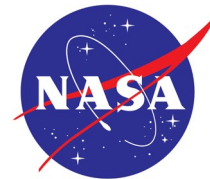# Artificial Intelligence and Large Language Models

*Kelly Lockhart and the ADS Team*

ADS Users Group Meeting, 16-17 Nov. 2023

# Large Language Models

- Data enrichment
  - Extension of our current machine learning efforts
- Data discovery
  - A new way of searching and synthesizing information

# Data Enrichment

LLM queries are one technique of several for current back-office data enrichment tasks.

- Planetary names (in development)
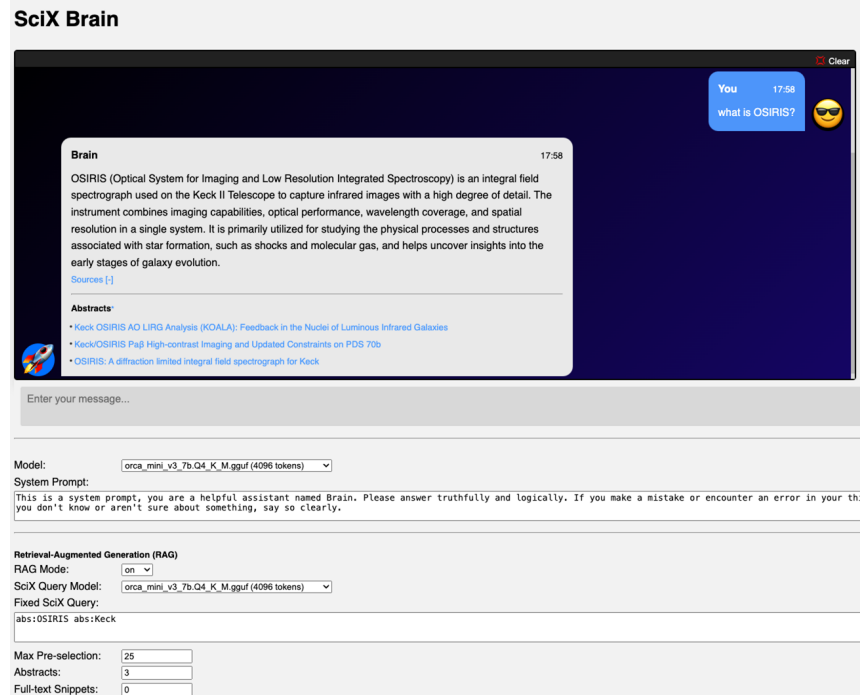- Named entity recognition (experimental)



Figure: Extract of an annotated paper from the 2022 DEAL Named Entity Recognition shared task corpus.

# Data Enrichment: avenues to explore

- Very short (1-2 sentence) paper summaries
  - Useful for myADS email notifications
- Paper summaries written at a level for the general public, for undergrads, or for K-12 education

# Data Discovery: SciX Brain chatbot

- Experimental, restricted access
  - Developed by Sergi Blanco-Cuaresma
- Test bed for LLM techniques
  - Retrieval augmented generation (RAG)
  - Comparison of various open-source LLMs
  - Architectures
  - Grammars
  - Natural language → structured Solr queries (student project)

# SciX Brain chatbot caveats

Concerns about opening this more widely:

- Hallucinations
    - → Though some LLMs are more resilient to this
- Accuracy, reputational impact
- Exposing protected content to users
    - → Though RAG mode could use only open access articles
- Cost (either purchasing GPUs or via AWS)

→ Maintaining trust, of publishers, of our users, is paramount

# SciX Brain demo