What is the role of data management in the age of AI?

Chris Bizon Renaissance Computing Institute University of North Carolina, Chapel Hill

May 4, 2025 - me to my family

Bulls game is on TV and they assigned a ball to the pitcher because he didnt make eye contact with the batter. Was confused about the rule so I searched it on google. The google AI starts strong but the answer took a somewhat unusual turn!

2:16 PM · MMS

May 4, 2025 - AI to me

Search Labs | Al Overview

In Major League Baseball, the new pitch clock rules emphasize the importance of eye contact between the pitcher and batter. The rule is designed to reduce game time and prevent stalling. Making eye contact signals that the pitcher is ready to throw, and the batter is expected to be in the box and attentive. Refusing to make eye contact can result in a "D-sist," and after five such instances, the pitcher can challenge the hitter to a duel.



New York Times, 2023



The New York Times

https://www.nytimes.com > athletic > 2023/06/14 > glan ...

Doug Glanville: MLB's new rules have turned umpires into ...

Jul 26, 2023 - ... pitcher or hitter refuses to make eye contact and gets a warning for it. After five Disses, a pitcher may challenge the hitter to a duel.

Garbage In, Garbage Out?

- Not really ... the original article just contains a joke
- Even if the joke (or the opinion) is bad, it's not technically garbage
- However google is pulling in this information is not able to understand its context and appropriate use
- Maybe an LLM should be able to tell the difference, maybe not, but there's also *no metadata* here to help.

What if you want your bot to respond with facts?

- Have an agent look at the question and find relevant facts
- Feed these facts to the LLM along with the user question to generate a response.
- You're using the LLM for what it's good at: text processing, and (hopefully) not what it's only so-so at: knowing things and reasoning
- (Caveat that might be changing over time)
- RAG: Retrieval Augmented Generation

It's up to the RAG-builder to populate this bot with the corpus high-quality documents that the bot will serve!

Put some Triples together, that's a Knowledge Graph



DugBot Knowledge Graph Agent

Courtesy Suparna Goswami (RENCI)

• Query Dug's knowledge Graph using biomedical concepts



There's more to AI than LLMs

ChatGPT o3-mini-high \vee

Analyze every possible combination of FDA approved drugs and known diseases to suggest the top 10 most likely drug repurposing opportunities

Reasoned for 44 seconds >

Below is a conceptual synthesis—drawing on systems pharmacology, transcriptomics, network analysis, and literature evidence—that outlines a set of ten promising repurposing hypotheses. In practice, "analyzing every possible combination" of FDA-approved drugs versus all known diseases is an enormous computational task. Instead, researchers typically integrate multi-omics data, drug– target networks, and epidemiological insights to highlight candidates with strong mechanistic rationale. The list below represents candidates that recurrently emerge from such approaches. (Keep in mind that these hypotheses require further experimental and clinical validation.)



Embeddings turn topology into geometry (points in space)



Embeddings can be optimized in different ways, reflecting different assumptions about the embedding space.

This is analogous to the way that large language models embed words

We can derive answers, but where do they come from?

井

This system is for research purposes and is not meant to be used by clinical service providers in the course of treating patients.

NCATS Translator	[프] Workspace	🕙 Search History	🦪 Send Feedback	? Help	Pref	erences	Log Out			
	What drugs may treat condition Immune-mediated Necrotizing I	is related to Myopathy 岱?								
Results Showing 1-10 of 510 Results ∓ Filters								1 2	34	51 >
NAME Mycophenolate Mofetil						evii	DENCE Dications (92) nical Trials (4)	PATHS ()	5.00)↓
Immunosuppressive Agent (18) Pro-agent (19) Ec 1.* (The morpholinoethyl ester of mycophenolic acid (MPA) with potent vivo, the active metabolite, MPA, reversibly inhibits inosine 5'-monophenolic acid (MPA) Simple active metabolite, MPA, reversibly inhibits inosine 5'-monophenolic acid (MPA)	xidoreductase) Inhibitor (40) Antimicrobial Agent (32) mmunosuppressive properties. Mycophenolate stops T-cell and B-ce hosphate dehydrogenase, an enzyme involved in the de novo synth	Il proliferation through select esis of guanine nucleotides. N	ive inhibition of the de novo IPA displays high lymphocyt	pathway of purine e specificity and cy	e biosynthe ytotoxicity	Sou sis. In lue to	urces (7)			
Lymphocyte Immune Globulin, Anti-t UMLS Semantic Type: STY:T116; UMLS Semantic Type: STY:T121; UM	hymocyte Globulin S Semantic Type: STY:T129				B	Sou	urces (1)	1 Path	5.00	~
Prednisone					e	Clir	olications (93) hical Trials (5) urces (7)	18 Paths	5.00	~

It's important to track provenance





ORION creates and exposes metadata

- Data Versions, of both input and output
- Tool Versions
- Processes
- Graph characteristics
- Evidence
- Provenance

Metadata exposed in machine readable formats

```
"graph id": "RobokopKG",
"graph name": "ROBOKOP KG",
"graph description ": "ROBOKOP (KG) is an open-source biomedical KG that supports the ROBOKOP application and
currently contains millions of biomedical relationships ... ",
"graph url": "http://robokopkg.renci.org/browser/ ",
"graph version": "4901b2bc764444ea",
"sources": [],
"subgraphs": [],
"build status": "stable",
"build time": "07-30-24 23:37:31",
"build error": null,
"final node count ": 9800836,
"final edge count ": 137235468,
"merged nodes": 280904,
"merged edges ": 4676433,
"qc results": {},
"neo4j dump":
"https://stars.renci.org/var/plater/bl-4.2.1/RobokopKG/4901b2bc764444ea/graph 4901b2bc764444ea.db.dump
```

A dirty little secret

- Very few of the sources that ORION ingests have sufficient metadata.
- Almost none of them have metadata in a machine readable format.
- None of them have metadata in a machine readable format that conforms to any common standard.

Croissant: A Metadata Format for ML-Ready Datasets

Mubashara Akhtar⁷, Omar Benjelloun⁴, Costanza Conforti⁴, Pieter Gijsbers¹², Joan Giner-Miguelez¹³, Nitisha Jain⁷, Michael Kuchnik⁸, Quentin Lhoest⁵, Pierre Marcenac⁴, Manil Maskey⁹, Peter Mattson⁴, Luis Oala³, Pierre Ruyssen⁴, Rajat Shinde¹⁰, Elena Simperl^{7,11}, Goeffry Thomas^{4,6}, Slava Tykhonov², Joaquin Vanschoren¹², Jos van der Velde¹², Steffen Vogler¹, Carole-Jean Wu^{8*} ^{*}Authors in alphabetical order ¹Bayer, ²DANS-KNAW, ³Dotphoton, ⁴Google, ⁵Hugging Face, ⁶Kaggle, ⁷King's College London, ⁸Meta, ⁹NASA, ¹⁰NASA IMPACT & UAH, ¹¹Open Data Institute, ¹²TUE & OpenML, ¹³Universitat Oberta de Catalunya



1	{ "@type": "sc:Dataset",
2	"name": "PASS",
3	"dct:conformsTo":
	"http://mlcommons.org/croissant/1.0",
4	"description": "PASS is a large-scale image dataset
	that does not include",
5	<pre>"citeAs": "@Article{asano21pass, title = \"PASS: An</pre>
	ImageNet replacement",
6	"license": "cc-by-4.0",
7	"url": "https://www.robots.ox.ac.uk/~vgg/data/pass/"}

Figure 3: Dataset metadata for the PASS dataset.

1	{	"@id": <mark>"pass0"</mark> ,
2		"@type": "cr:FileObject",
3		"contentUrl":
		"https://zenodo.org/6615455/PASS.0.tar",
4		"sha256": "0be3a104d6257d83296460b419f82c71",
5		<pre>"encodingFormat": "application/x-tar"},</pre>
6	{	"@id": <mark>"image-files"</mark> ,
7		"@type": "cr:FileSet",
8		"containedIn": {"@id": <mark>"pass0"</mark> }
9		"includes": "*.jpg",
10		"encodingFormat": "image/jpeg"},
11		{ "@id": "metadata",
12		"@type": "cr:FileObject",
13		"contentUrl":
		"https://zenodo.org/6615455/pass_metadata.csv",
14		"sha256": "0b033707ea49365a5ffdd14615825511",
15		<pre>"encodingFormat": "text/csv"}</pre>

Figure 4: Definitions of Resources for the PASS dataset.

1	1
2	"@id": "images",
3	"@type": "cr:RecordSet",
4	"key": "images/hash",
5	"field": [
6	{
7	"@id": "images/image content".
8	"@type": "cr:Field"
	"dataType", "sc:ImageObject"
10	"source", f
10	"fileSet", ("Aid", "image_files")
11	"inteset . { eid . image lifes },
12	extract . (Therroperty : content y
13	
14	3,
15	1
16	"@id": "images/hash",
17	"@type": "cr:Field",
18	"dataType": "sc:Text",
19	"source": {
20	"fileSet": {"@id": "image-files"},
21	"extract": {"fileProperty": "filename"},
22	"transform": {"regex": "([^\\/]*)\\.jpg"}
23	},
24	"references": {
25	"fileObject": {"@id": "metadata"},
26	"column": "hash"
27	}
28	},
29	{
30	"@id": "images/coordinates",
31	"@type": "cr:Field",
32	"dataType": "sc:GeoCoordinates",
33	"subField": [
34	{
35	"@id": "images/coordinates/latitude",
36	"@type": "cr:Field",
37	"source": {
38	"fileObject": {"@id": "metadata"},
39	"column": "latitude"
40	}
41	},
42	(
43	"@id": "images/coordinates/longitude".
44	"@type": "cr:Field",
45	"source": {
46	"fileObject": {"@id": "metadata"}.
47	"column": "longitude"
48	}
49	}
50	1
51	3
52)
34	,

Figure 5: Definition of a RecordSet that joins images and structured metadata from the PASS dataset.

What is the role of Data Management in the age of AI? What is the role of AI in the age of Data Management?

We are developing a prototype "Auto-croissant" pipeline





Thanks to

RENCI

Suparna Goswami Yaphet Kebede Yogesh Kale Gaurav Vaidya Evan Morris James Chung Kathleen Carter Ola Olasunkanmi

<u>UNC</u>

Alex Tropsha Marcello DeLuca Kushal Koirala Elvin Ding Melissa Haendel Niko Matentzoglu

Every Cure

David Fajgenbaum Pascal Brokmeier Lee Lancashire Elliott Sharpe

Funding provided by

NHLBI (DugBot) ARPA-H (MATRIX)