# SciSpace Copilot: Empowering Researchers through Intelligent Reading Assistance

**Trinita Roy, Asheesh Kumar, Daksh Raghuvanshi, Siddhant Jain, Goutham Vignesh, Kartik Shinde, Rohan Tondulkar**

SciSpace

(trinita,asheesh,daksh,siddhant,goutham,kartik,rohan)@typeset.io

## Abstract

We introduce SciSpace Copilot, an AI research assistant that helps in understanding and reading research papers faster by providing a plethora of features. Answering questions from a document has recently become popular using the Retrieval Augmented Generation (RAG) approach. Our tool uses an advanced question-answering pipeline to give accurate answers while also citing sources from the paper. We also provide other valuable features on scientific text, including generating explanations, generating summaries, adding notes and highlights, and finding related papers from our 280 million corpus. Our tool supports 75+ languages, making research more accessible across language barriers. Thousands of users use SciSpace Copilot on a daily basis by uploading their articles to understand research faster and better. Our tool can be accessed at this link: https://typeset.io.

Figure 1: Copilot feature window preview

## Introduction

Researchers today are facing a big problem of keeping up-to-date with latest research as approx. 100 papers are published daily in machine learning domain itself. They want to extract important insights from a publication and get answers to some important questions from the publication. Another problem is the trouble of understanding complex parts of research papers which requires referring to multiple resources. This increases the overall time of reading a paper, leading to less research data being consumed. Available existing tools don't address these problems.

Recent advancements in large language models (LLMs) have achieved remarkable progress in various tasks, displaying enhanced performance and versatility in answering questions across diverse domains. LLMs have been trained on terabytes of open access knowledge but are known to generate misinformation with inaccurate citations making it difficult to trust outputs from their own knowledge. The approach outlined in the paper (Lewis et al. 2020; Cai et al. 2022) introduces a novel approach that employs retrieval-augmented generation (RAG) models. These models combine pre-trained parametric and non-parametric memory elements with the objective of improving language generation. As a result, they achieve state-of-the-art results on open domain QA tasks responses that are more precise, varied,
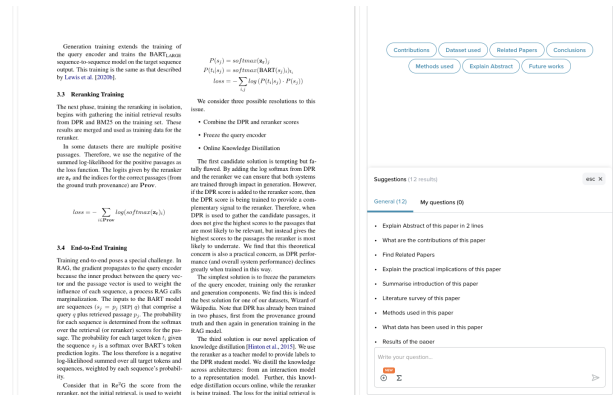
and fact-based. (Nguyen, Le, and Nguyen 2022) developed a "retriever-reader" framework for querying scientific texts, employing a sliding window technique to manage large context scales in documents by dividing them into manageable text blocks.

Our paper demonstrates an AI-powered Copilot assistant which has the following characteristics:

- Provide evidence-based responses and clarify intricate research elements, including texts and equations, facilitating comprehensive understanding.

- Streamline research efforts through content summarization, support for regional languages, and improved engagement with highlighting, note-taking, and personalized recommendations.

## Methodology

This section gives details about the core components forming the building blocks of features in SciSpace Copilot:

- **Parsing**: We use various parsing tools like Grobid[1], PyMuPDF[2] and Science-Parse[3] to get a clean textual representation for research pdfs.

---

[1] https://github.com/kermitt2/grobid
[2] https://github.com/pymupdf/PyMuPDF
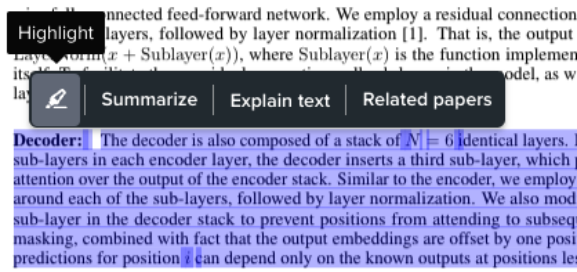[3] https://github.com/allenai/science-parse

Figure 2: Functionalities on text selection.

- **Embedding generation**: We use state-of-the-art open-source embedding models to generate embeddings on chunks of text data from the parsing output. These embeddings are stored in a vector database for later use in various features.

- **Context extraction**: We use structured and unstructured approaches to extract relevant context for a given question. We parse through the parsing outputs to give importance to the structure of the paper and also use vector search to get relevant text.

- **LLMs and prompt engineering**: We use generative LLMs to generate required output for each use case. We use tools such as OpenAI (OpenAI 2023), Anthropic [4], Bard [5], etc., as well as open source models for certain use cases. Prompt engineering (White et al. 2023; Oppenlaender, Linder, and Silvennoinen 2023; Sorensen et al. 2022) is done to get the desired outputs. Relevant contexts are passed from the research paper to get factual responses and reduce the possibilities of hallucination (Li et al. 2023; McKenna et al. 2023; Mündler et al. 2023).

## Core Functionalities

### Features on Text Selection

SciSpace Copilot gives the options of Explain text, Summarize, Related papers and highlight to users on selecting a piece of text in the paper as shown in Figure 2.

- **Summarize**: The "Summarize" option will generate a short summary for a long piece of text to help users to read papers faster by getting summaries of each section.

- **Explaining complex text and equations**: Users can use the "Explain text" option to understand research paper text and use the snip tool to understand formulae, equations, and tables by selecting the area containing the required information as shown in Figure 3.

- **Related papers**: This option recommends papers from a corpus of 200+ million research papers similar to the selected text which helps in discovery of more papers.

- **Highlight**: Users can apply colorful highlights and add accompanying notes to any uploaded research paper.

---

[4]https://www-files.anthropic.com/production/images/Model-Card-Claude-2.pdf

[5]https://blog.google/intl/en-africa/products/explore-get-answers/an-important-next-step-on-our-ai-journey/
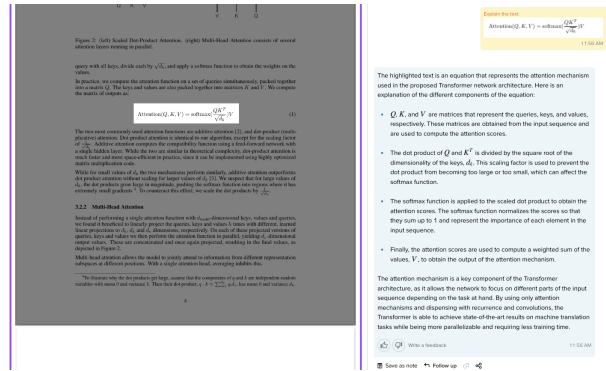


Figure 3: Explanations based on selected texts, tables and equations.

## Question Answering

Users inquire about research paper concepts, receiving answers with precise citations via SciSpace Copilot. Users can locate the source in paper as shown in Figure 4. Follow-up questions are supported, and answers can be stored as Notes for future reference.

## Regional Language Support

Copilot enables language-agnostic research paper comprehension. Users can ask and receive answers, as well as read papers, in their regional language among 75+ supported options.

## Conclusion

SciSpace Copilot is a generative AI tool that streamlines the process of consuming research by rapidly distilling insights from papers and making complex content accessible to newer researchers. It saves time and increases research consumption with features like regional language support, which is particularly beneficial for users in Asia and Africa. The tool keeps researchers up-to-date and speeds up their work.
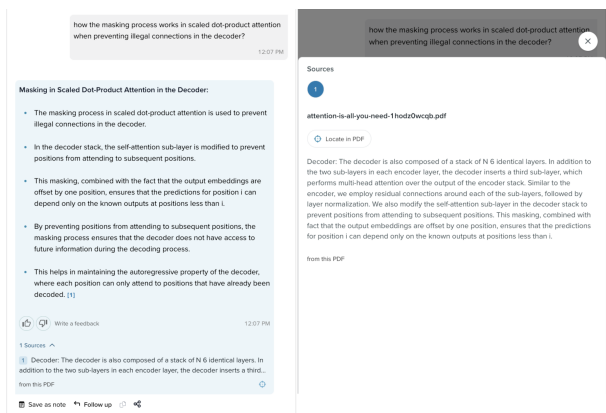


Figure 4: Questions answering with navigable citations.

# References

Cai, D.; Wang, Y.; Liu, L.; and Shi, S. 2022. Recent advances in retrieval-augmented text generation. In *Proceedings of the 45th International ACM SIGIR Conference on Research and Development in Information Retrieval*, 3417–3419.

Lewis, P.; Perez, E.; Piktus, A.; Petroni, F.; Karpukhin, V.; Goyal, N.; Küttler, H.; Lewis, M.; Yih, W.-t.; Rocktäschel, T.; et al. 2020. Retrieval-augmented generation for knowledge-intensive nlp tasks. *Advances in Neural Information Processing Systems*, 33: 9459–9474.

Li, J.; Cheng, X.; Zhao, W. X.; Nie, J.-Y.; and Wen, J.-R. 2023. HaluEval: A Large-Scale Hallucination Evaluation Benchmark for Large Language Models. *arXiv e-prints*, arXiv–2305.

McKenna, N.; Li, T.; Cheng, L.; Hosseini, M. J.; Johnson, M.; and Steedman, M. 2023. Sources of Hallucination by Large Language Models on Inference Tasks. *arXiv preprint arXiv:2305.14552*.

Mündler, N.; He, J.; Jenko, S.; and Vechev, M. 2023. Self-contradictory Hallucinations of Large Language Models: Evaluation, Detection and Mitigation. *arXiv preprint arXiv:2305.15852*.

Nguyen, D.-H.; Le, N.-K.; and Nguyen, M. L. 2022. Exploring Retriever-Reader Approaches in Question-Answering on Scientific Documents. In *Asian Conference on Intelligent Information and Database Systems*, 383–395. Springer.

OpenAI. 2023. GPT-4 Technical Report. *arXiv preprint arXiv:2303.08774*.

Oppenlaender, J.; Linder, R.; and Silvennoinen, J. 2023. Prompting ai art: An investigation into the creative skill of prompt engineering. *arXiv preprint arXiv:2303.13534*.

Sorensen, T.; Robinson, J.; Rytting, C. M.; Shaw, A. G.; Rogers, K. J.; Delorey, A. P.; Khalil, M.; Fulda, N.; and Wingate, D. 2022. An information-theoretic approach to prompt engineering without ground truth labels. *arXiv preprint arXiv:2203.11364*.

White, J.; Fu, Q.; Hays, S.; Sandborn, M.; Olea, C.; Gilbert, H.; Elnashar, A.; Spencer-Smith, J.; and Schmidt, D. C. 2023. A prompt pattern catalog to enhance prompt engineering with chatgpt. *arXiv preprint arXiv:2302.11382*.