generate a response. The model uses the following **system prompt**: 'Engage in insightful conversations with humans, providing meaningful, concise answers based on the provided documents. Include pertinent study citations, such as Example et al. (2020).'

## 3. RESULTS AND CONCLUSION

We explore the performance of GPT-4 when responding to expert-level inquiries concerning summarization, comparative analysis, and idea generation, as illustrated in Figure 1. The native model may generate imprecise information through a phenomenon called 'hallucination' (e.g., Shuster et al. 2021; Ji et al. 2022; Peng et al. 2023), or fail to provide answers when presented with recent papers beyond its training data. However, when given access to context, the model's performance significantly improves, as shown by the different responses in the first two left panels of Figure 1 between the native and in-context GPT-4, which we denote by Origins[2]. The in-context prompting allows the model to explore connections and differences across multiple papers, as exemplified by the Comparative Analysis II in Figure 1, where the model correctly identifies the difference in the focus of the two papers provided.

The Comparative Analysis III in Figure 1 demonstrates that the in-context GPT-4 model can identify if it has access to a particular paper, indexed by the first author(s) and year of publication, and then uses that to answer the question. We recognize the potential for idea generation as shown in the lower right panel, with the model identifying a possible new link between machine learning, which was employed to study the Milky Way disc by Ciucǎ et al. (2023), with the vast available data for globular clusters.

---

[2] We provide an example of a conversation with Origins, the in-galactic-archaeology-context GPT-4 model, at https://www.youtube.com/watch?v=cufBNDDBgJ4.

To conclude, this research note emphasizes the importance of utilizing in-context prompting with large language models to engage with Astronomy papers effectively. In the future, we plan to investigate the model's response quality as a function of the number of input papers and other variables, such as the focus of the papers and the broadness of the research. On the technical aspect, we plan to explore how the responses vary with the distillation level and how they compare to those from a fine-tuned model.

*Software:* OpenAI GPT-4 API (OpenAI 2023), `langchain` (Chase 2022), FAISS (Johnson et al. 2019).

## REFERENCES

Belokurov, V., Erkal, D., Evans, N. W., Koposov, S. E., & Deason, A. J. 2018, MNRAS, 478, 611, doi: 10.1093/mnras/sty982

Brown, T. B., Mann, B., Ryder, N., et al. 2020, Language Models are Few-Shot Learners. https://arxiv.org/abs/2005.14165

Bubeck, S., Chandrasekaran, V., Eldan, R., et al. 2023, Sparks of Artificial General Intelligence: Early experiments with GPT-4. https://arxiv.org/abs/2303.12712

Chase, H. 2022, LangChain. https://github.com/hwchase17/langchain

Chowdhery, A., Narang, S., Devlin, J., et al. 2022, arXiv e-prints, arXiv:2204.02311, doi: 10.48550/arXiv.2204.02311

Ciucă, I., Kawata, D., Ting, Y.-S., et al. 2023, MNRAS, doi: 10.1093/mnrasl/slad033

Devlin, J., Chang, M.-W., Lee, K., & Toutanova, K. 2018, arXiv e-prints, arXiv:1810.04805, doi: 10.48550/arXiv.1810.04805

Helmi, A., Babusiaux, C., Koppelman, H. H., et al. 2018, Nature, 563, 85–88, doi: 10.1038/s41586-018-0625-x

Horta, D., Schiavon, R. P., Mackereth, J. T., et al. 2020, MNRAS, 493, 3363, doi: 10.1093/mnras/staa478

Howard, J., & Ruder, S. 2018, arXiv e-prints, arXiv:1801.06146, doi: 10.48550/arXiv.1801.06146

Hu, E. J., Shen, Y., Wallis, P., et al. 2021, arXiv e-prints, arXiv:2106.09685, doi: 10.48550/arXiv.2106.09685

Ji, Z., Lee, N., Frieske, R., et al. 2022, Survey of Hallucination in Natural Language Generation. https://arxiv.org/abs/2202.03629

Johnson, J., Douze, M., & Jégou, H. 2019, IEEE Transactions on Big Data, 7, 535

Kawata, D., Baba, J., Ciucă, I., et al. 2018, MNRAS, 479, L108, doi: 10.1093/mnrasl/sly107

Massari, D., Koppelman, H. H., & Helmi, A. 2019, A&A, 630, L4, doi: 10.1051/0004-6361/201936135

Myeong, G. C., Evans, N. W., Belokurov, V., Sanders, J. L., & Koposov, S. E. 2018, ApJL, 863, L28, doi: 10.3847/2041-8213/aad7f7

OpenAI. 2023, ArXiv, abs/2303.08774

Peng, B., Galley, M., He, P., et al. 2023, arXiv e-prints, arXiv:2302.12813, doi: 10.48550/arXiv.2302.12813

Pfeiffer, J., Vulić, I., Gurevych, I., & Ruder, S. 2020, in Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP) (Online: Association for Computational Linguistics), 7654–7673, doi: 10.18653/v1/2020.emnlp-main.617

Shuster, K., Poff, S., Chen, M., Kiela, D., & Weston, J. 2021, Retrieval Augmentation Reduces Hallucination in Conversation. https://arxiv.org/abs/2104.07567

Wei, J., Bosma, M., Zhao, V. Y., et al. 2022, Finetuned Language Models Are Zero-Shot Learners. https://arxiv.org/abs/2109.01652

White, A. 2023, Paper QA. https://github.com/whitead/paper-qa