

### Abstract

The AI for Curation project aims to integrate advanced LLM models into various aspects of our data curation workflow, enhancing the efficiency and accuracy from submission to user interaction. This initiative will impact multiple touchpoints, including data ingestion, curation processes, and user engagement with our curated datasets, affecting various centers, domains, and a broad user base. First, we are developing tools capable of parsing incoming data across all formats and utilizing LLMs to convert unstructured data into structured, standardized formats. This process streamlines curation by converting data into community-standard formats like ISA-Tab, significantly reducing the time and effort required by curators. This allows curators to allocate more resources to scientific analysis rather than data formatting tasks. Second, we are implementing AI/ML models to automate and enhance the accuracy of data validation and verification. These models ensure that the curated data adheres to high standards of quality and reliability, benefiting researchers and data users by providing them with rigorously verified datasets. Finally, we are developing a conversational agent (chatbot) that interfaces with our extensive repository of curated scientific studies on the Open Science Data Repository (OSDR). The chatbot enhances data discoverability by assisting users in navigating the knowledge base and referencing relevant studies. This improvement in accessibility makes scientific data more available to the community, thereby promoting the principles of open science.

### Data Alignment & Verification

Generative AI can be applied throughout the entire curation pipeline, offering numerous use cases, especially in data curation. One of the initial challenges addressed was the alignment of data received from researchers. Data often arrives in various formats and with disparate terminologies. To tackle this, we have developed tools that map and align these terms to standardized ontologies. For instance, while fields like "sex" and "genotype" may be consistent, others may use different terms for the same concepts.

Combining programmatic techniques with generative AI provides a robust toolkit for curating datasets. For example, we utilize similarity metrics such as the Levenshtein distance to match strings that are similar, handling a significant portion of the workload. For terms that use different words but are semantically similar, we leverage generative AI to understand the context and find the appropriate matches, enhancing the precision and utility of our curated datasets.

To further streamline the data alignment process, our graphical user interface (GUI) provides a suite of intuitive tools. These GUI capabilities are integral to ensuring that the data not only meets standardization requirements but also aligns perfectly with the users' specific research needs and contexts.

### Enhancing Study Discoverability

We've implemented a chatbot to improve how researchers and users interact with our data. The chatbot integrates directly with our databases, pulling information in real-time to ensure that responses are always up-to-date and accurate. This facilitates dynamic interaction based on the latest data available. It also employs natural language processing to summarize intricate research findings or explain complex scientific concepts, making the information more accessible and easier to understand for a broader audience.

Additionally, the chatbot incorporates a Retrieval-Augmented Generation (RAG) model, which enhances its ability to generate precise and contextually relevant responses. The RAG model first retrieves relevant information from a vast database of scientific literature before producing a coherent reply, thereby improving the chatbot's effectiveness in navigating and discovering scientific studies. These integrated features ensure that critical information is not only more discoverable but also more understandable and secure within the scientific community.

This chatbot initiative aligns with our commitment to open science, facilitating broader access to knowledge and fostering a collaborative environment for innovation.

We extend our thanks to the Ames Life Sciences Data Archive (ALSDA), NASA Biological Institutional Scientific Collection (NBISC), and NASA GeneLab, integral projects of the Open Science Data Repository (OSDR), for their invaluable contributions to this project. Additionally, we are grateful to the AI for Life Science (AI4LS) group for their innovative approaches that have significantly enhanced our data analysis processes. Our appreciation also goes out to the Analysis Working Group (AWG), a team of dedicated volunteers, whose expertise and commitment have been instrumental in refining our research outputs.

# AI Curation Methods for NASA Scientific Data

## Walter Alvarado<sup>1</sup>, Lauren Sanders<sup>1</sup>, Hari Parthasarathy<sup>2</sup>, Harlan Phillips<sup>2</sup>, Kayvon Crenshaw<sup>3</sup>, Samrawit Gebre<sup>1</sup>, Sigrid Reinsch<sup>1</sup>, James Casaletto<sup>4</sup>, and Sylvain Costes<sup>1</sup>

<sup>1</sup>NASA Ames Research Center, Mountain View, CA, <sup>2</sup>University of California, Berkeley, CA, <sup>3</sup>University of Illinois at Urbana-Champaign, Champaign, IL, <sup>4</sup>Blue Marble Space Institute of Science, Seattle, WA

### Structured Outputs from LLM

In our efforts to enhance the AI for Curation project, we've developed advanced tools for automated tagging that significantly streamline the curation process. These tools are adept at identifying patterns within diverse datasets, automatically suggesting context-relevant metadata tags. This capability not only line is not include facilitates the initial organization of newly ingested data but also allows for the retrospective curation and categorization of existing datasets. Once data is aligned with community stanwas done on astronauts, as not done on dards, these tools apply contextual tags to meticulously organize the information, rendering it AI-ready for subsequent training processes. For example, they can generate specific tags for studies involving astronauts or particular experimental con-Pydantic, a Python library for data validation and settings management, facilitates the creation of schemas that ensures responses from large language ditions, ensuring a high level of precision and relevance in data models (LLMs) conform to a specified structure. Here, a Pydantic class is used to define the structure for automated tagging of OSDR studies, illustrathandling. This method greatly enhances the accessibility and ing how data is organized and standardized for AI processing. utility of the curated data, supporting more efficient research and analysis.

Column in File1	Best Match in File2	Similarity Score	Match Status				
Euthanasia method	Euthanasia method	100	Above Threshold				
Sample preservation method	Sample preservation method	100	Above Threshold				
Biospecimen category	Biospecimen category	100	Above Threshold				
Sex	Sex	100	Above Threshold				
Genotype	Genotype	100	Above Threshold	C.	last Format.		
Absorbed dose rate	Absorbed dose rate	97	Above Threshold	56	elect Format:		
Number of fractions	Number of fractions	97	Above Threshold		Metadata		
Time point of sacrifice post irradiation	****Time point of sacrifice post-irradiation	93	Above Threshold				
Additional notes	Additional Notes	91	Above Threshold	1	Ireatment 1	Ireat	
Absorbed dose rate (Unit)	Absorbed dose rate unit	88	Above Threshold	2	Status	Statu	
Total absorbed dose (Unit)	Total absorbed dose units	88	Above Threshold	_	D I.		
Time between exposures (Unit)	Time between exposures (Hour)	86	Above Threshold	3	Barcode	Box II	
Absorbed dose per fraction	Total absorbed dose per fraction	86	Above Threshold	4	Chronic Exposure Time	Chro	
Time between fractions (Days)	Time between fractions value	84	Above Threshold		Linear Francis	<b>F</b>	
Sample storage temperature	Sample storage temperature (Celsius)	84	Above Threshold	5	Linear Energy	Enerç	
Energy (Unit)	Energy unit	83	Above Threshold	6	Sex	Sex	
Absorbed dose per fraction (Unit)	Total absorbed dose per fraction units	82	Above Threshold		0	0	
Dose 1 (non-radiation)	No Match Found	0	Below Threshold		Species	Spec	
				8	Vendor	Subj€	
Matching: Dose 1 (non-r	adiation)			9	Strains	Strair	
LLM Match: Treatment 1	dose						
Match found!							
Dose 1 (non-radiation) <> Treatment 1 dose					Enable Al Matching		
Matching: Status (stts_	name)						
LLM Match: Status of sa	mples (Available, Not-availab	le)					
Match found!							
Match found!	-> Status of samples (Availab	le. Not-avai	ilable)				

Users can choose from predefined templates, which promote consistency across different datasets. Once a template is selected, the system automatically identifies and matches the fields according to the template structure. This automated matching significantly reduces the time and effort required to align data.

class Study(Bas	seModel):	
Title: str	= Field(description	="Study title.")
Accssion:	<pre>str = Field(descript</pre>	ion="Study accession
Cells: <mark>str</mark>	= Field(description	="Metadata values ·
		<pre>for the study's :</pre>
		in the metadata,
Astronaut:	<pre>bool = Field(descri</pre>	ption="True if the
		False if the stu
		astronauts.")

NASA Metadata	Sta	ndardization
	Ν	IBISC Intake Form
nent 1	1	Treatment 1
s of samples (Available, No	2	Sex
formation for Shipment (N	3	Status of samples (Available, Not-available)
ic exposure duration (Hou	4	Energy
у	5	Box Information for Shipment (Number, barcode, etc.) se
	6	Chronic exposure duration (Hours)
es/Type	7	Strain
ct supplier	8	Species/Type
	9	Subject supplier
Upload	Ex	Cel Export Excel
Find M	atcl	nes



Experimental Conter Bill Dynan Emory Universit Effect of GCR Sim Brookhaven/NSR Radiation Exp. Mus Musculus C57BL6/J Single Ion | Si - 28

To enhance the presentation and understanding of experimental data, we levearge generative AI to create graphical abstracts. Users can manually input experimental parameters using a user-friendly interface, ensuring direct control over the data visualization process. Additionally, the system supports file imports, as well as intelligent extraction techniques using generative AI to parse and interpret data directly from scientific papers. This capability automates the visualization of complex experiments, transforming dense textual information into clear, informative graphical abstracts.

### References

1. Colvin, S., Jolibois, E., Ramezani, H., Garcia Badaracco, A., Dorsey, T., Montague, D., Matveenko, S., Trylesinski, M., Runkle, S., Hewitt, D., & Hall, A. (2024). Pydantic (Version 2.9.0) https://github.com/pydantic/pydantic 2. Gao, Yunfan, et al. "Retrieval-augmented generation for large language models: A survey." arXiv preprint arXiv:2312.10997 (2023). 3. Dubey, Abhimanyu, et al. "The llama 3 herd of models." arXiv preprint arXiv:2407.21783 (2024).

Center. (2024).





### Auto-Generated Graphical Abstracts

		Experimenta	al Design	Clickable Tissue Button	
f Possible)	м	Month 1	Month 6	(Redirects to NBISC Collection)	Month 18
	÷ 2	Received 5 Sub	eceived 5 Su	bjects	
0.40 cGy	<u> </u>	Received 5 Sub	pjects Received 5 Su	bjects	
Tre	Patment: NR	Received 5 Sub	pjects Received 5 Su	Received 180	Subjects
		Received 5 Sub	pjects Received 5 Su	bjects	
0.75 cGy	<b>P</b>	Received 5 Sub	ojects	Received 180	Subjects
*		Received 5 Sub	pjects Received 5 Su	bjects	

4. Data are courtesy of the NASA Open Science Data Repository (OSDR). NASA Ames Research

