



THEORY AND PRACTICE
OF DIGITAL LIBRARIES

BIRDS 2025

Beyond Retrieval: A Vision of **Intelligent** Digital Libraries in the Large Language Model Era

DR. JIAN WU

ASSOCIATE PROFESSOR OF COMPUTER SCIENCE

OLD DOMINION UNIVERSITY, VIRGINIA, UNITED STATES

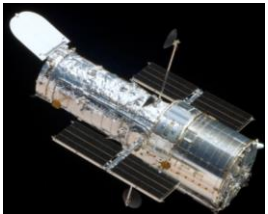


OLD DOMINION
UNIVERSITY®

Self-Introduction



2004: B.S. in
Physics and
Astronomy

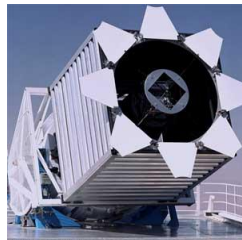


Hubble Space
Telescope



2011: Ph.D. in
Astronomy and
Astrophysics

Big Data



Sloan Digital Sky
Survey



2011-2017:
Postdoctoral
fellow, Information
Sciences and
Technology



2017-2018:
Assistant Teaching
Professor,
Information
Sciences and
Technology



2018-2025:
Assistant
Professor (tenure
track), Computer
Science



2025-: Associate
Professor
(effective July 25),
Computer Science

Scholarly Big Data + AI

CiteSeer^x

Conference Proceedings

ETDs

KGs

Tables

Journals

Technical Drawings

Keyphrases

Figures

BIRDS@TPDL'25

X: @FANCHYNA

BSKY: @FANCHYNA

Acknowledgments

CiteSeer^x

- Dr. C. Lee Giles (Pennsylvania State University)
 - PI of CiteSeerX
 - Eminent David Reese Professor of Information Sciences and Technology
 - Professor of Supply Chain and Information Systems
 - Director of the Intelligent Systems Research Laboratory
- Dr. Edward A. Fox (Virginia Tech)
 - Director of NDLTD
 - Eminent Professor of Computer Science
 - Director, Digital Library Research Laboratory

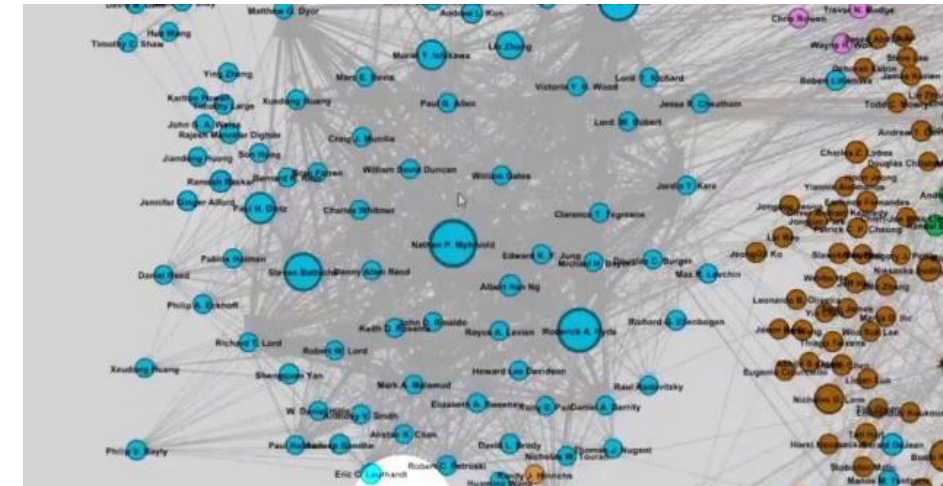




Is this the library we want to visit?

© 2006 The Authors
Journal compilation © 2006 Blackwell Publishing Ltd

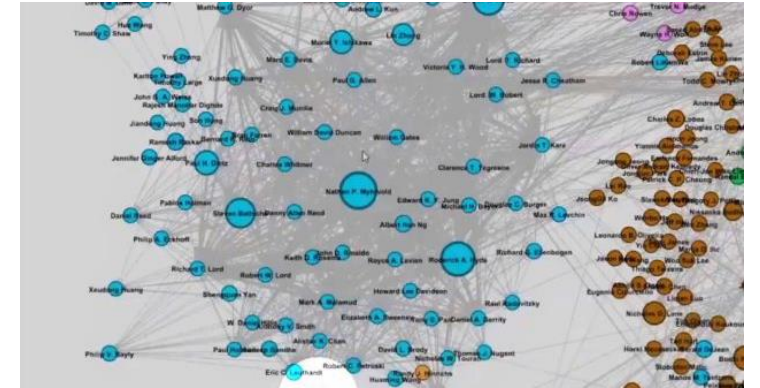
- Digital Repositories
- Digital Library Search Engines
- Intelligent Digital Libraries (IDLs)





- A simple catalog of documents with basic metadata records with or without the full text and a simple boolean query interface
- Can be implemented using a relational database with a single lookup interface
 - University ETD repositories
 - arXiv and sister repos (bioRxiv, medRxiv, PsyArXiv)
 - etc.

Digital Library Search Engines

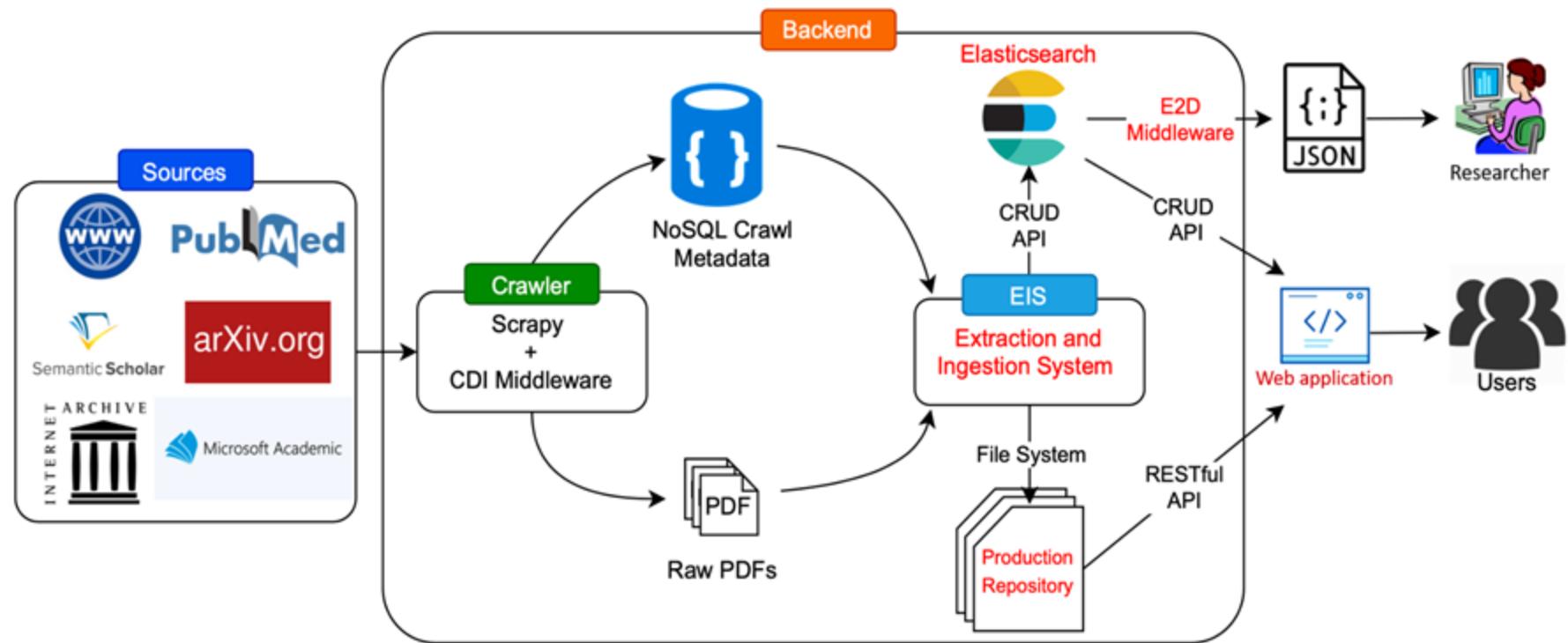


- Documents are organized in a **connected** manner, by inverted index, citation networks, co-author networks, or other structures so that they are more findable, navigable, and usually provide meta-level knowledge.
 - CiteSeerX
 - Google Scholar
 - Semantic Scholar
 - etc.
- Require
 - Backend (crawling, extraction and ingestion, index, repository)
 - Frontend (search, browsing, download)

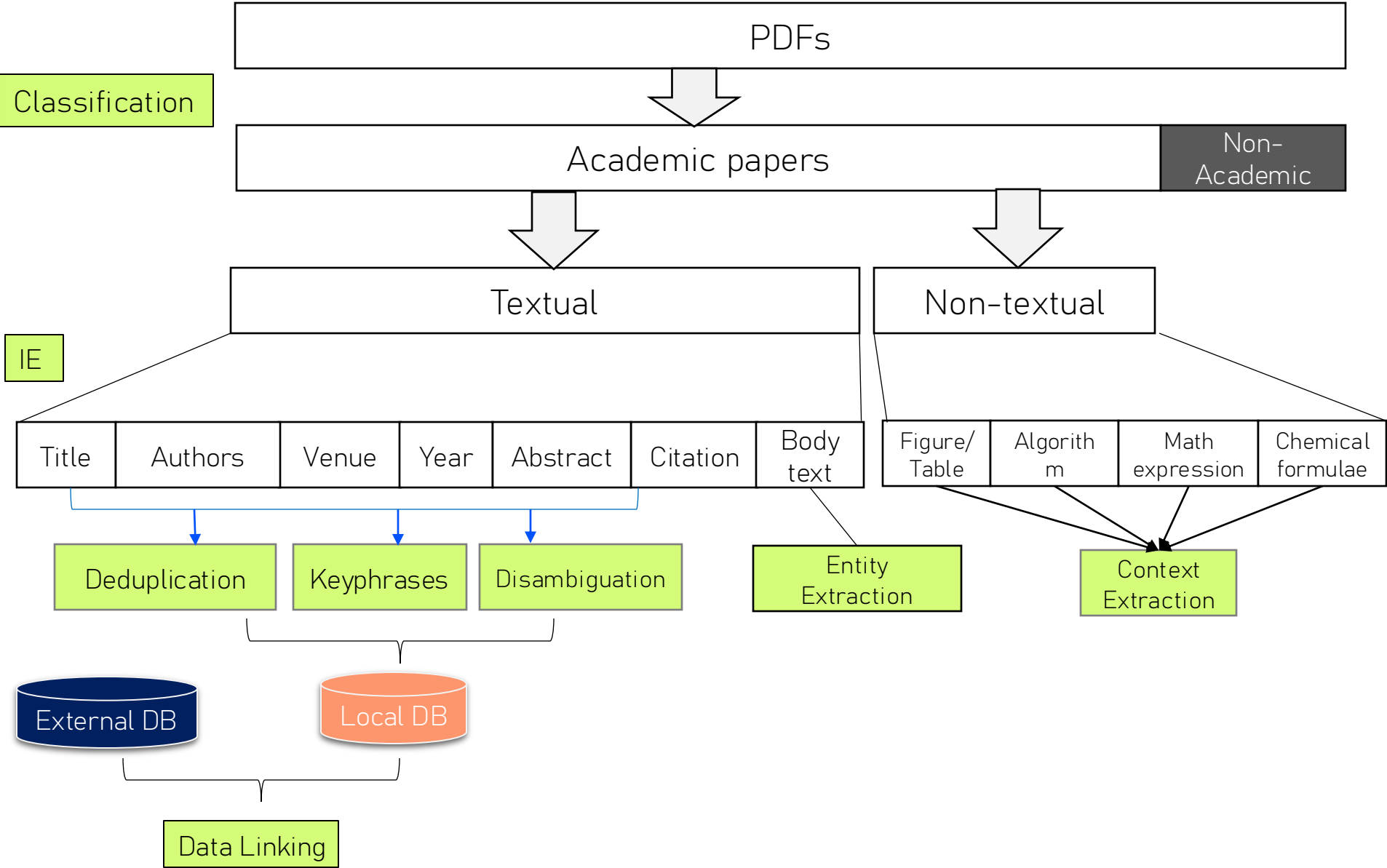


The Architecture of CiteSeerX (2022-present)

- A focused web crawler for open access academic PDFs
- Elasticsearch for a search engine and a data storage
- Integrated and **parallelized** extraction and ingestion
- Redundancy for high **availability** and resilience



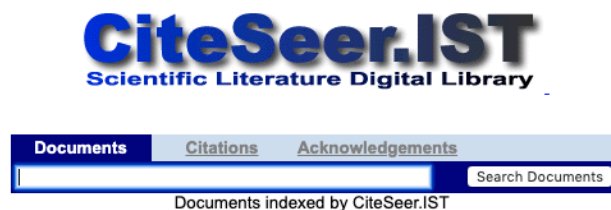
AI in CiteSeerX



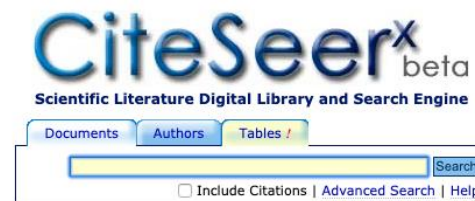
Retrospect and Status of CiteSeerX



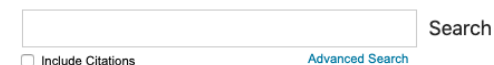
2003



2006



2011



2015



- 15+ million full text English documents and metadata.
- 1 billion hits and 180 million downloads annually.
- Googling "CiteSeerX OR CiteSeer" returns 6 million results (Sep 6, 2025).



Semanti Scholar

- Basic functionalities
 - Search
 - Browse
 - Download

Semantic **Scholar**

- AI-powered functionalities are built based on documents in the repository
 - TLDR summarization
 - Citation intent and influence classifications
 - Field of study classification
 - Paper recommendation
 - Metrics (most influential citations, etc.)

Recently Emerged Digital Library Search Engines

- AllSci: hypothesis-centric, AI-powered
 - Atomized more than 12 million scientific **hypotheses** into a knowledge graph
 - Using AI-guided tools to help researchers formulate better hypotheses
- Scite: using citation context for QA and Table search
- Consensus: search engine + QA
- ResearchRabbit
 - Allowing users to visualize networks

The AllSci logo features the word "allsci" in a bold, black, sans-serif font. A small orange flame icon is positioned above the letter "i".The Scite_ logo consists of the word "scite_" in white, lowercase, sans-serif font, centered within a solid blue rectangular background.The Consensus logo includes a circular icon on the left, composed of two interlocking shapes in blue and green. To the right of the icon, the word "Consensus" is displayed in a bold, black, sans-serif font.

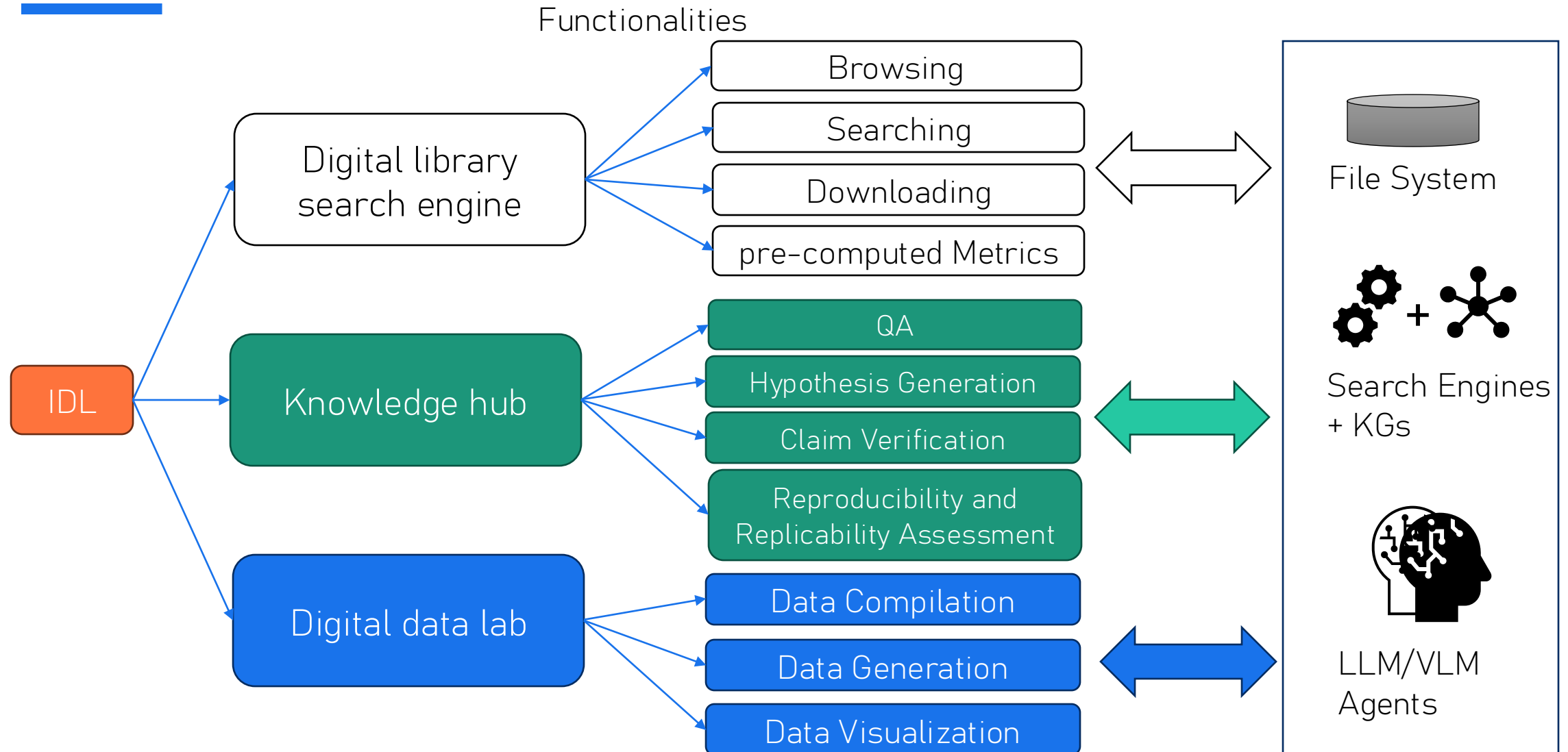
<https://guides.pnw.edu/AISearch>



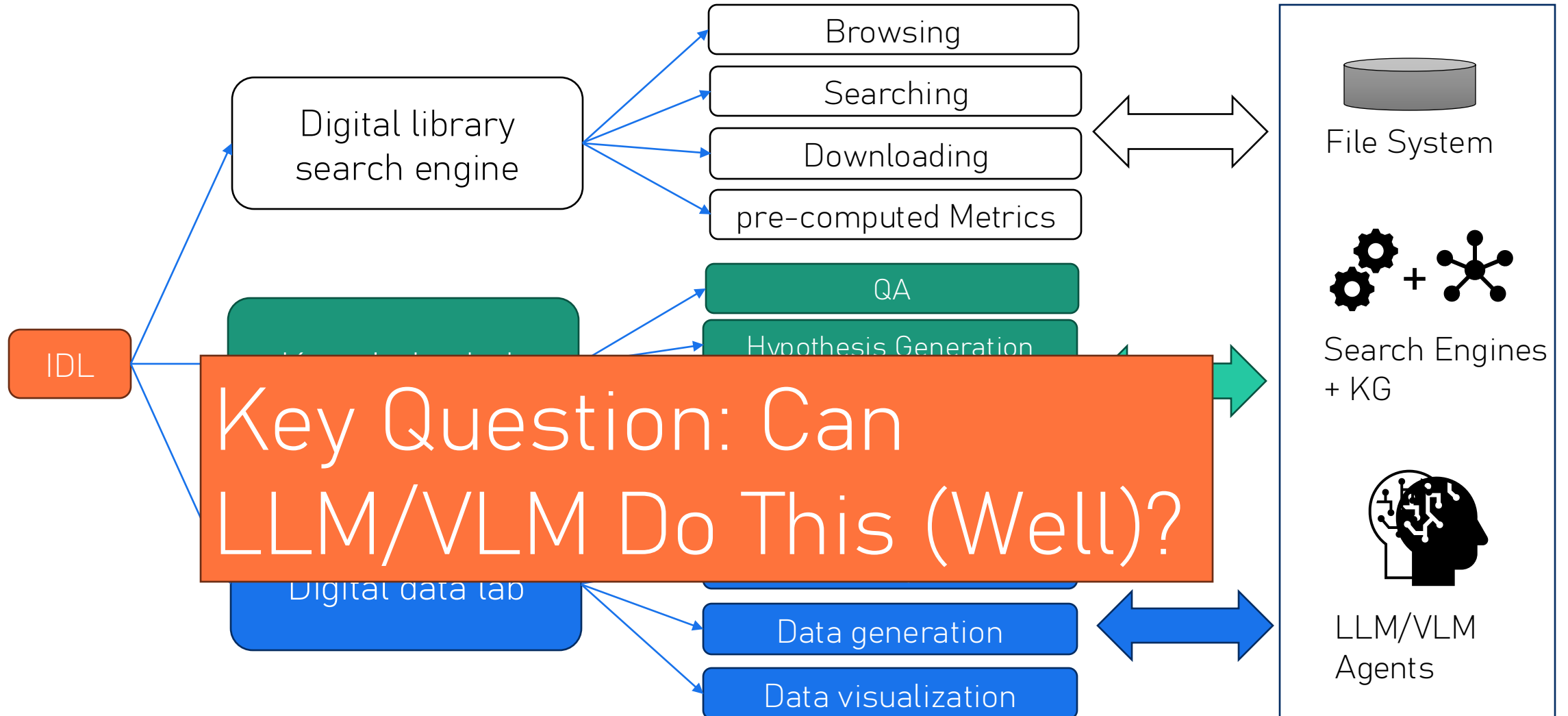
What's Next? Intelligent Digital Libraries (IDLs)

- Providing digital services that require semantic, multigranularity reasonings about a variety of documents and their derivatives.
- Implemented by
 - Backend (document acquisition, IE, index, repository (document, data, software, etc.), LLM/VLM agents, KGs, internal and external connections)
 - Frontend (more interfaces will be created to support user-user, user-document, user-data, user-knowledge **interactions**)

A Vision of the Intelligent Digital Libraries (IDLs)



A Vision of the Intelligent Digital Libraries (IDLs)



IDL as a Knowledge Hub

IDL as a Digital Lab

Recent Research on LLM/VLM-based Functionalities

QA – Offloading Reading to Bots

- Single document QA
 - Answer questions after reading a single scholarly document provided by a user
- Multi-document QA
 - Answer questions after reading multiple documents provided by a user
- Users
 - Any one to need to read papers (students, professors, researchers, reviewers, etc.)



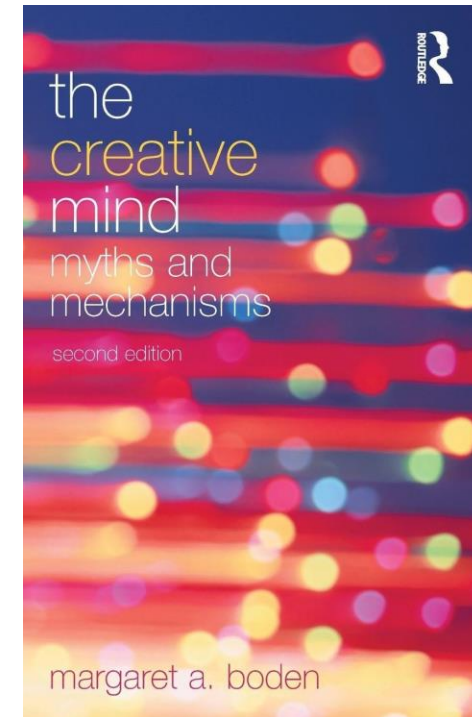
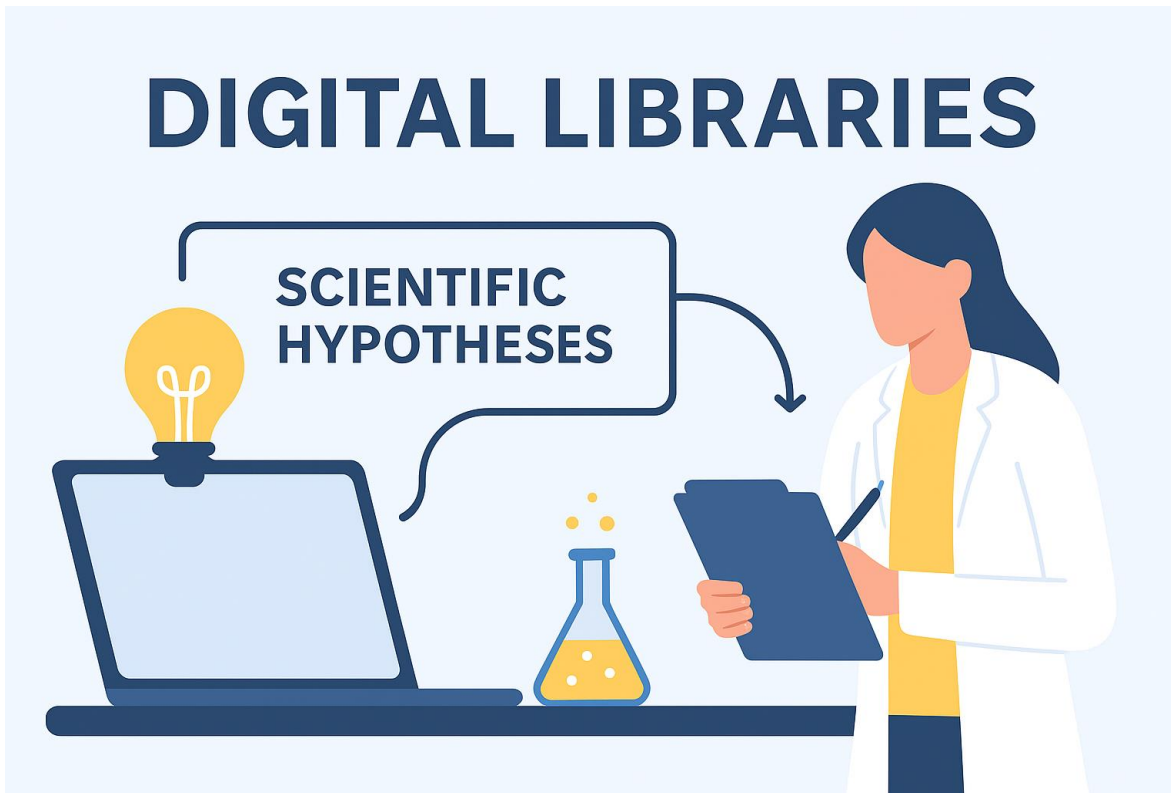


Limitations of Existing Tools

- Limitation: users must know which documents to use to find the answers
- Example: what evidence has been reported to refute that Covid was produced in the lab?
- Vision: automatically **reason** on hundreds of millions of documents to answer the questions (as opposed to simply search for **relevant** documents)

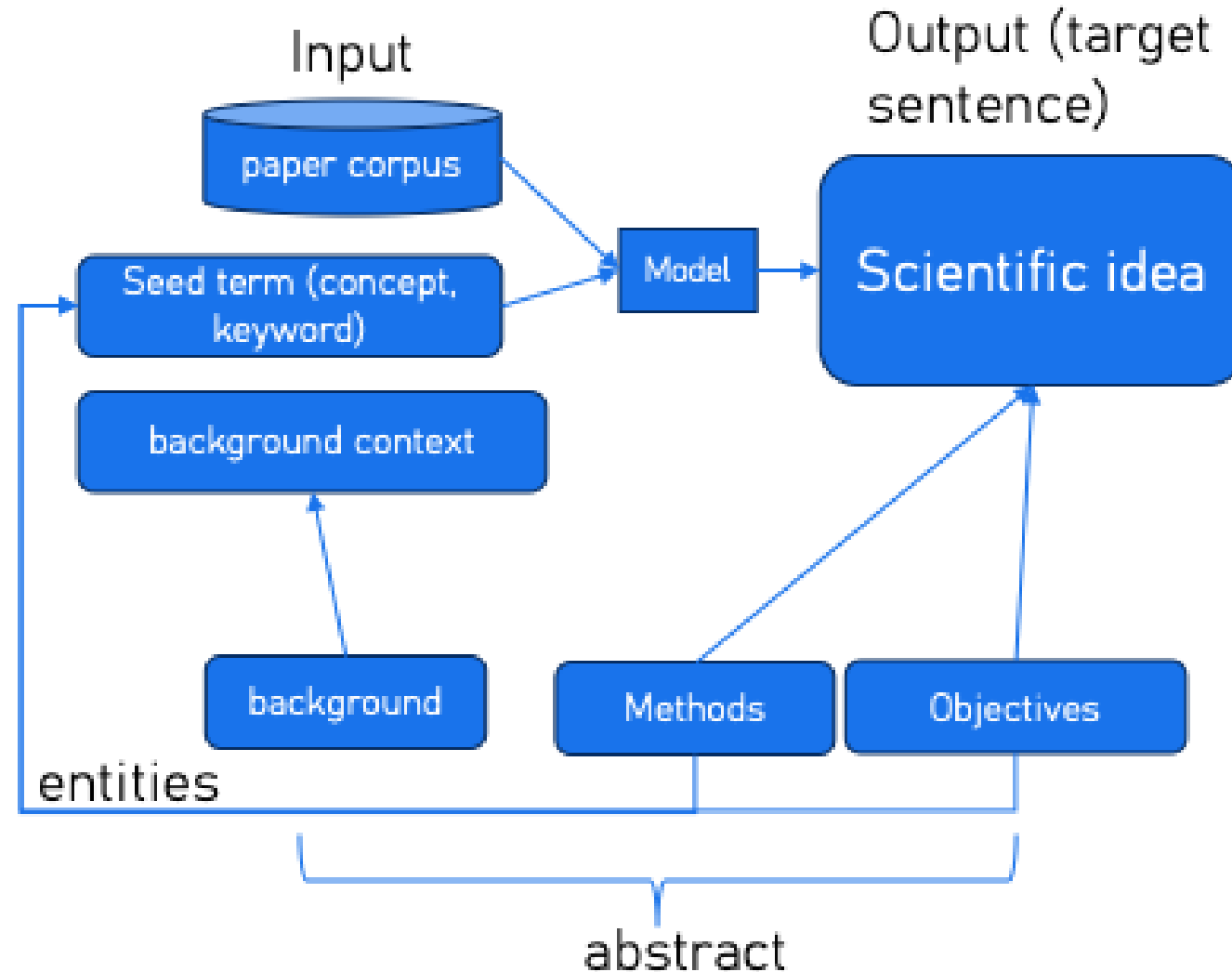
Hypothesis Generation

- Can IDLs help scientist propose novel and feasible scientific hypotheses and then plan experiments to verify the them?

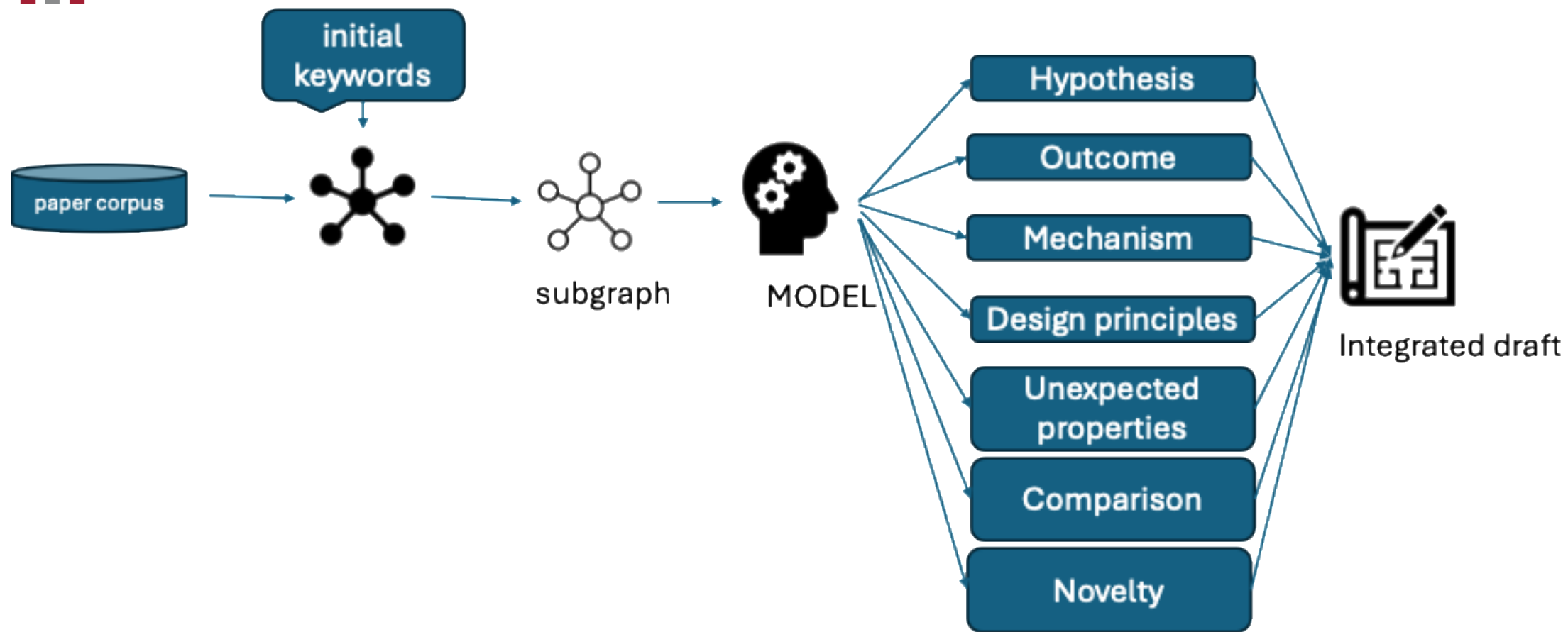


Boden, M. A. (2004).
The Creative Mind:
Myths and
Mechanisms

SciMon (Wang et al. ACL 2024)



SciAgents (Ghafarollahi et al., Advanced Materials, 2025)



What Are Research Hypotheses?

Jian Wu^{1,*}, Sarah Rajtmajer^{2,*}

¹Old Dominion University

²The Pennsylvania State University



2025

Hypothesis Generation

Title	Input	Output	Reference
SciMon	Background + Keywords	Ideas	Wang et al. (ACL'24)
SciAgents	Keywords	Proposal	Ghafarollahi et al. (Advanced Materials, 2025)
Hypothesis generation with large language models	Data	Hypothesis	Zhou et al. (NLP4Science'24)
AI-CoScientist	Research Goal	Proposal	Gottweis et al., (arXiv:2502.18864)

Evaluation is a major challenge!

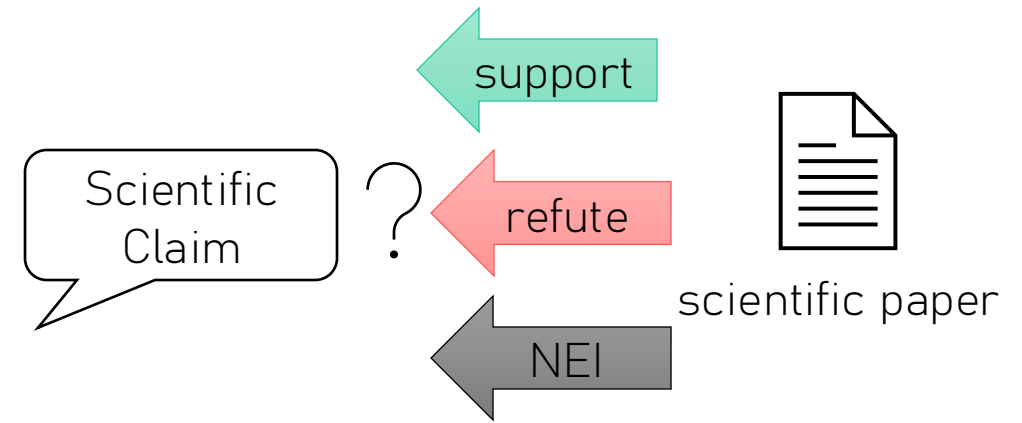
- Models have different input and output.
- A lack of benchmark and expert evaluation standard.

Our Proposed Work

- We propose generating highly novel and interdisciplinary citation-enriched hypothesis proposals through iterative interactions between expert LLMs.
- One challenge is evaluation. How to **automatically evaluate** the quality of the generated hypotheses?
 - Novelty
 - Feasibility



Scientific Claim Verification



- Problem definition: Given a claim and a scientific paper, can AI tell us if the paper supports or refute claim (or does not provide enough information (NEI))?
- Example A (claim in a scientific paper): Is there an association between social media use and bad mental health outcomes?
- Example B (claim in news or social media): Use of hand sanitizer can seriously mess with breath alcohol test results.
- Use cases
 - Scientific review
 - Misinformation and disinformation

Discern Claims (hypotheses) in Scientific Papers

Can Large Language Models Discern Evidence for Scientific Hypotheses? Case Studies in the Social Sciences

Sai Koneru¹, Jian Wu², Sarah Rajtmajer¹

¹ Pennsylvania State University, State College, PA

² Old Dominion University, Norfolk, VA

{sdk96, smr48}@psu.edu, j1wu@odu.edu

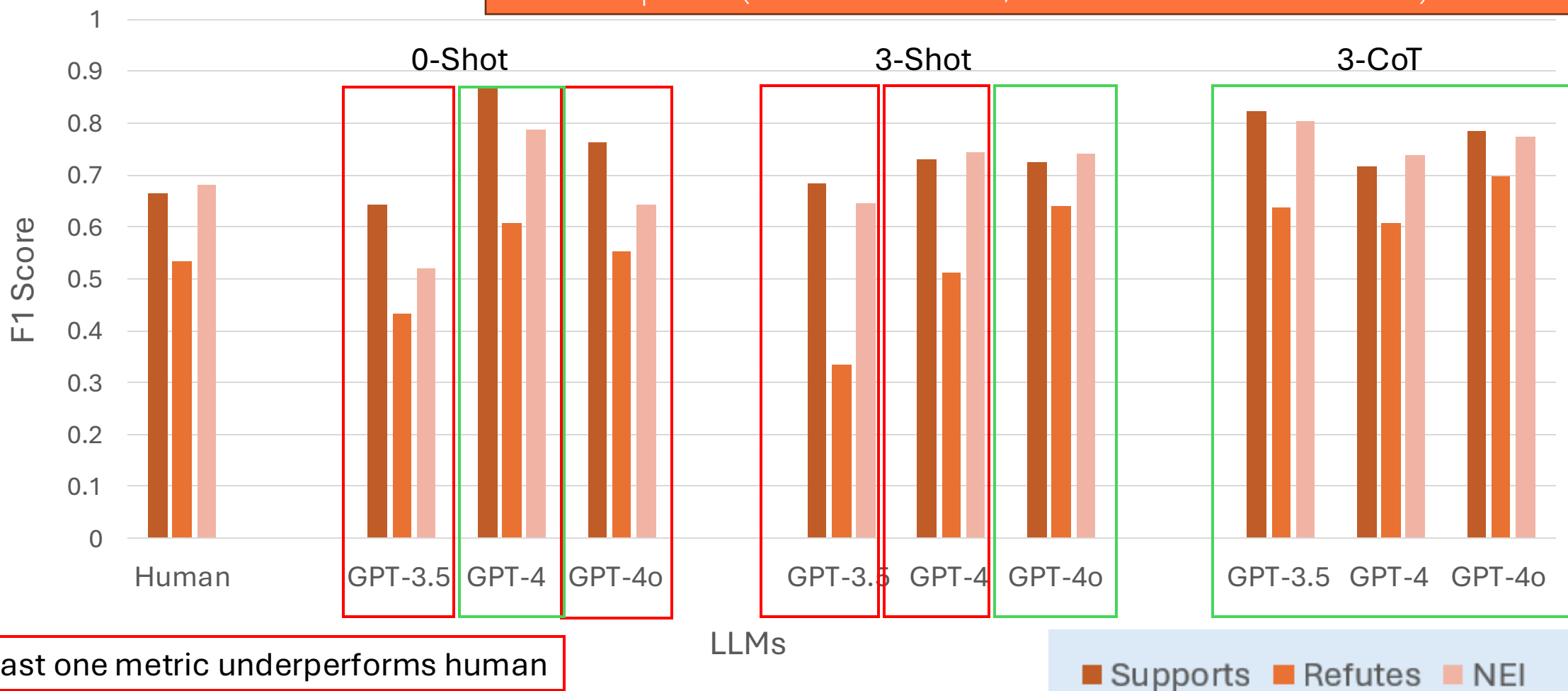


AI underperforms domain experts on discerning claims (hypotheses) in scientific papers.

Best model of each type	Accuracy	Macro F1
embedding + supervised classification	70.31%	0.615
Transfer learning	67.97%	0.523
GPT3.5 few-shot	66.57%	0.576
PaLM2	62.87%	0.536

Discern Claims in News and Social Media posts

- AI beats college students on discerning claims from news and social media posts. (Evans et al. 2024, Dzhaman et al. 2025 NCUR)



Automatic Reproducibility and Replicability Assessment

- Reproducibility: same data, same method
- Replicability: different data, same method
- Reproducibility and replicability crisis in
 - Social and Behavioral Science (SBS) (Camerer 2016 Nature; Camerer 2018 Nature)
 - Computer Science (Moraila et al. 2014 PloS; Collberg et al. 2016)
 - Artificial Intelligence (Raff et al. 2019 NeurIPS; Gundersen et al. 2018 AAAI; Haibe-Kains et al. 2020 Nature; Ajayi et al. 2023 ICDAR)
 - Biomedical Science (Gentleman et a. 2005)

The Challenge of Reproducibility and Replicability

- Manually reproducing reported results is time-consuming and not scalable
- Average time to reproduce the main results in one paper
 - Table Structure Recognition (an AI task): **8 hours** (using code and data provided by the original authors; Ajayi et al. 2023 ICDAR)
 - General AI tasks: **53.5 days** (using re-implemented codes and data provided by the original authors; Raff 2023 AACL)
 - Social and Behavioral Science: months – **up to 1 year** (using the same methods and new data collected from new user studies)

Can IDL tell me the reproducibility and replicability of a paper?

Citation Context Sentiments vs. Reproducibility Scores

- Correlation analysis between sentiment and reproducibility scores for ML papers
- Positive sentiment correlates with higher reproducibility scores
- Negative sentiment correlates with lower reproducibility scores

SHORT: Can citations tell us about a paper's reproducibility? A case study of machine learning papers

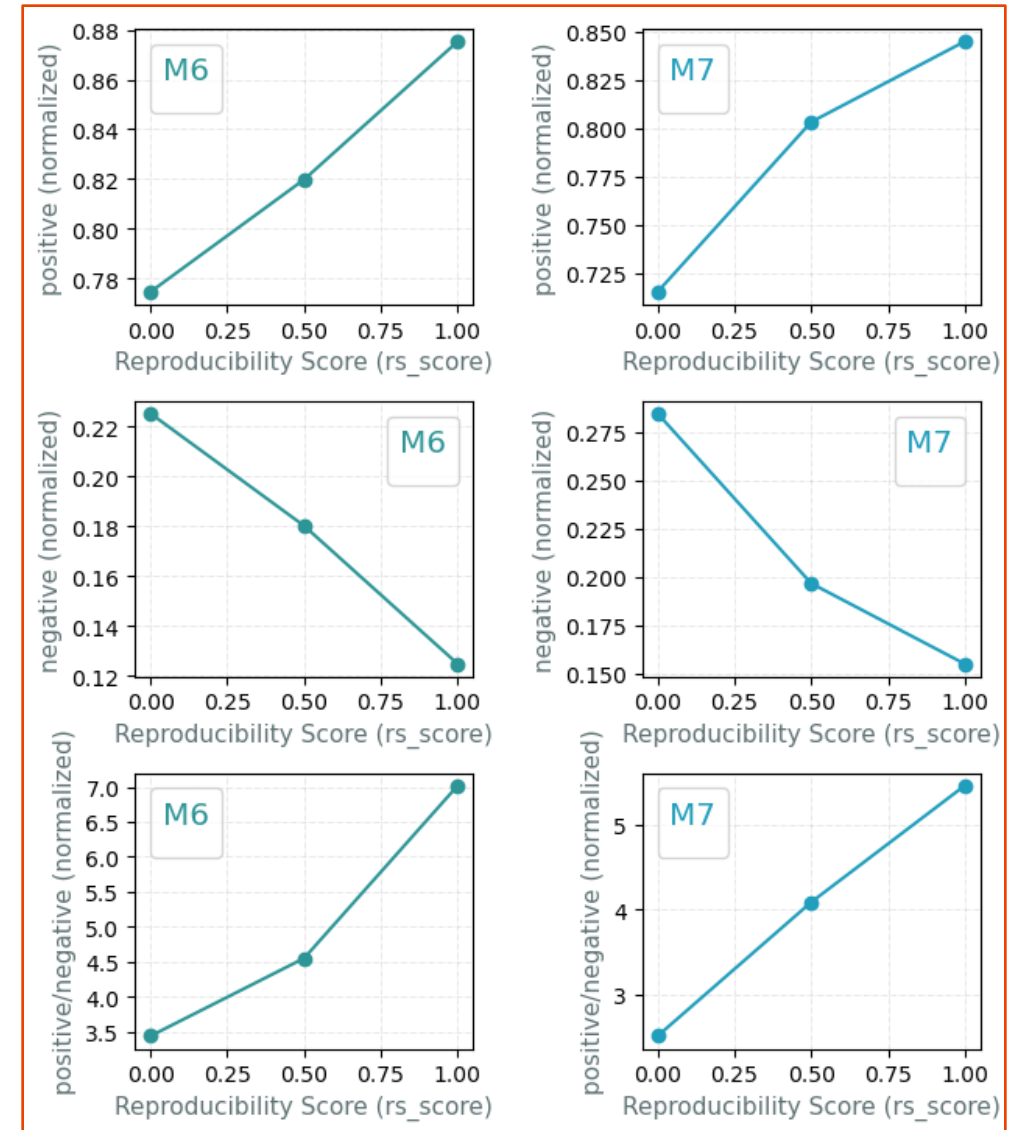
Rochana R. Obadage
Old Dominion University
Norfolk, VA, USA
oruma001@odu.edu

Sarah M. Rajtmajer
IST, Pennsylvania State University
University Park, PA, USA
smr48@psu.edu

Jian Wu
Old Dominion University
Norfolk, VA, USA
j1wu@odu.edu



ACM REP '24





Automatic Reproducibility and Replicability Assessments

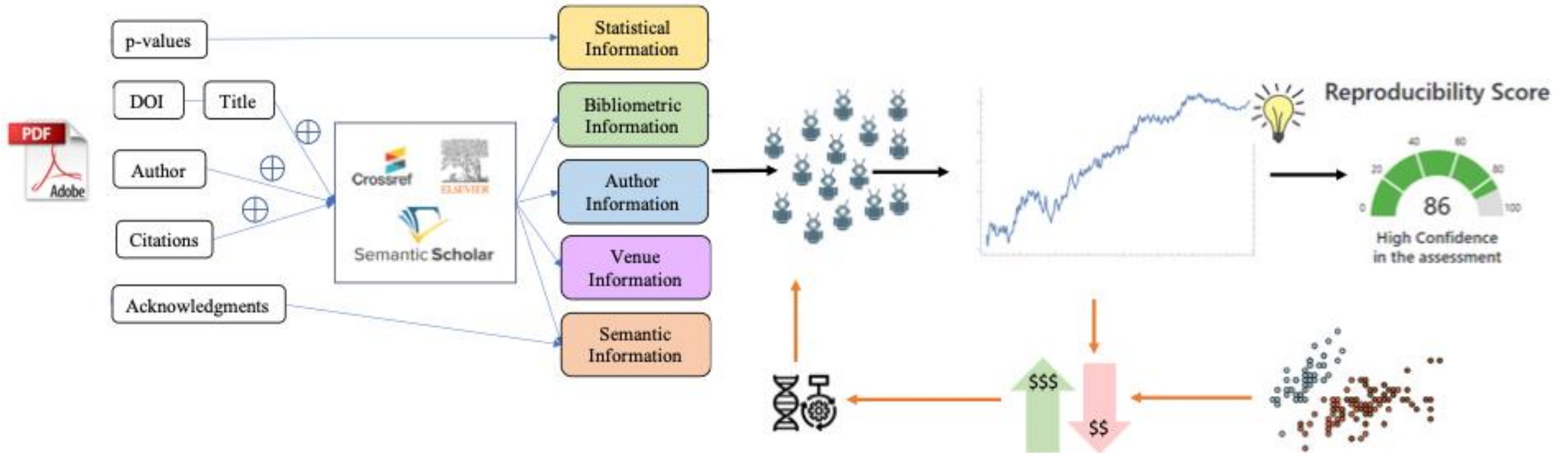
- Machine Learning (Wu et al. 2021, arxiv)
- Deep Learning (Yang et al. 2020 PNAS; Wu et al. 2023 PNAS)
- Prediction Market (Viganola et al. 2021 R. Soc. Open Sci.)
- Synthetic Prediction Market (Rajtmajer et al. 2022 AAAI; Chakravorti et al. 2023 HHAI)
- Question: can we use LLM agents to automatically reproduce/replicate the claims?

Machine Learning vs. Neural Models vs. Prediction Market

Model Type	Model	F1	Accuracy
ML (Our Work)	Random Forest	0.61	0.65
	XGBoost	0.57	0.62
	Naïve Bayesian	0.52	0.62
	MLP	0.54	0.59
	Logistic Regression	0.51	0.56
	SVM	0.46	0.56
Neural	word2vec on full text (Yang et al. 2020 PNAS)	NA	0.65-0.78
Prediction market	Prediction market (Chakravorti et al. 2023)	NA	0.71
Human	Human survey (Chakravorti et al. 2023)	NA	0.58

** The test samples of ML, Neural, and Prediction Market overlap but are not exactly the same but are all in social and behavioral sciences.

A **Synthetic Prediction Market** for Estimating Confidence in Published Work (Rajtmajer et al. 2022 AAAI)



Synthetic prediction markets— *Prediction markets populated by **artificial** agents (trader-bots), trained and updated within human-expert prediction markets, but deployable "offline".*

- Trader-bots will represent atomic (human-interpretable) properties of relevant signals, including full text of scientific papers, metadata for specific papers, and metadata about the community and the field.
- Bots will learn trading patterns from subject matter experts engaged in prediction markets, but unlike their human counterparts, will have comprehensive, unbiased view of the existing literature and metadata.

Synthetic Prediction Market (Our Work)

A Synthetic Prediction Market for Estimating Confidence in Published Work

Sarah Rajtmajer,¹ Christopher Griffin,¹ Jian Wu,² Robert Fraleigh,¹ Laxmaan Balaji,¹ Anna Squicciarini,¹ Anthony Kwasnica,¹ David Pennock,³ Michael McLaughlin,¹ Timothy Fritton,¹ Nishanth Nakshatri,¹ Arjun Menon,¹ Sai Ajay Modukuri,¹ Rajal Nivargi,¹ Xin Wei,² C. Lee Giles¹

¹The Pennsylvania State University

²Old Dominion University

³Rutgers University

{smr48,cxg286,rdf5090,lpb5347,acs20,amk17,mvm7085,tjf115,nzn5185,amm8987,svm6277,rfn5089,clg20}@psu.edu

{j1wu,xwei001}@odu.edu

david.pennock@rutgers.edu



Caveats:

- Feature extraction is noisy
- Bots may not always converge (hybrid market)

Can LLM directly reproduce/replicate the work in the paper?

Results on scored papers. Our system provides a confidence score for 68 of 192 (35%) of the papers in our set. On the set of scored papers, accuracy is 0.894, precision is 0.917, recall is 0.903, and **F1 is 0.903** (macro averages). A sizeable un-scored subset of data (65%) is the trade-off for high accuracy on the scored subset of the data. A test point is un-scored when the system has determined it has insufficient information to evaluate it.

LLM-based Reproducibility and Replicability Assessment

Features	Core-Bench  PRINCETON UNIVERSITY	PaperBench  OpenAI
References	Siegel et al. (2024 arxiv2409.11363)	Starace et al. (2025 arXiv2504.01848)
Nature of Work	Reproducibility	Replicability (same data, same method, new implementation)
Domain	CS, Bio, and Social	ML
#Papers	90	20
#Tasks	270	8,316
Agent Scaffolding	Adapted AutoGPT, orchestration: Python subprocess	Inspect AI basic agent, orchestration: nanoeval
Input	Code, Original Data, Paper	Paper only
Human comparison	No	Yes (Top PhD)
Conclusion	GPT-4o-based CORE-agent achieves the best reproducibility rate	Claude 3.5 Sonnet (New) with open-source scaffolding achieves the highest score

Our Ongoing Research

- Goal: building a new benchmark for LLM agents to replicate claims in Social and Behavioral Science papers
- Data : 200+ papers, pre-registrations, and human replication study reports from the SCORE project (Nosek et al. 2021 ARP)
- Replicability (new data), multi-difficulty level, multi-stage, no-human involved



PennState®



OLD DOMINION
UNIVERSITY®

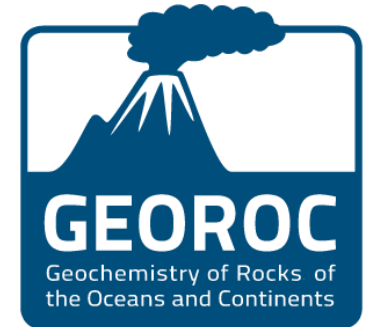
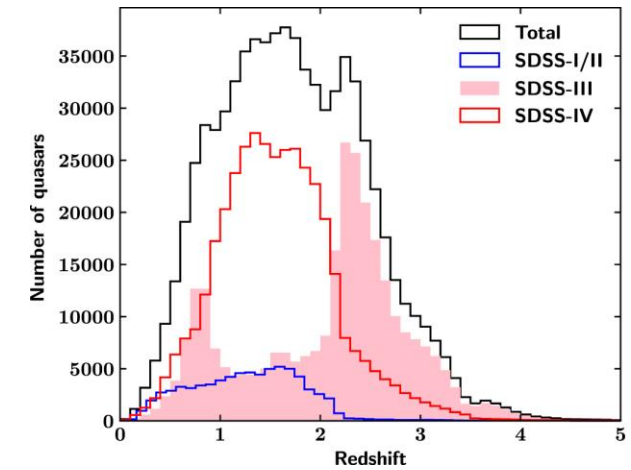
IDL as a Knowledge Hub

IDL as a Digital Lab

Recent Research on LLM/VLM-based Functionalities

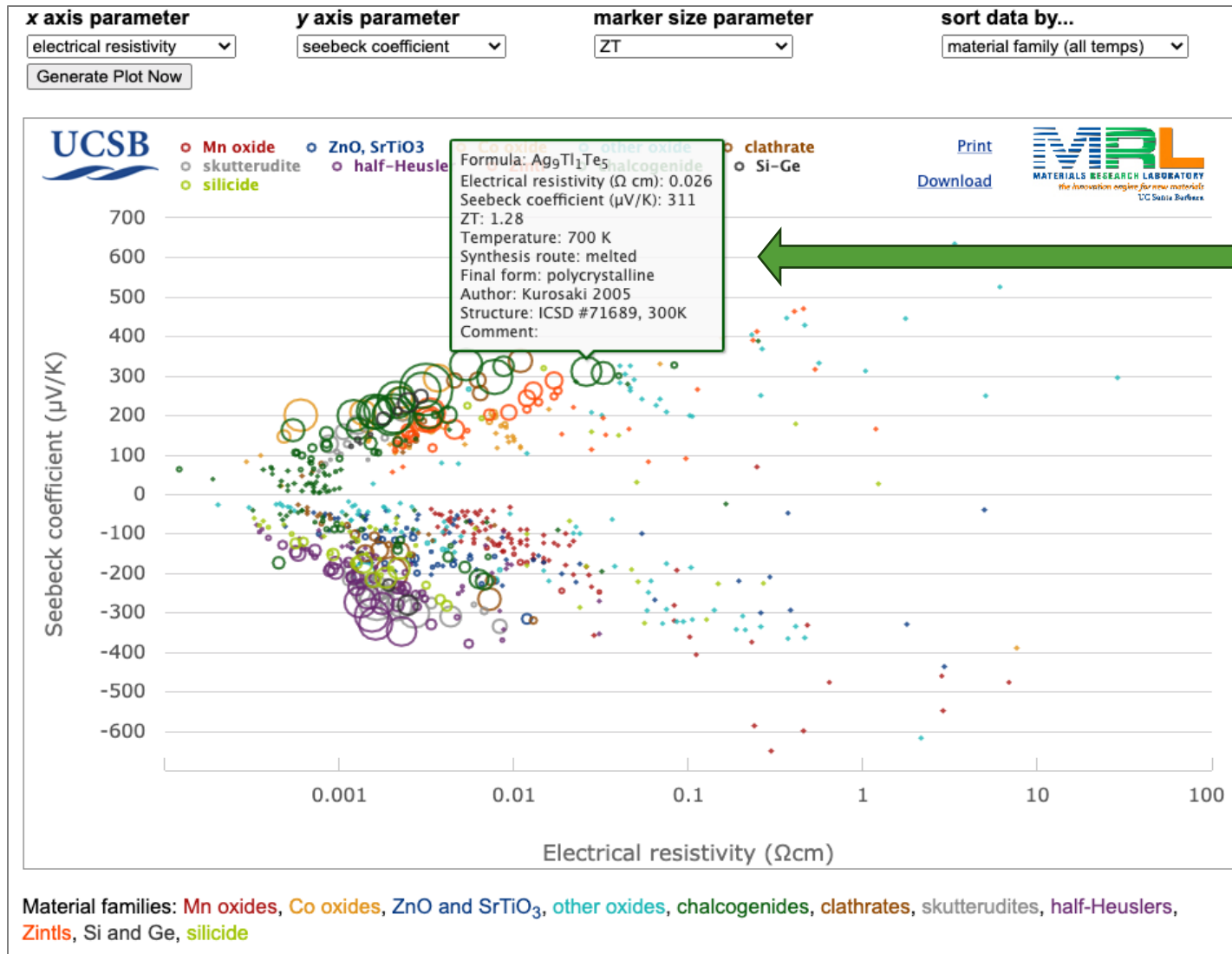
Data Compilation

- Why is it a problem for domain scientists?
 - Tons of data are published in PDFs
 - Data can be published in multiple modalities (text, figures, tables)
 - Compiled data published in papers can be very useful for analysis
 - Astronomy: Catalog. SDSS quasar catalog, Véron-Cetty & Véron Quasar Catalog, etc.
 - Geoscience: Geochemical database of elemental abundances: GEOROC, and EarthChem
 - Materials Science: the CALPHAD database of thermodynamic parameters (e.g., Gibbs free energy functions, interaction parameters) for elements, compounds, and phases.
- Question: can we automatically and faithfully compile data from scholarly papers into a pre-defined table schema for downstream analysis?



An Example in Materials Science

<http://www.mrl.ucsb.edu:8080/datamine/thermoelectric.jsp>



source paper

ysics Letters

Ag₉TlTe₅: A high-performance thermoelectric bulk material with extremely low thermal conductivity

Cite as: Appl. Phys. Lett. **87**, 061919 (2005); <https://doi.org/10.1063/1.2009828>
Submitted: 21 April 2005 . Accepted: 01 July 2005 . Published Online: 04 August 2005

Ken Kurosaki, Atsuko Kosuga, Hiroaki Muta, Masayoshi Uno, and Shinsuke Yamanaka


<https://aip.scitation.org/doi/10.1063/1.2009828>

Can LLM/VLM Do this ?

Department of Chemistry and the Radiation Center, Oregon State University,
Corvallis, Oregon 97331



	10085, 796	10085, 803	10085, 814	10085, 820	10085, 824	10085, 1002	10085, 1012	10085, 1015	10085, 1020	10085, 1032	10085, 1037	10085, 1042	10085, 1046	10085, 1050	10085, 1054	10085, 1058	10085, 1062	10085, 1066	10085, 1070	10085, 1074	10085, 1078	10085, 1082	10085, 1086	10085, 1090	10085, 1094	10085, 1098	10085, 1102	10085, 1106	10085, 1110	10085, 1114	10085, 1118	10085, 1122	10085, 1126	10085, 1130	10085, 1134	10085, 1138	10085, 1142	10085, 1146	10085, 1150	10085, 1154	10085, 1158	10085, 1162	10085, 1166	10085, 1170	10085, 1174	10085, 1178	10085, 1182	10085, 1186	10085, 1190	10085, 1194	10085, 1198	10085, 1202	10085, 1206	10085, 1210	10085, 1214	10085, 1218	10085, 1222	10085, 1226	10085, 1230	10085, 1234	10085, 1238	10085, 1242	10085, 1246	10085, 1250	10085, 1254	10085, 1258	10085, 1262	10085, 1266	10085, 1270	10085, 1274	10085, 1278	10085, 1282	10085, 1286	10085, 1290	10085, 1294	10085, 1298	10085, 1302	10085, 1306	10085, 1310	10085, 1314	10085, 1318	10085, 1322	10085, 1326	10085, 1330	10085, 1334	10085, 1338	10085, 1342	10085, 1346	10085, 1350	10085, 1354	10085, 1358	10085, 1362	10085, 1366	10085, 1370	10085, 1374	10085, 1378	10085, 1382	10085, 1386	10085, 1390	10085, 1394	10085, 1398	10085, 1402	10085, 1406	10085, 1410	10085, 1414	10085, 1418	10085, 1422	10085, 1426	10085, 1430	10085, 1434	10085, 1438	10085, 1442	10085, 1446	10085, 1450	10085, 1454	10085, 1458	10085, 1462	10085, 1466	10085, 1470	10085, 1474	10085, 1478	10085, 1482	10085, 1486	10085, 1490	10085, 1494	10085, 1498	10085, 1502	10085, 1506	10085, 1510	10085, 1514	10085, 1518	10085, 1522	10085, 1526	10085, 1530	10085, 1534	10085, 1538	10085, 1542	10085, 1546	10085, 1550	10085, 1554	10085, 1558	10085, 1562	10085, 1566	10085, 1570	10085, 1574	10085, 1578	10085, 1582	10085, 1586	10085, 1590	10085, 1594	10085, 1598	10085, 1602	10085, 1606	10085, 1610	10085, 1614	10085, 1618	10085, 1622	10085, 1626	10085, 1630	10085, 1634	10085, 1638	10085, 1642	10085, 1646	10085, 1650	10085, 1654	10085, 1658	10085, 1662	10085, 1666	10085, 1670	10085, 1674	10085, 1678	10085, 1682	10085, 1686	10085, 1690	10085, 1694	10085, 1698	10085, 1702	10085, 1706	10085, 1710	10085, 1714	10085, 1718	10085, 1722	10085, 1726	10085, 1730	10085, 1734	10085, 1738	10085, 1742	10085, 1746	10085, 1750	10085, 1754	10085, 1758	10085, 1762	10085, 1766	10085, 1770	10085, 1774	10085, 1778	10085, 1782	10085, 1786	10085, 1790	10085, 1794	10085, 1798	10085, 1802	10085, 1806	10085, 1810	10085, 1814	10085, 1818	10085, 1822	10085, 1826	10085, 1830	10085, 1834	10085, 1838	10085, 1842	10085, 1846	10085, 1850	10085, 1854	10085, 1858	10085, 1862	10085, 1866	10085, 1870	10085, 1874	10085, 1878	10085, 1882	10085, 1886	10085, 1890	10085, 1894	10085, 1898	10085, 1902	10085, 1906	10085, 1910	10085, 1914	10085, 1918	10085, 1922	10085, 1926	10085, 1930	10085, 1934	10085, 1938	10085, 1942	10085, 1946	10085, 1950
--	---------------	---------------	---------------	---------------	---------------	----------------	----------------	----------------	----------------	----------------	----------------	----------------	----------------	----------------	----------------	----------------	----------------	----------------	----------------	----------------	----------------	----------------	----------------	----------------	----------------	----------------	----------------	----------------	----------------	----------------	----------------	----------------	----------------	----------------	----------------	----------------	----------------	----------------	----------------	----------------	----------------	----------------	----------------	----------------	----------------	----------------	----------------	----------------	----------------	----------------	----------------	----------------	----------------	----------------	----------------	----------------	----------------	----------------	----------------	----------------	----------------	----------------	----------------	----------------	----------------	----------------	----------------	----------------	----------------	----------------	----------------	----------------	----------------	----------------	----------------	----------------	----------------	----------------	----------------	----------------	----------------	----------------	----------------	----------------	----------------	----------------	----------------	----------------	----------------	----------------	----------------	----------------	----------------	----------------	----------------	----------------	----------------	----------------	----------------	----------------	----------------	----------------	----------------	----------------	----------------	----------------	----------------	----------------	----------------	----------------	----------------	----------------	----------------	----------------	----------------	----------------	----------------	----------------	----------------	----------------	----------------	----------------	----------------	----------------	----------------	----------------	----------------	----------------	----------------	----------------	----------------	----------------	----------------	----------------	----------------	----------------	----------------	----------------	----------------	----------------	----------------	----------------	----------------	----------------	----------------	----------------	----------------	----------------	----------------	----------------	----------------	----------------	----------------	----------------	----------------	----------------	----------------	----------------	----------------	----------------	----------------	----------------	----------------	----------------	----------------	----------------	----------------	----------------	----------------	----------------	----------------	----------------	----------------	----------------	----------------	----------------	----------------	----------------	----------------	----------------	----------------	----------------	----------------	----------------	----------------	----------------	----------------	----------------	----------------	----------------	----------------	----------------	----------------	----------------	----------------	----------------	----------------	----------------	----------------	----------------	----------------	----------------	----------------	----------------	----------------	----------------	----------------	----------------	----------------	----------------	----------------	----------------	----------------	----------------	----------------	----------------	----------------	----------------	----------------	----------------	----------------	----------------	----------------	----------------	----------------	----------------	----------------	----------------	----------------	----------------	----------------	----------------	----------------	----------------	----------------	----------------	----------------	----------------	----------------



Automatically Compile the Apollo Basalt Database

Lunar sample age and errors in billion years

Elemental abundance pattern

Age (by)	Age error (by)	Dating method	Cone age ref.	SiO2	TiO2	Al2O3	Cr2O3	FeO	MnO	MgO	CaO	Na2O	K2O	P2O5	BaO	S	SO3	NiO	ZrO2
NA	NA	NA	NA	NA	9.7	9.5	0.266	19.5	0.248	7	11.1	0.396	0.073	NA	NA	NA	NA	NA	NA
NA	NA	NA	NA	NA	11.1	7.3	0.342	20.5	0.218	9	10.5	0.462	0.28	NA	NA	NA	NA	NA	NA
NA	NA	NA	NA	NA	11.5	7.7	0.348	20.1	0.220	8	10.9	0.474	0.28	NA	NA	NA	NA	NA	NA
NA	NA	NA	NA	NA	10.2	8.6	0.289	20.2	0.223	8	10.8	0.527	0.33	NA	NA	NA	NA	NA	NA
3.83	0.016	Ar-Ar	14	41.2	8.6	11.8	0.3	19	0.249	8	11.5	0.37	0.09	0.28	NA	NA	NA	NA	NA
3.896	0.019	Ar-Ar	14	NA	NA	NA	0.24	18.8	0.275	NA	10.7	0.38	0.094	NA	NA	NA	NA	NA	NA
NA	NA	NA	NA	NA	NA	NA	NA	21.05	NA	NA	12.74	0.393	NA	NA	NA	NA	NA	NA	NA
NA	NA	NA	NA	NA	NA	NA	NA	16.29	NA	NA	16.58	0.54	NA	NA	NA	NA	NA	NA	NA
NA	NA	NA	NA	NA	NA	NA	NA	21.49	NA	NA	NA	0.29	NA	NA	NA	NA	NA	NA	NA
NA	NA	NA	NA	NA	NA	NA	NA	24.66	NA	NA	NA	0.40	NA	NA	NA	NA	NA	NA	NA
NA	NA	NA	NA	NA	NA	NA	NA	20.04	NA	NA	16.14	0.48	NA	NA	NA	NA	NA	NA	NA

Two Subtasks of Data Compilation (tabular data)

Table Image

Cell Coordinates

XML

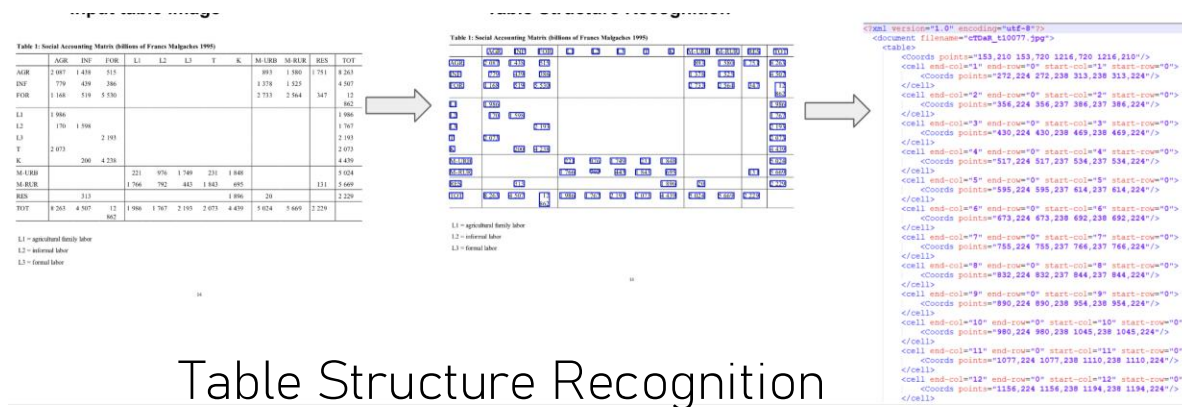


Table Structure Recognition

1. **Data Extraction:** faithfully extract data from tables in PDF documents

- Step 1: Table Structure Recognition
- Step 2: OCR
- VLM-powered OCR or use a VLM to directly extract data
- Improve data correction efficiency using Uncertainty Quantification (UQ)

2. **Data Population:** insert data into the correct position of a table with a pre-defined schema

- Not extensively evaluated to our best knowledge



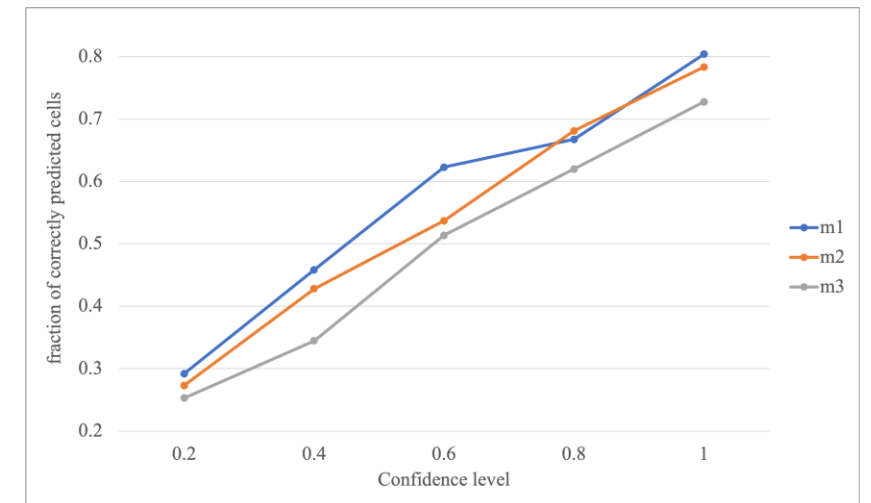
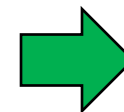
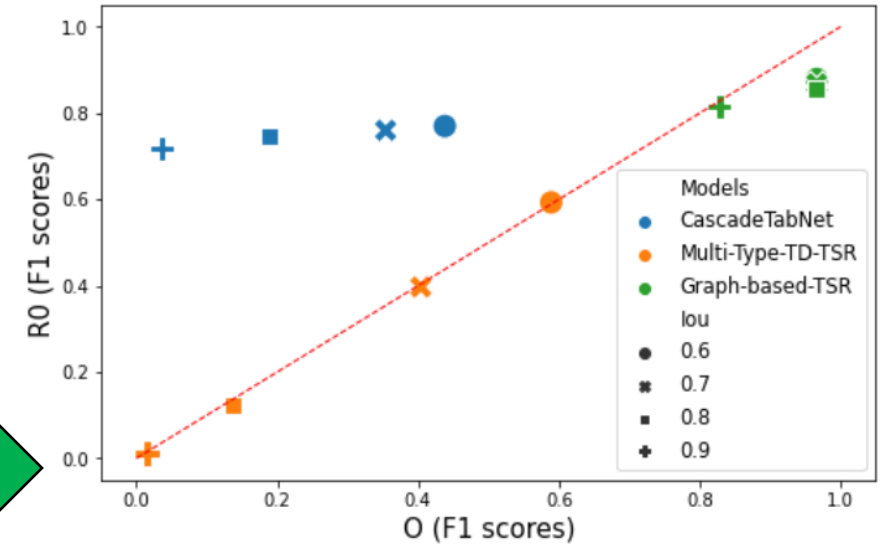
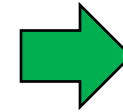


Existing Performance on Table Data Extraction

- CHATEXTRACT (Polak et al. Nature): extracting material data from research papers
 - precision = 90.8%, recall = 87.7%
- CHARTLLAMA (Han et al. arXiv:2311.16483): Chart data extraction
 - precision = 84.92%
- CHATGPT (Brown et al. 2020): extract structured data from clinical notes
 - precision of 77–99%, recall of 76–91%

Our Work on Table Data Extraction

1. Neural models can achieve 75%–85% F1-scores on table structure recognition. However, reproducibility is a big concern for many results reported (Ajayi et al. 2023 ICDAR).
2. The errors in extracted data are correlated with the uncertainty (or confidence levels) of the extraction results (Ajayi et al. 2024 IRI).
3. Conformal prediction is an effective UQ method for table data extraction and can reduce the manual effort of data correction by 50% (Ajayi et al. 2025 ICDAR).

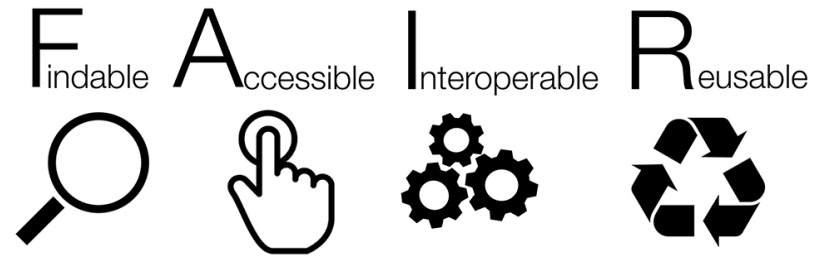


Our Proposed Work

- Automatic data compilation:
 - Use VLMs to extracting tabular data and attributes from PDFs
 - Use LLM/VLM to populate data into a table or a KG (given the schema or ontology)
- Requirement: Ensure fidelity and trustworthy without fine-tuning (very limited training data)



Data Generation



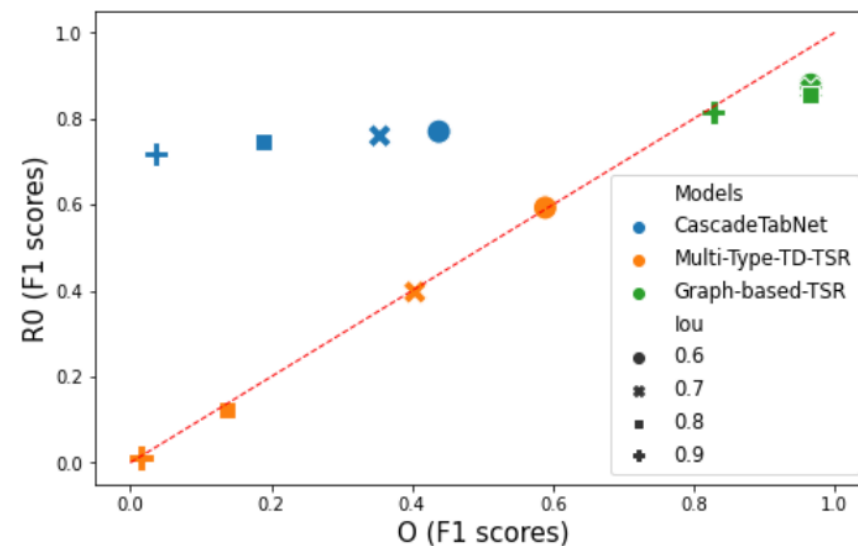
- Motivation: as a solution of data scarcity problem
- Research question: can you use LLMs to generate close-to or better than human generated data?
- Example: web data has depleted for LLM training, can we generate new data?
- Problem:
 - Data security and data quality (beyond the FAIR principle)
 - Integrate data generation modules to data extraction and compilation

Data Visualization

- Motivation: Same data can be illustrated in different ways, making conclusions more evident.

TSR Model	Data	IoU	$F1(O)$	$F1(R_0)$	Δ_0
CascadeTabNet	ICDAR 2019	0.6	0.438	0.770	0.332
CascadeTabNet	ICDAR 2019	0.7	0.354	0.760	0.406
CascadeTabNet	ICDAR 2019	0.8	0.190	0.745	0.555
CascadeTabNet	ICDAR 2019	0.9	0.036	0.718	0.682
Multi-Type-TD-TSR	ICDAR 2019	0.6	0.589	0.593	0.004
Multi-Type-TD-TSR	ICDAR 2019	0.7	0.404	0.397	-0.007
Multi-Type-TD-TSR	ICDAR 2019	0.8	0.137	0.124	-0.013
Multi-Type-TD-TSR	ICDAR 2019	0.9	0.015	0.012	-0.003
Graph-based-TSR	ICDAR 2019	0.6	0.966	0.879	-0.087
Graph-based-TSR	ICDAR 2019	0.7	0.966	0.868	-0.098
Graph-based-TSR	ICDAR 2019	0.8	0.966	0.856	-0.110
Graph-based-TSR	ICDAR 2019	0.9	0.828	0.815	-0.130

VS



Generative AI for Visualization

Four tasks

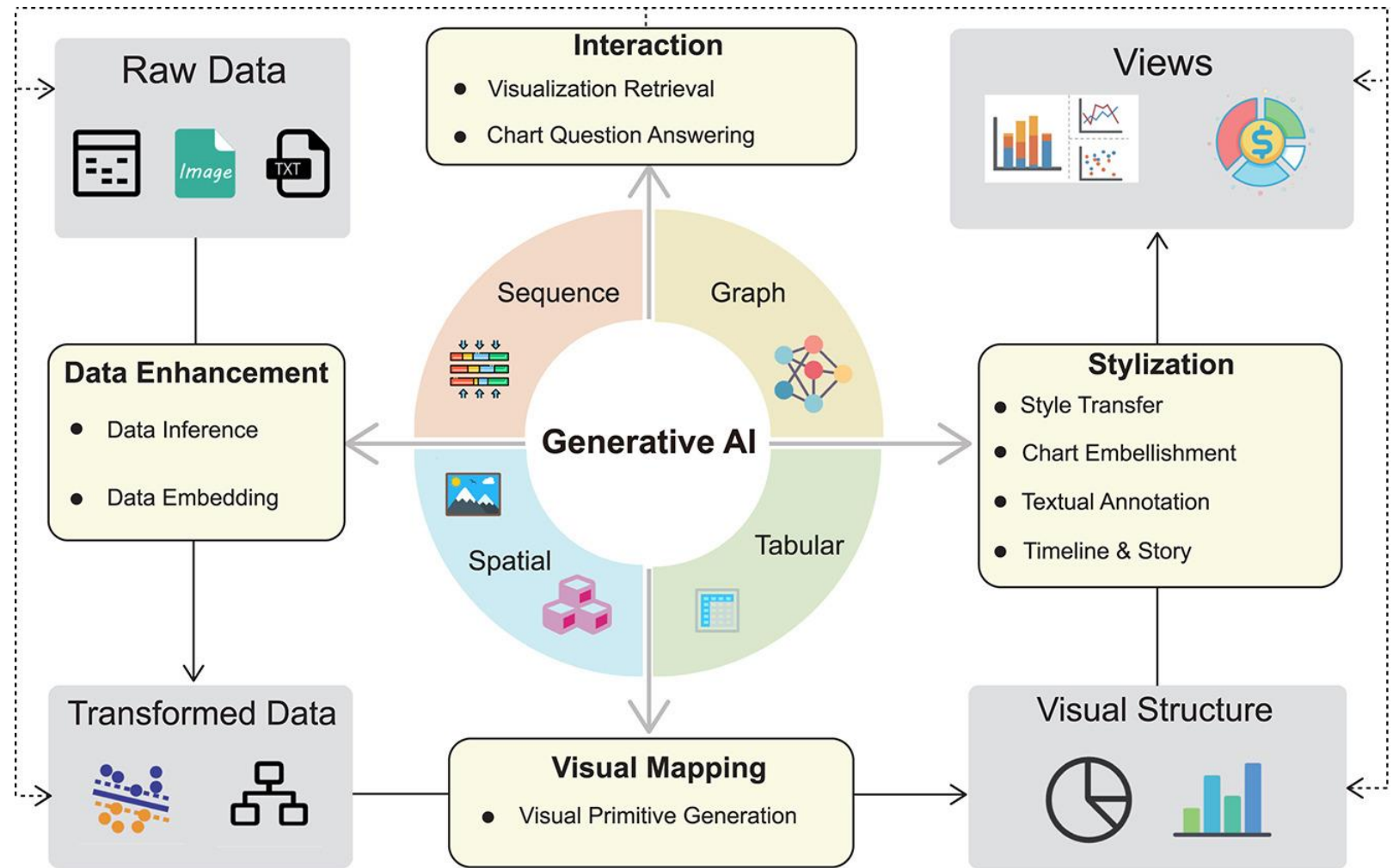
1. data enhancement
2. visual mapping generation
3. stylization
4. interaction

Four types of methods by data structure

- sequence
- tabular
- spatial
- graph

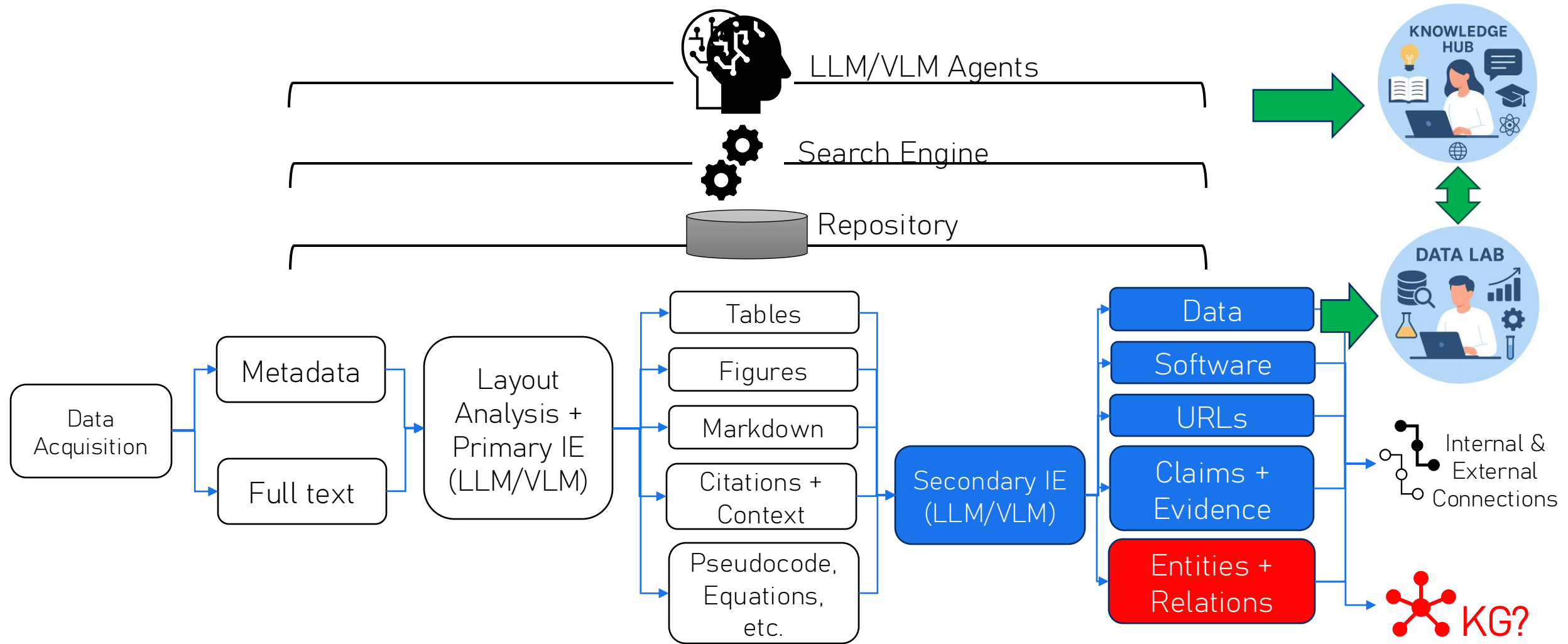
Three challenges

- evaluation
- dataset
- the gap between end-to-end GenAI and VIS.



Ye et al. 2024 Visual Informatics

How to Integrate Everything?



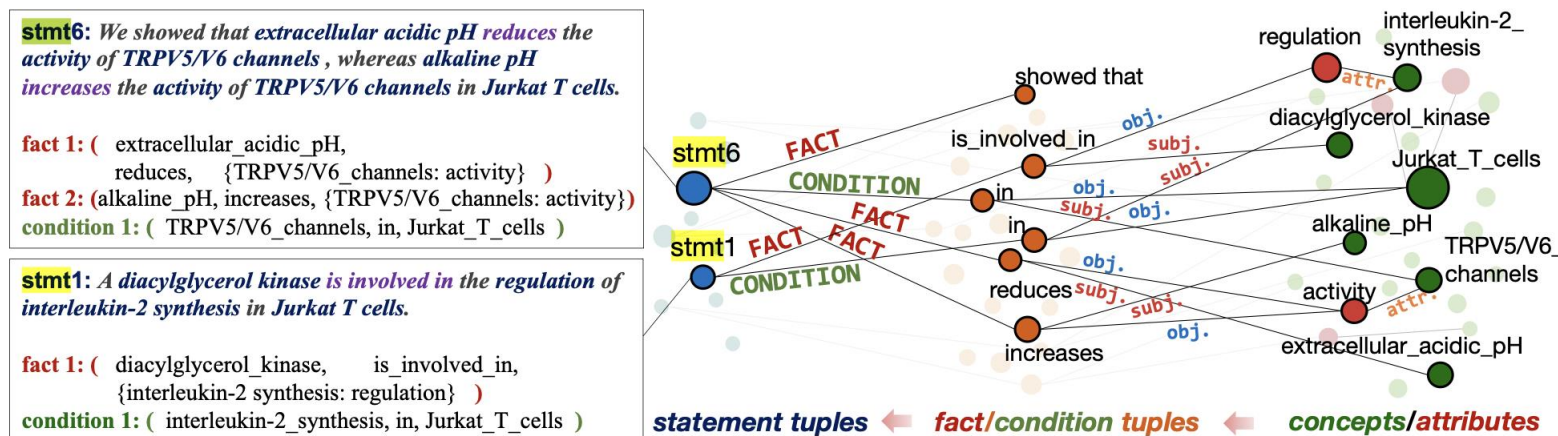
KG?

- We still need KGs because of
 - interpretability: easier to explain facts
 - scalability: relatively easy to add new knowledge
 - provenance document: track entities and relations from original documents
 - metric calculation: e.g., citation and author networks

- However, scientific KGs have intrinsic limitations
 - Conditional knowledge representation: hard to presentation complicated conditional knowledge in scientific statements (see below)
 - Unified ontology: hard to build a unified ontology for all scientific domains
 - Reliable relation recognition is still challenging

- Solution

- Focused on meta-level KGs and KGs for specific domains with well established ontologies, e.g., SNOMED CT.





Do We Still Need Digital Library Search Engines?

- Yes!
- IDLs will be built on existing functionalities of digital library search engines to
 - meet users' needs (e.g., browsing, searching, downloading, metrics)
 - trace provenance of derivatives (e.g., figures, tables, claims)
 - provide explainability (e.g., references, pre-extracted features and metrics)
 - mitigate metadata-level hallucination (authors, publication year, DOIs, URLs, citations, references, etc.)

Summary

- Digital Library Search Engines still holds important roles in future digital libraries
- Intelligent DLs (IDLs) will add more value and facilitate **AI for Science** through Knowledge Hub and Data Lab
- Key Challenges
 - [Technical] How to build a multidomain scientific AI from general purpose AI?
 - Is reading papers enough? How about textbooks and knowledge graphs?
 - [Financial] How to persuade funding agencies to invest research and education? Justify the commercial value.
 - [Infrastructure] Computing power! – light-weight LLMs



Abstract

- Since 2023, there has been a surge of public and research interest in large language models (LLMs), which has significantly shifted the paradigm of information retrieval from returning keyword-based search results to the generation of natural language responses. This shift brings both challenges and opportunities for traditional digital libraries, which have served as a core infrastructure for browsing, searching, and accessing scholarly content. A critical question emerges: What role should digital libraries play in this LLM era? In this keynote, we share our vision of digital libraries in the LLM era. We argue that digital libraries are still indispensable, not only as repositories for digital preservation and provenance but also for trustworthy metadata discovery and verification. We explore how digital libraries can evolve by integrating LLMs and structured knowledge to support advanced services such as automatic data extraction, scholarly comparison, review generation, and science communication for broader audiences. We share preliminary work in this direction, including initiatives on preserving endangered open-access datasets and software, complex table data extraction, scientific claim verification, and assessing research reproducibility.