

# ToxDataCommons: Driving toxicology research forward through (meta)data sharing

Rance Nault

Department of Pharmacology & Toxicology

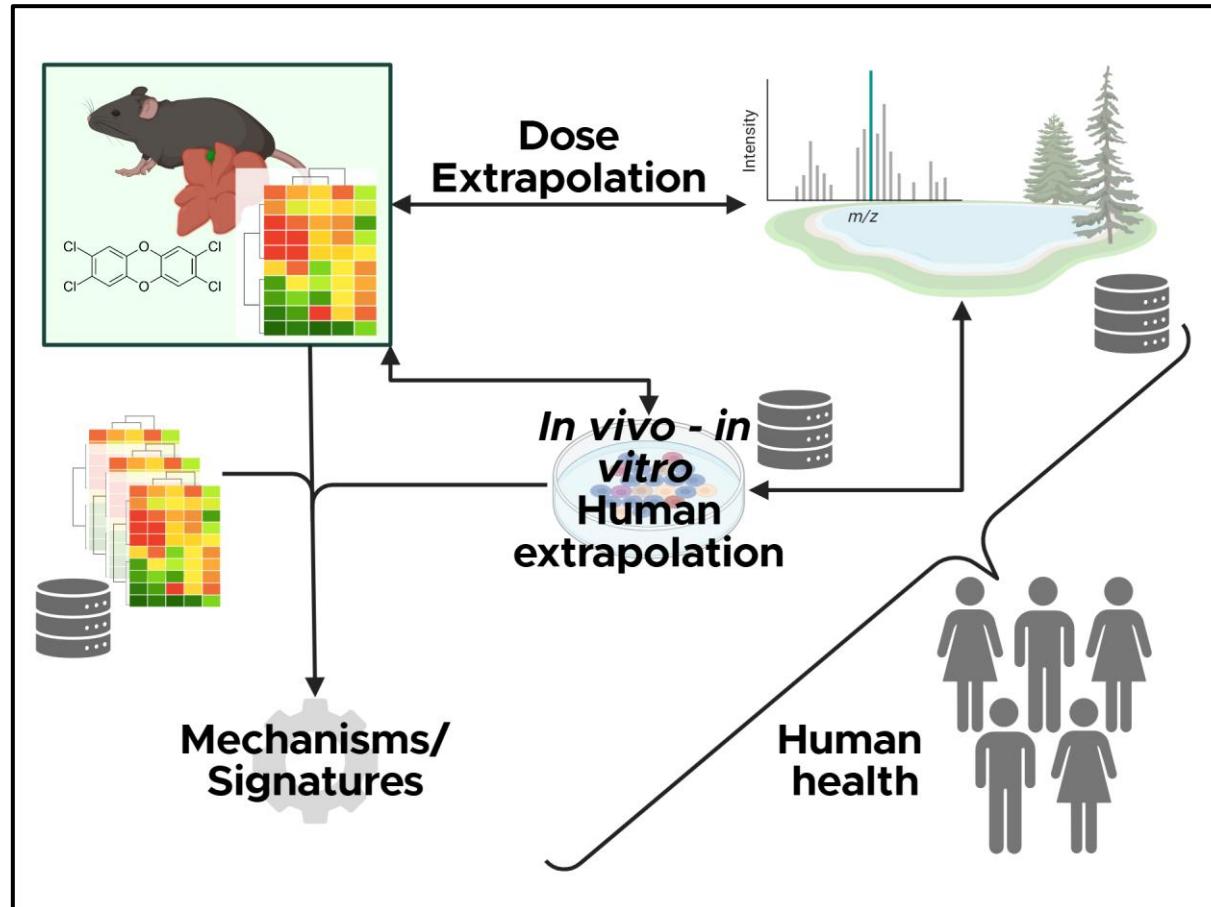
Institute for Integrative Toxicology

Superfund Research Center Data Management and Analysis Core

Michigan State University

[naultran@msu.edu](mailto:naultran@msu.edu)

# WHY CARE ABOUT DATA SHARING



We generate huge amounts of data.

We also **NEED** more high-quality data to complete the picture.



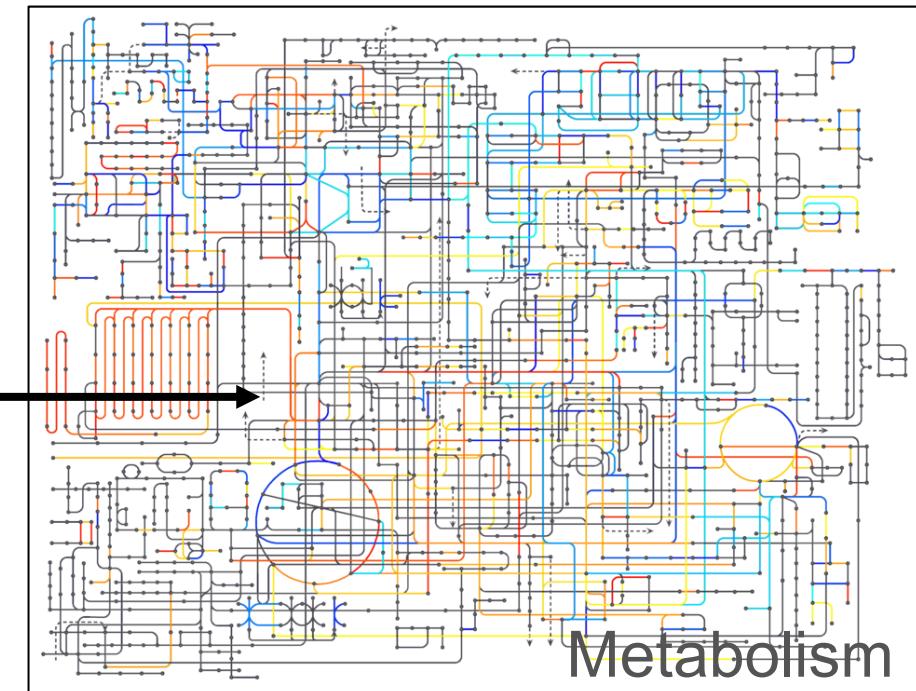
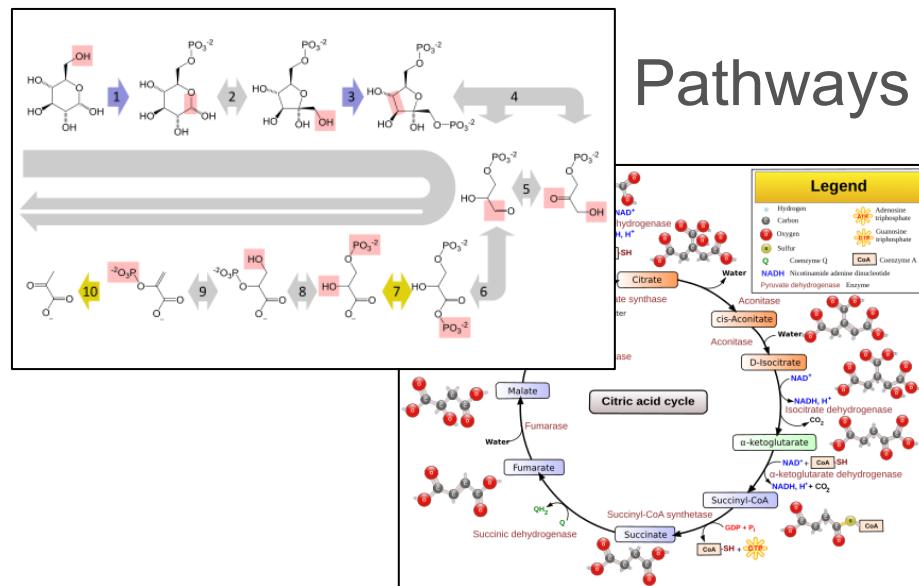
*Image created with ChatGPT (DALL·E), 2025.*

# WHY CARE ABOUT DATA SHARING

Reproducibility	Validation (e.g., NAMs)	Discovery
<ul style="list-style-type: none"><li>• Can we get to the same answer using the same dataset?</li><li>• Do the same study designs produce the same outcomes?</li><li>• Can we explain differences in outcomes when metadata is richer?</li></ul>	<ul style="list-style-type: none"><li>• Benchmarking <i>in vitro</i> results to <i>in vivo</i> results (reducing animal use)</li><li>• Can we build statistical/mathematical models to capture address differences between human, rodents, cells, ...</li></ul>	<ul style="list-style-type: none"><li>• Leverage data from all studies of chemical X to refine safety limits, mechanisms, ...</li><li>• Using data hungry methods such as artificial intelligence to gain new knowledge</li></ul>

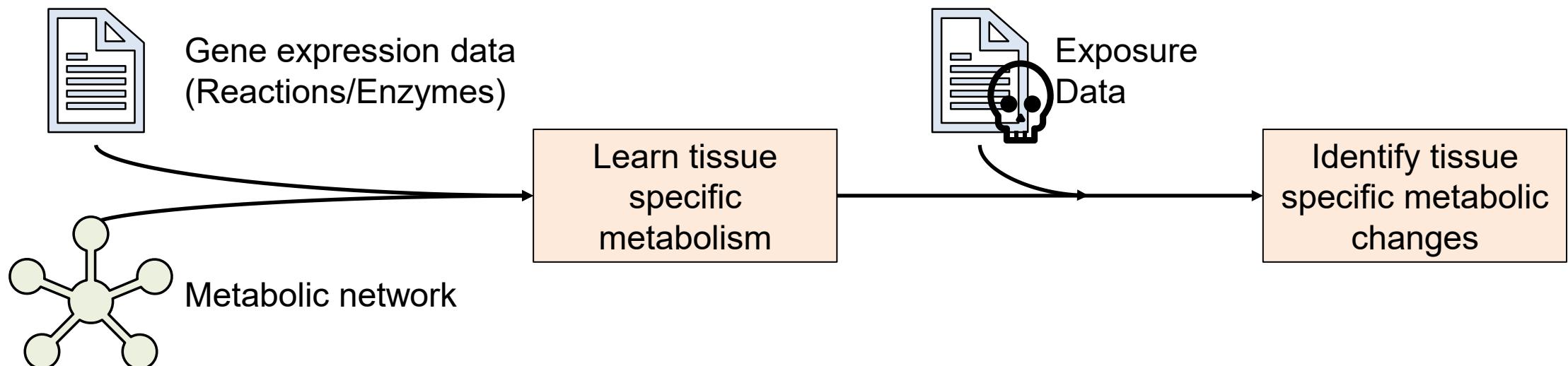
# ROLE OF DATA IN AI APPLICATIONS

**Research Question:** Can we better characterize metabolic pathway changes by environmental toxicants considering metabolism as a network of reactions.



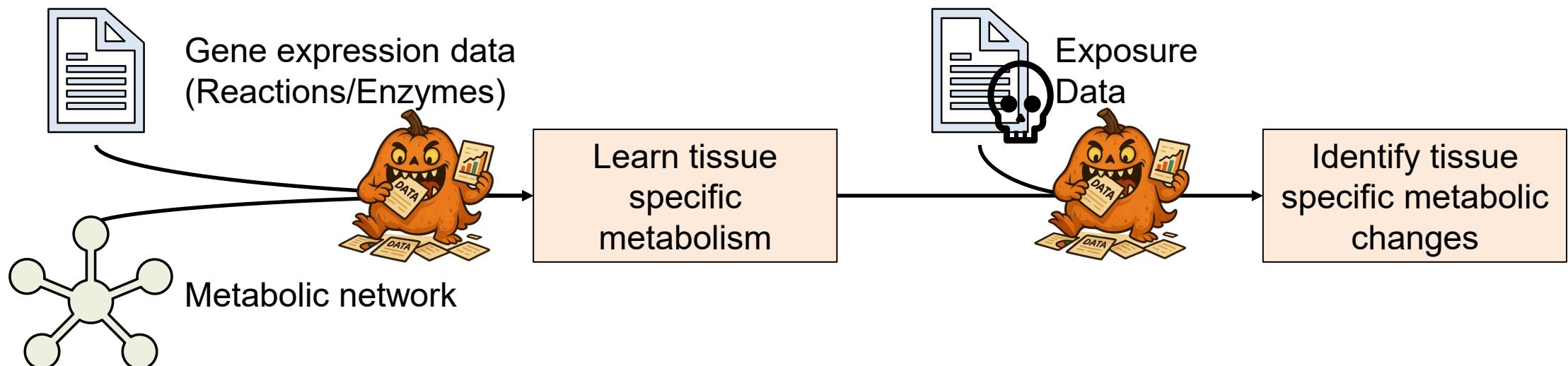
# TRAINING A GRAPH NEURAL NETWORK

A Graph Neural Network (GNN) is a type of artificial intelligence applied to data that can be represented as a graph (e.g., social networks, chemical structures)

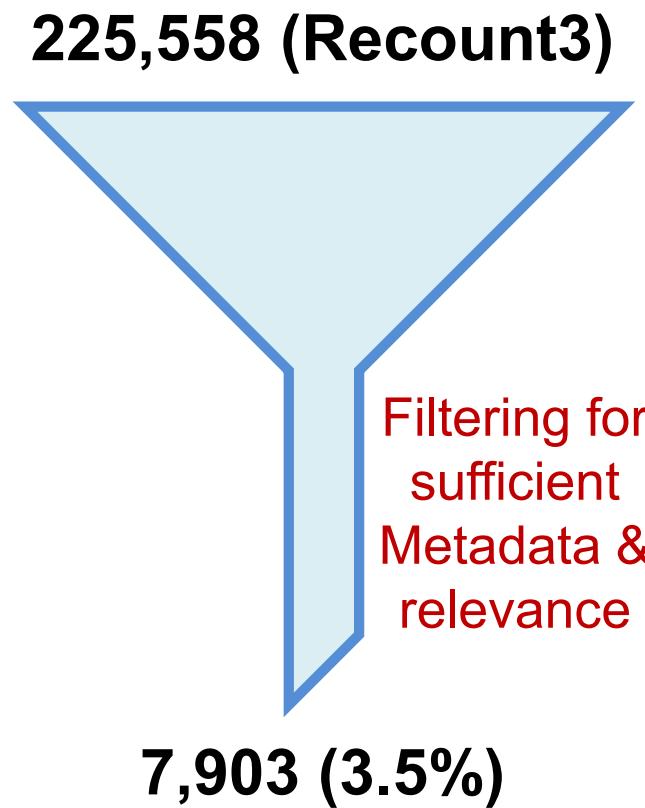


# TRAINING A GRAPH NEURAL NETWORK

A Graph Neural Network (GNN) is a type of artificial intelligence applied to data that can be represented as a graph (e.g., social networks, chemical structures)



# TRAINING AI - THE METADATA PROBLEM

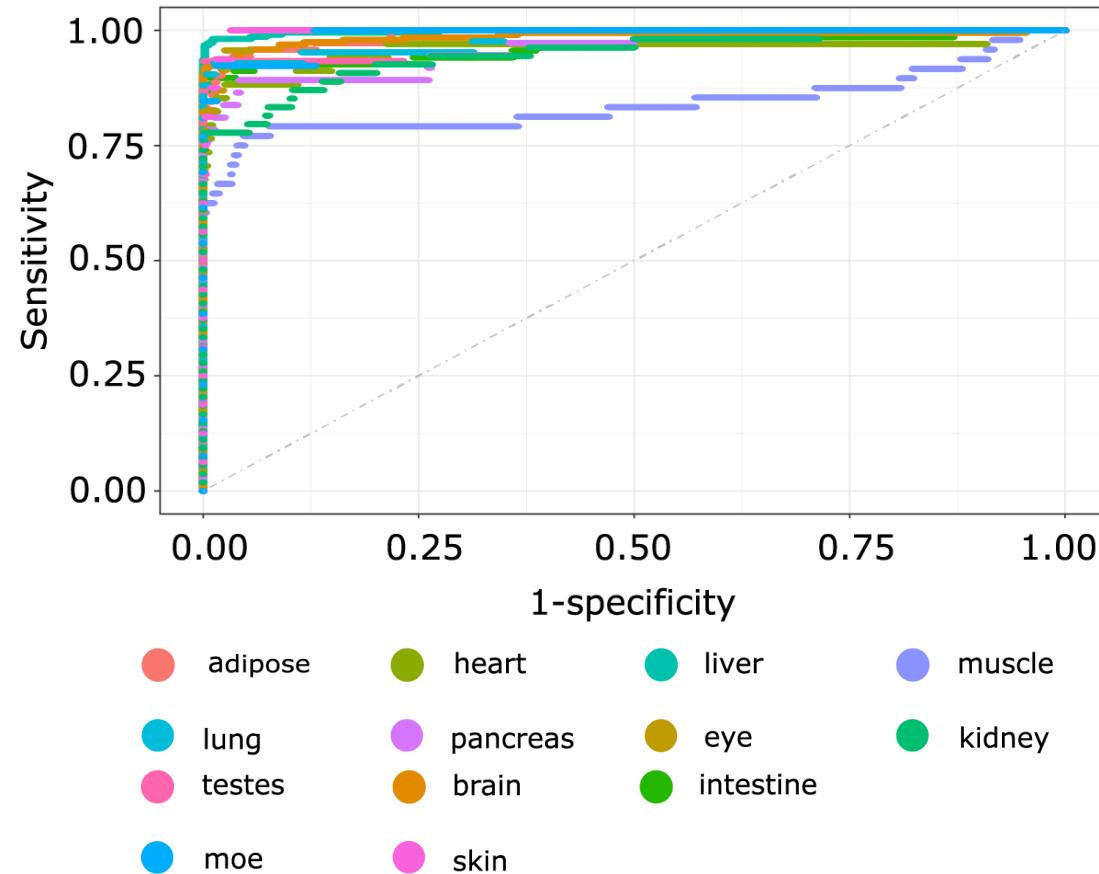


We started by downloading from GEO all metadata for 225,558 mouse samples.

- Many were excluded for relevance (e.g., *in vitro*, *genetically modified organisms*, *developmental (E0-E19)*, ...).
- A significant portion we couldn't determine the relevance (*missing strain, sex, treatment details*, ...).

Toxicant perturbations samples were even more challenging to keep!

# CLASSIFICATION OF TISSUE METABOLISM



Despite keeping only 3.5% of samples, there was sufficient data to accurately classify tissues based on expression of genes in the metabolic network.

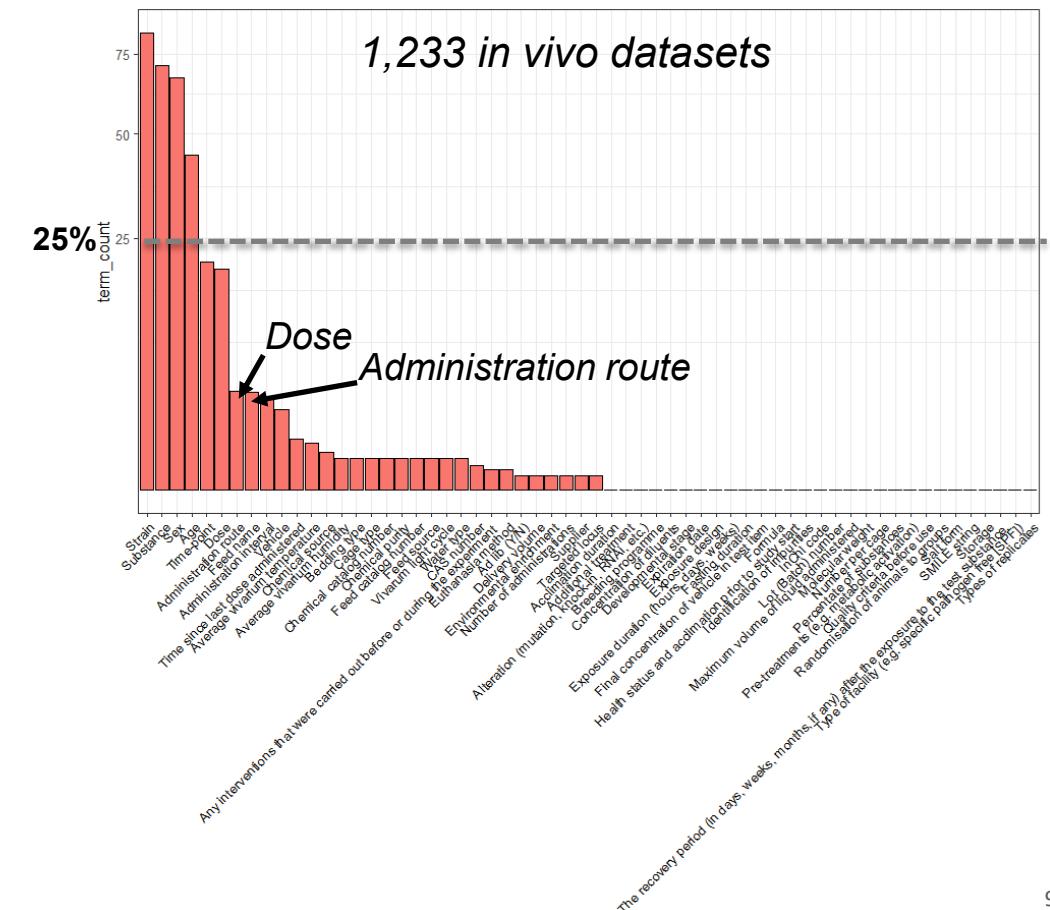
**This was much more challenging for toxicant-driven gene expression changes.**

# UNDER-REPORTING OF METADATA

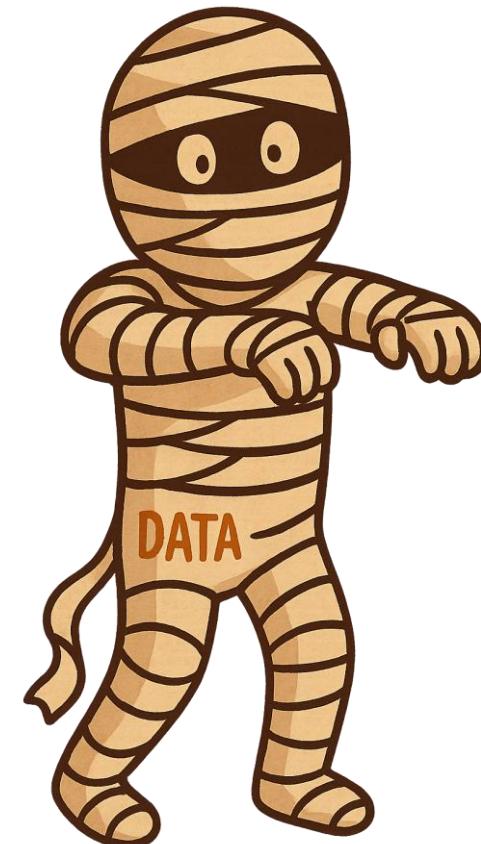
Several toxicant exposure gene expression studies could be identified.

Significant gaps in metadata reporting making it more challenging to provide AI enough information.

These challenges are also important in other data integration and reuse activities.

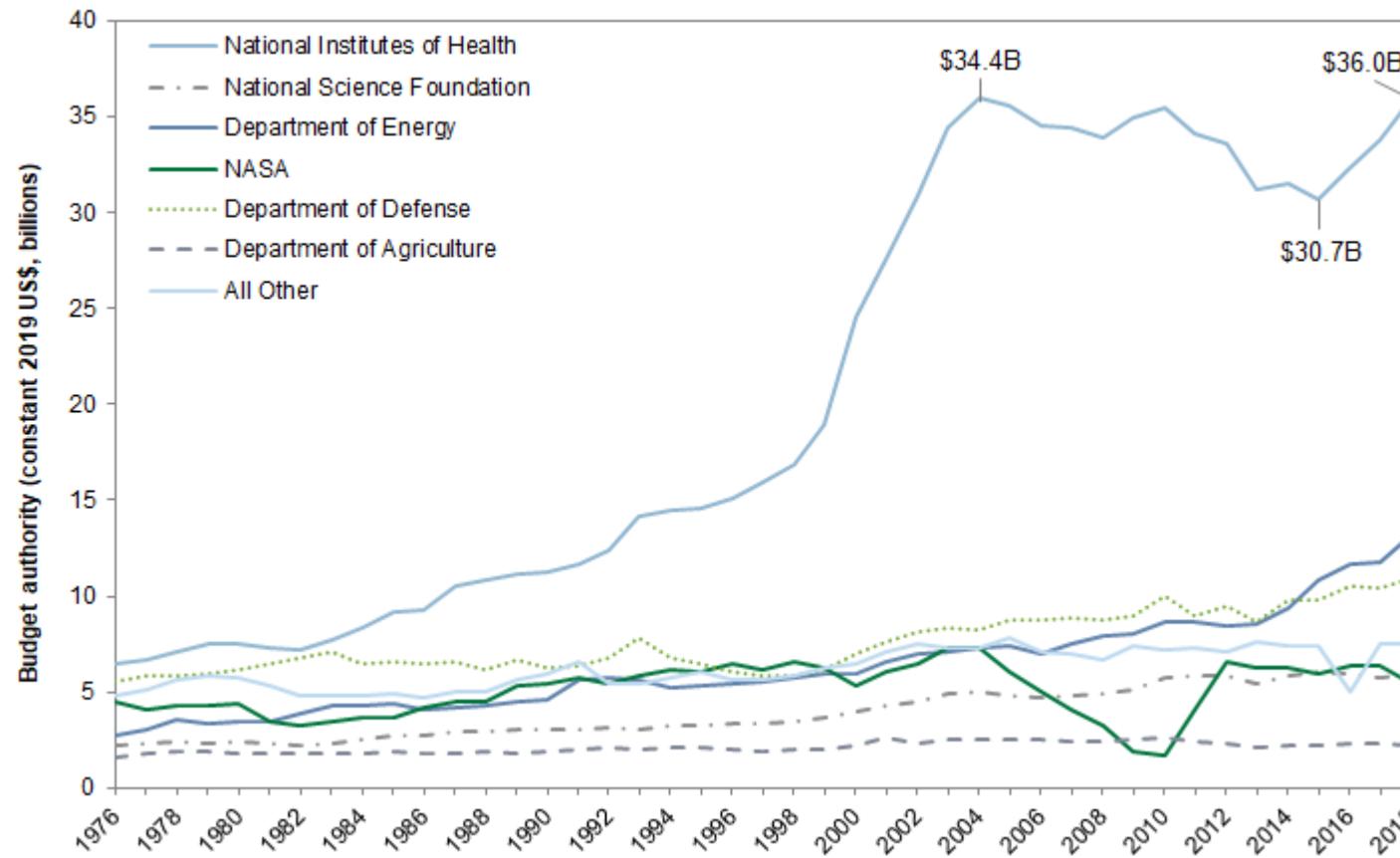


# UNRAVELING DATA MANAGEMENT



*Image created with ChatGPT (DALL·E), 2025.*

# DATA HAS VALUE – BUT WE CAN ADD MORE



Billions of dollars are given to support the generation of new research data annually.

How can we maximize the return on investment and minimize unnecessary repetition.

*Reproducibility is still important!*

# WIDESPREAD INITIATIVES

**FAIR Principles** have been central to new NIH data management policies.

*“NIH encourages data management and data sharing practices **consistent with the FAIR data principle**”*

NIH aims to require data sharing of all research data.

*“... regardless of whether the data are used to support scholarly publications”*

SRP ... DMAC to support the management and integration of data assets.



National  
Science  
Foundation



National Institute of  
Environmental Health Sciences  
*Superfund Research Program*

# FAIR DATA PRINCIPLES



Image: <https://www.linkedin.com/pulse/why-your-data-must-fair-david-crosswell/>

FAIR principles represents a set of broad concepts, not a specific set of rules.

Open Access | Published: 15 March 2016

## The FAIR Guiding Principles for scientific data management and stewardship

Mark D. Wilkinson, Michel Dumontier, IJsbrand Jan Aalbersberg, Gabrielle Appleton, Myles Axton, Arie Baak, Niklas Blomberg, Jan-Willem Boiten, Luiz Bonino da Silva Santos, Philip E. Bourne, Jildau Bouwman, Anthony J. Brookes, Tim Clark, Mercè Crosas, Ingrid Dillo, Olivier Dumon, Scott Edmunds, Chris T. Evelo, Richard Finkers, Alejandra Gonzalez-Beltran, Alasdair J.G. Gray, Paul Groth, Carole Goble, Jeffrey S. Grethe, Jaap Heringa, Peter A.C. 't Hoen, Rob Hooft, Tobias Kuhn, Ruben Kok, Joost Kok, Scott J. Lusher, Maryann E. Martone, Albert Mons, Abel L. Packer, Bengt Persson, Philippe Rocca-Serra, Marco Roos, Rene van Schaik, Susanna-Assunta Sansone, Erik Schultes, Thierry Sengstag, Ted Slater, George Strawn, Morris A. Swertz, Mark Thompson, Johan van der Lei, Erik van Mulligen, Jan Velterop, Andra Waagmeester, Peter Wittenburg, Katherine Wolstencroft, Jun Zhao & Barend Mons  Show fewer authors

[Scientific Data](#) 3, Article number: 160018 (2016) | [Cite this article](#)

487k Accesses | 4537 Citations | 2021 Altmetric | [Metrics](#)

# FAIR DATA PRINCIPLES: ROLE OF METADATA



Identifiable



compliant



Versioned



Metadata

Image: <https://www.liris.fr/>

## Box 2 | The FAIR Guiding Principles

### To be Findable:

- F1. (meta)data are assigned a globally unique and persistent identifier
- F2. data are described with rich metadata (defined by R1 below)
- F3. metadata clearly and explicitly include the identifier of the data it describes
- F4. (meta)data are registered or indexed in a searchable resource

### To be Accessible:

- A1. (meta)data are retrievable by their identifier using a standardized communications protocol
  - A1.1 the protocol is open, free, and universally implementable
  - A1.2 the protocol allows for an authentication and authorization procedure, where necessary
- A2. metadata are accessible, even when the data are no longer available

### To be Interoperable:

- I1. (meta)data use a formal, accessible, shared, and broadly applicable language for knowledge representation.
- I2. (meta)data use vocabularies that follow FAIR principles
- I3. (meta)data include qualified references to other (meta)data

### To be Reusable:

- R1. meta(data) are richly described with a plurality of accurate and relevant attributes
  - R1.1. (meta)data are released with a clear and accessible data usage license
  - R1.2. (meta)data are associated with detailed provenance
  - R1.3. (meta)data meet domain-relevant community standards

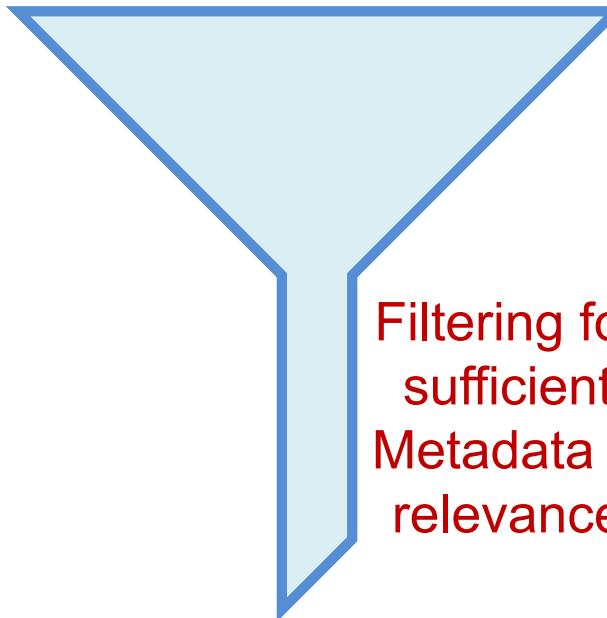
set of  
c set

2

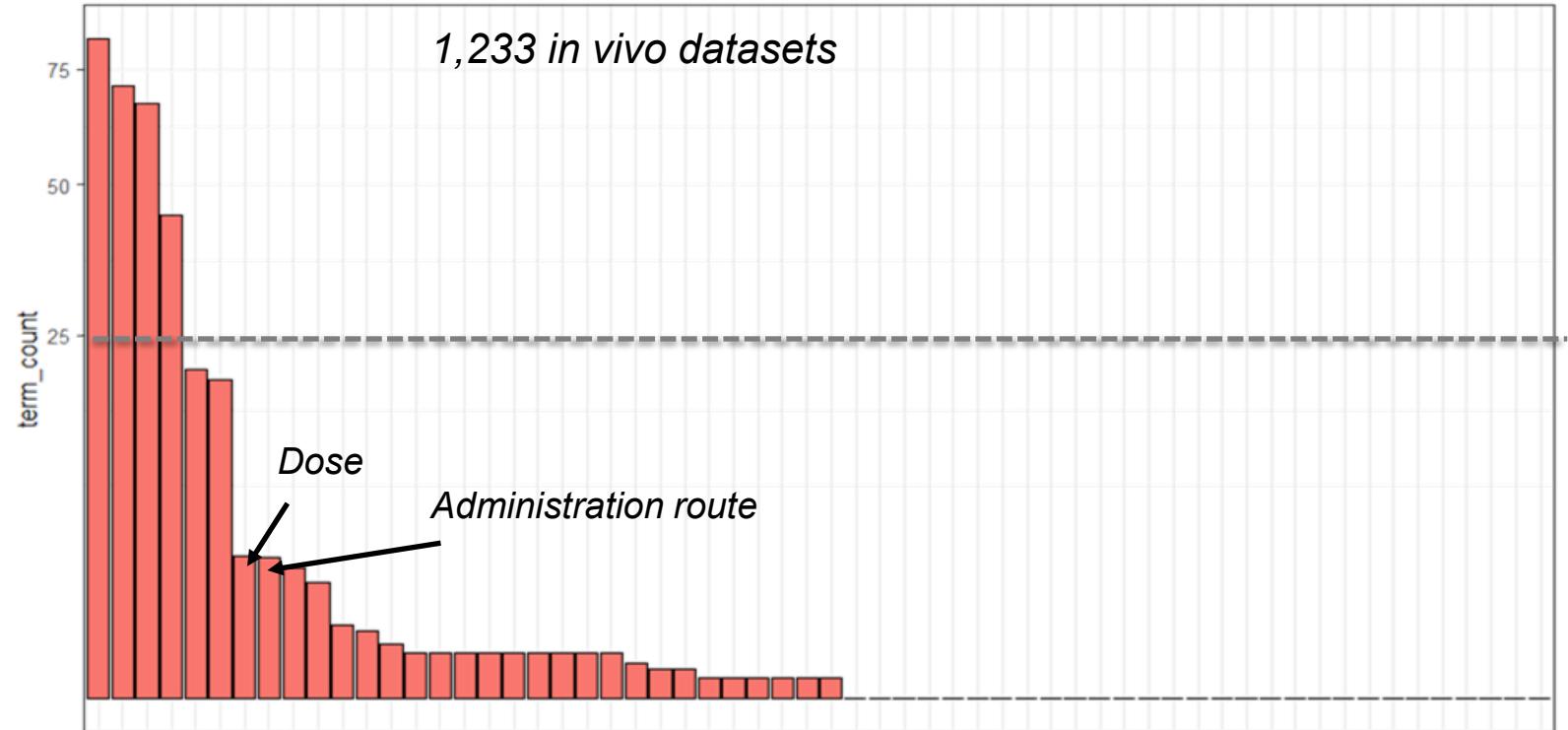
yles Axton, Arie  
e. Jildau Bouwman  
ds. Chris T. Evelo  
le. Jeffrey S. Grethe  
J. Lusher, Maryann E.  
is. René van Schaik  
Morris A. Swertz  
ter. Peter  
os.

# ARE WE ACHIEVING FAIR DATA?

**225,558 (Recount3)**



**7,903 (3.5%)**



# WE ARE GOOD AT REPORTING METADATA

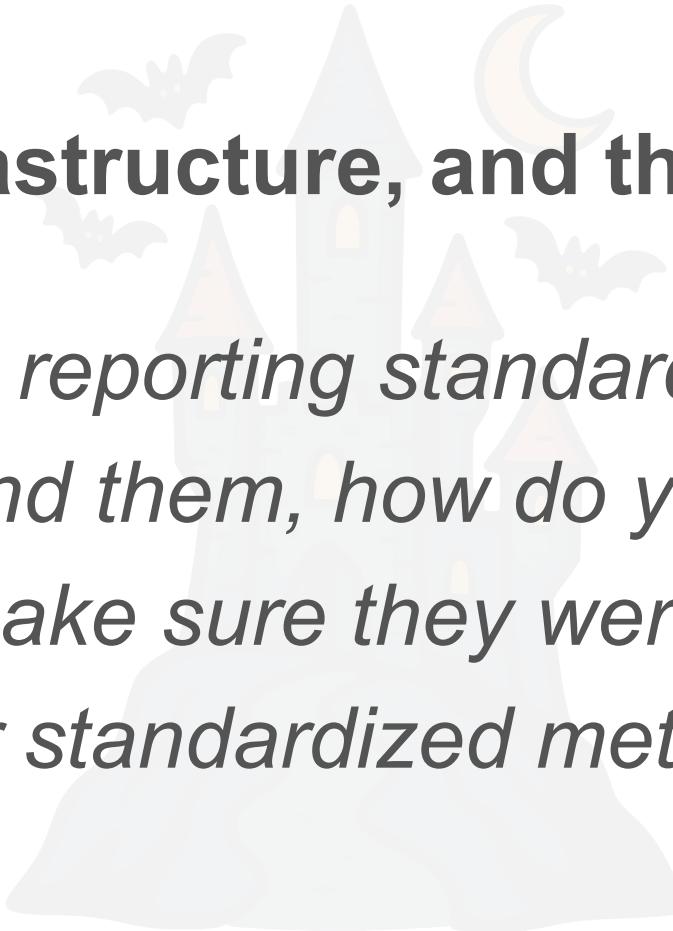
These studies had *poor reporting* in our analysis of GEO datasets

Study Name	Publication (manually identified terms)	GEO metadata terms
<b>GSE116653</b> : Analysis of the role of long non-coding RNAs involved in silica-induced pulmonary fibrosis of rat	38	0
<b>GSE18858</b> : Transcriptional Biomarkers to Predict Female Mouse Liver Tumors in Rodent Cancer Bioassays - A 26 Chemical Set	45	0
<b>GSE70583</b> : Response of rat liver to GW7647	17	1

# HOW DID WE GET HERE ...



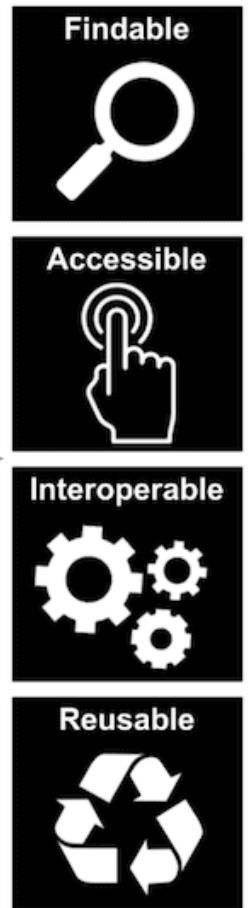
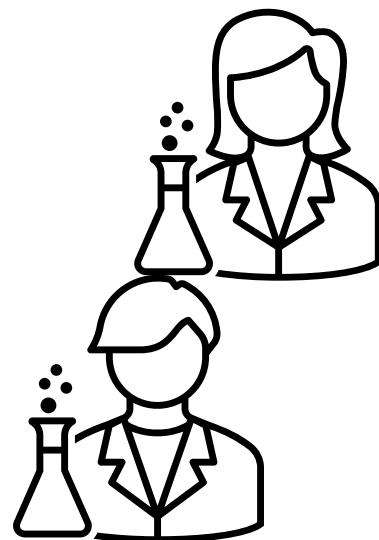
# HOW DID WE GET HERE ...



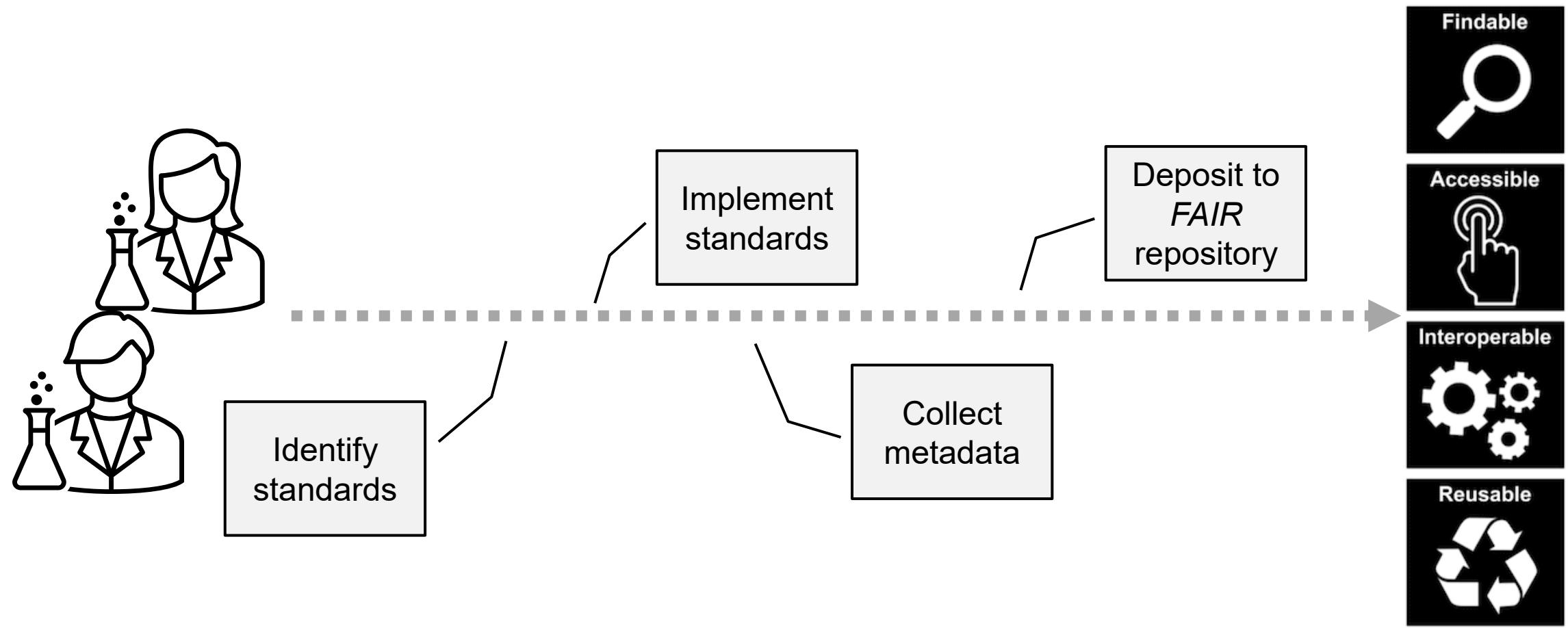
**It's the tools, infrastructure, and their accessibility ...**

- 1) *Where are the reporting standards?*
- 2) *When/If you find them, how do you apply them?*
- 3) *How do you make sure they were applied correctly?*
- 4) *You have your standardized metadata ... now what?*

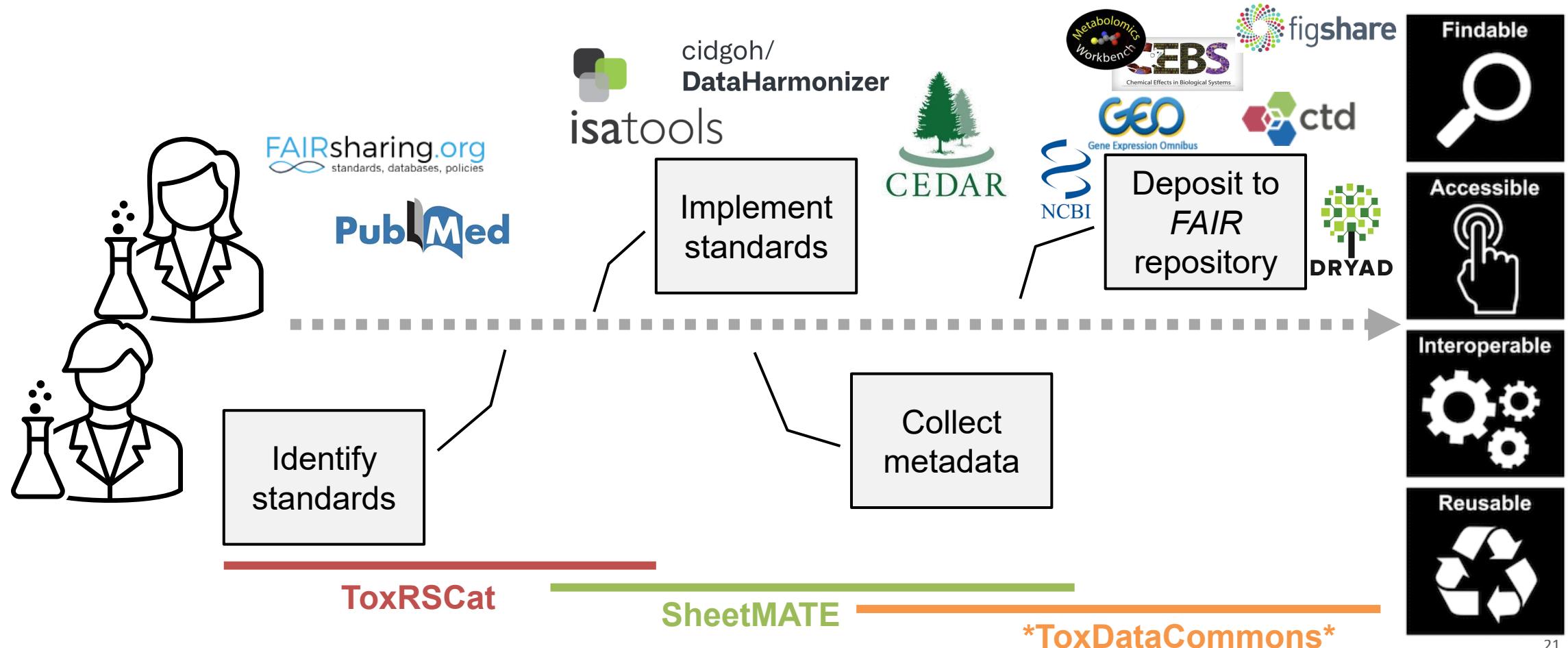
# FROM THE LAB TO FAIR



# FROM THE LAB TO FAIR

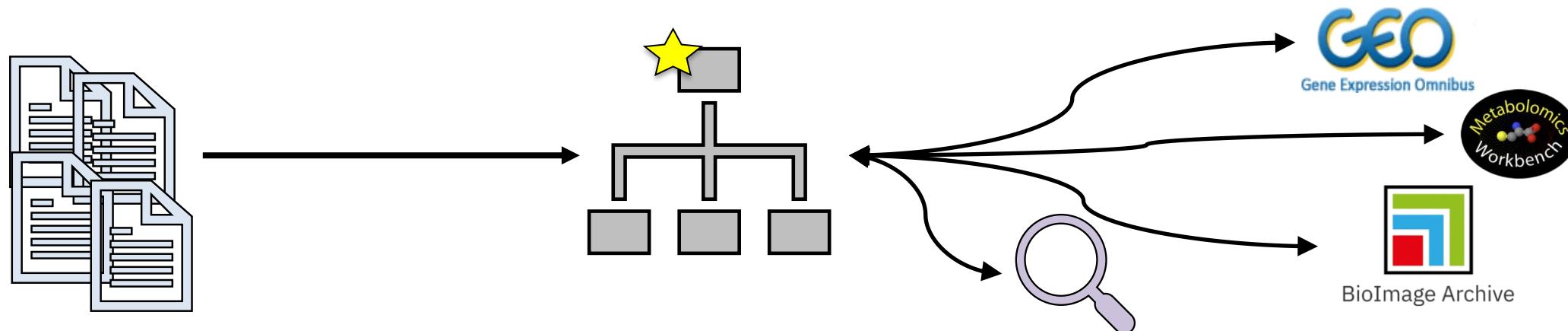


# FROM THE LAB TO FAIR



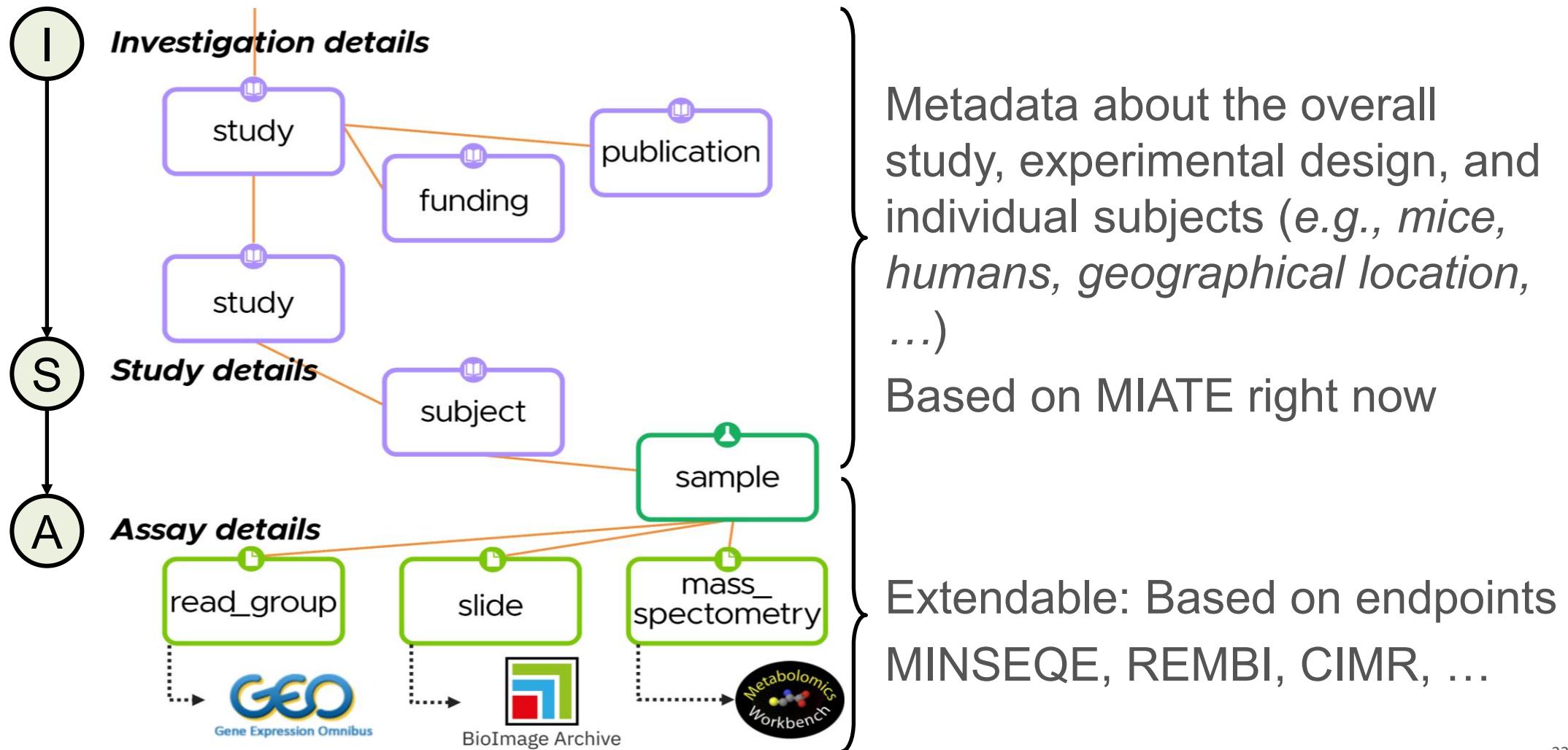
# WHAT IS ToxDataCommons?

ToxDataCommons is a FAIR resource aimed at helping the environmental health research community collect, share, and discover relevant AI-ready datasets.



We're not replacing existing repositories, but allowing researchers to add value to their public datasets

# UNDERLYING DATA MODEL



# UNDERLYING DATA MODEL

administrative

subject description of subject.

				JSON	TSV	Close
subject						
sex	<ul style="list-style-type: none"> <li>male</li> <li>female</li> <li>Not applicable</li> <li>Missing</li> <li>Not collected</li> <li>Not provided</li> <li>Restricted access</li> </ul>	No	No Description			
start_date	• string	No	The date the subject entered the study, which may differ from experiment_start_date (e.g., for acclimation).			
start_date_age	• number	No	Age of the subject at the start of the study, in days. Used to normalize for developmental stage.			
strain	<ul style="list-style-type: none"> <li>Not applicable</li> <li>Missing</li> <li>Not collected</li> <li>Not provided</li> <li>Restricted access</li> <li>C57BL/6NCrl</li> <li>C57BL/6J</li> <li>C57BL/6NCrl_PkmΔDRE</li> </ul>	No	No Description			

Lists all expected metadata for a subject

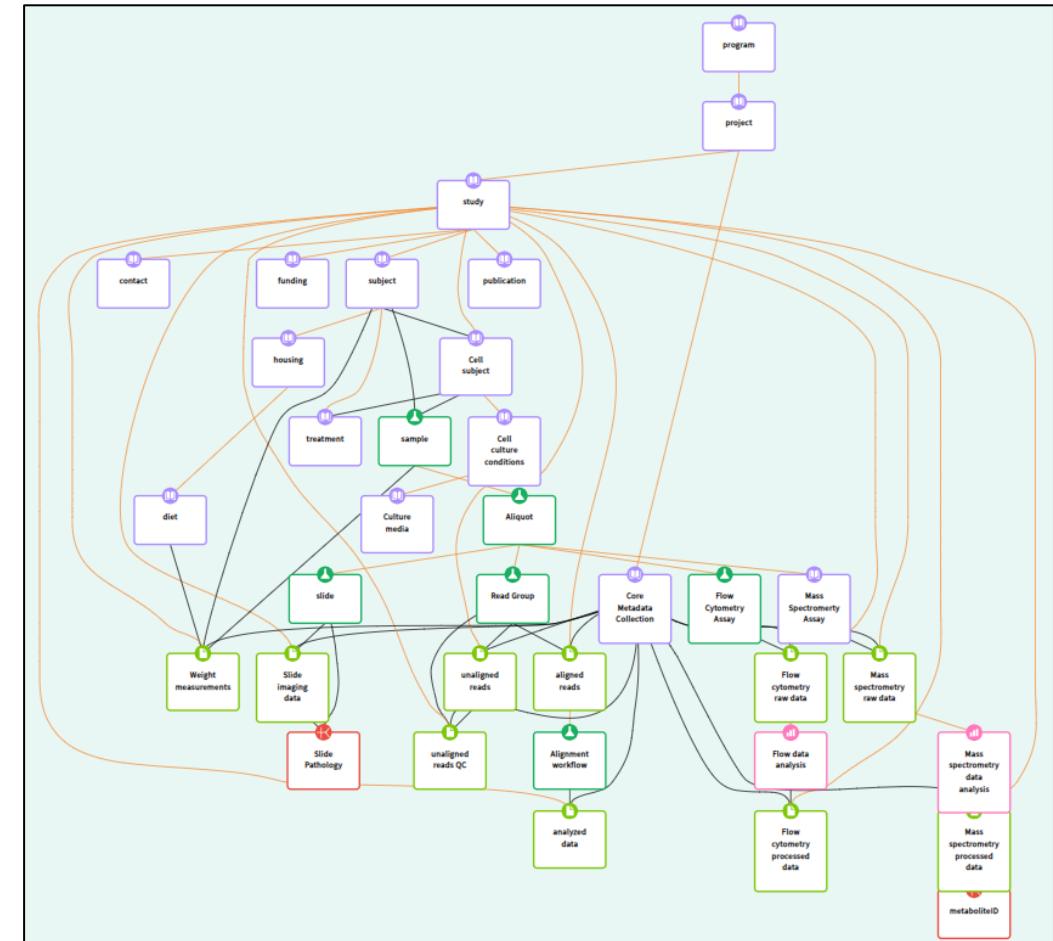
Prevents typos (e.g., BABL/C vs. BALB/C)

# UNDERLYING DATA MODEL

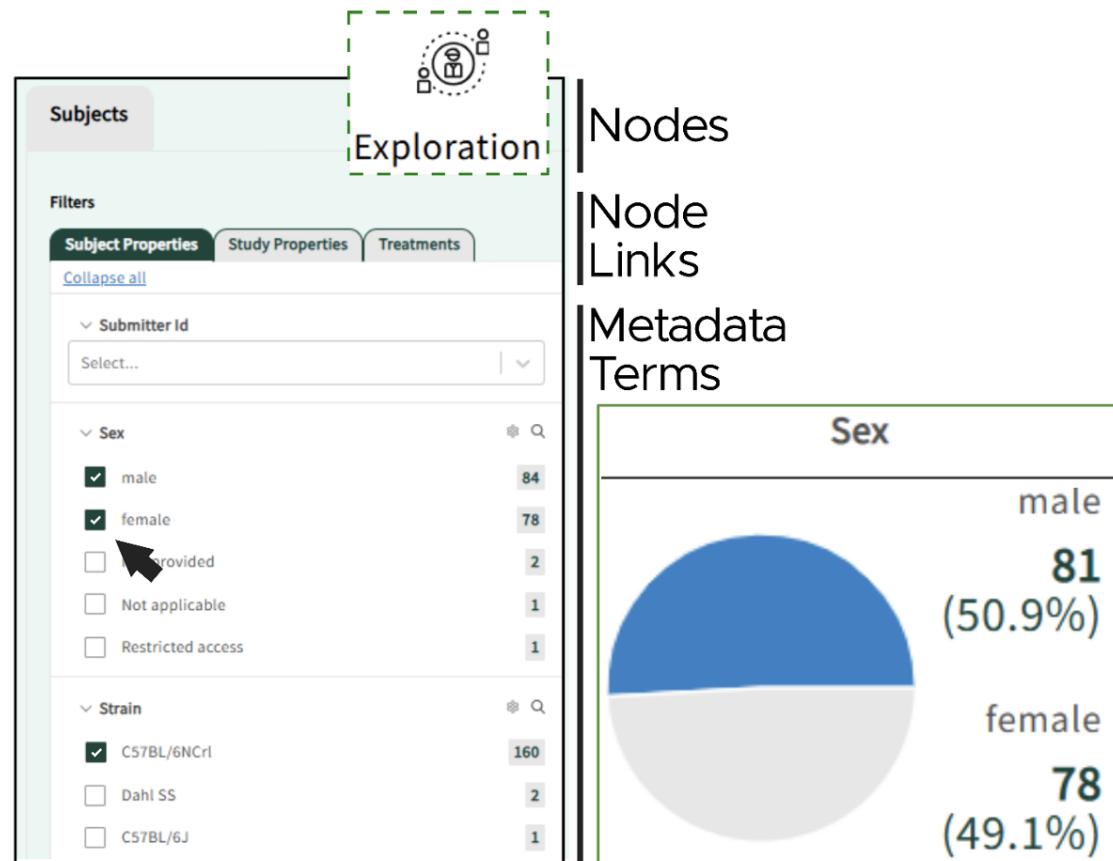
Our data model is evolving, we need user input for prioritization.

*In progress:*

- Addition of pup – dam/sire relationships
- Integrating a preliminary environmental sampling module
- Integrating a community survey module



# IMPACT OF ROBUST METADATA



Data can now be explored in a very structured manner knowing exactly what the '*rules*' are.

No need to search multiple terms:  
 ("TCDD" OR "Dioxin" OR "2,3,7,8-tetrachlorodibenzo-p-dioxin") AND ("liver" OR "hepatic" OR "hepatocyte")  
 ...

Both a user interface and API can be used to find these data.

# BUILDING COHORTS OF DATASETS

A citable persistent ID (PID) is assigned to each dataset

Project Id	Source Node	study_id	File Name	File Size	GUID
TRAINING-training001	slide_image	study_test001	PPGTEST_PPG001_MAP.txt	5.66 KB	<a href="#">dg.TDC/25d854ec-bfbd-4e85-be11-686f0d8d9d8b</a>
TRAINING-training001	slide_image	study_test001	test.txt	9 B	<a href="#">dg.TDC/169cb934-9939-4157-b36a-e9b6a434af25</a>
TRAINING-training001	slide_image	study_test001	PRJ129_bodyweights.txt	12.58 KB	<a href="#">dg.TDC/54ad8253-dfa7-4eb5-a549-f0fe7ce21337</a>
TRAINING-training001	slide_image	study_test001	PRJ129_sampleweights.txt	5.15 KB	<a href="#">dg.TDC/98481fc5-d300-47c7-a0e3-ff86aa979a6</a>
TRAINING-training001	slide_image	study_test001	test.txt	9 B	<a href="#">dg.TDC/4f04343d-5bfa-46b6-b025-003284753d63</a>
TRAINING-training001	slide_image	study_test001	PPGTEST_PPG001_SBP.txt	5.6 KB	<a href="#">dg.TDC/93277b6a-f6eb-4a71-8843-b743d8671e64</a>

◀ ▶

Previous	Page	1	of 1	20 rows ▾	Next
----------	------	---	------	-----------	------

# WHY WE CREATED ToxDATACOMMONS

We found a gap between existing infrastructure and the FAIR sharing of *all* research data.

Existing resources are siloed. Mapping expression data to metabolomic data, imaging data, weight data... was extremely difficult.

Gen3 commons showed a lot of promise

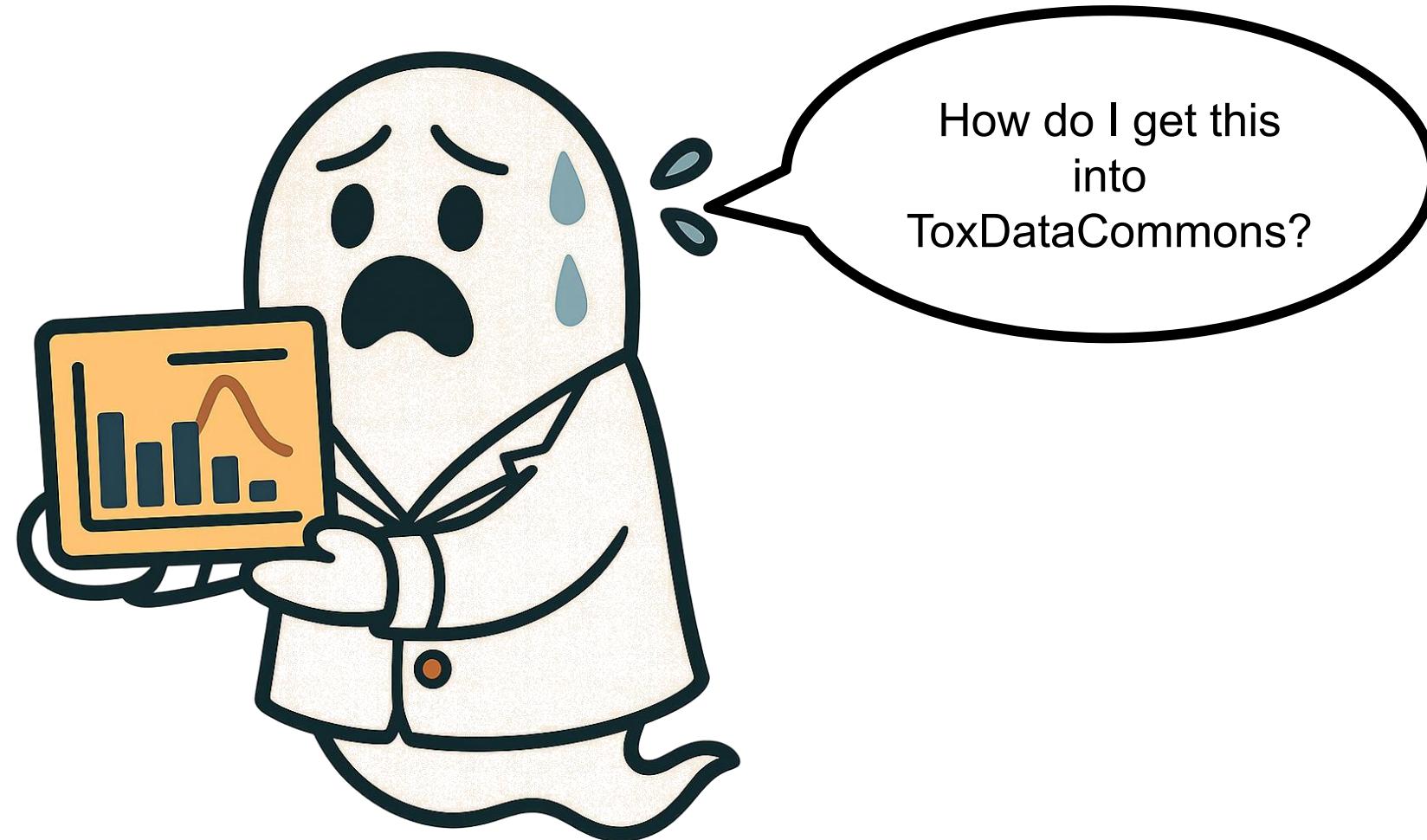


# PROBLEM SOLVED!



*Image created with ChatGPT (DALL·E), 2025.*

# PROBLEM SOLVED!



*Image created with ChatGPT (DALL·E), 2025.*

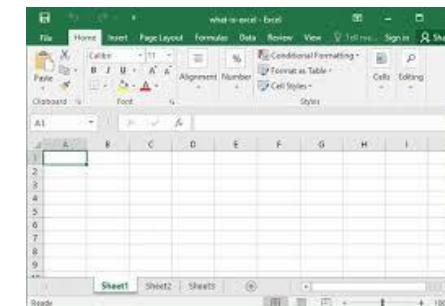
# MAKING METADATA COLLECTION ACCESSIBLE

The research community has aimed to make metadata collection more accessible for years through ***Templatization***:

- Who creates the templates?
- Who has access to the templates?
- Do we use web forms or spreadsheets?
- ...



cidgoh/  
DataHarmonizer

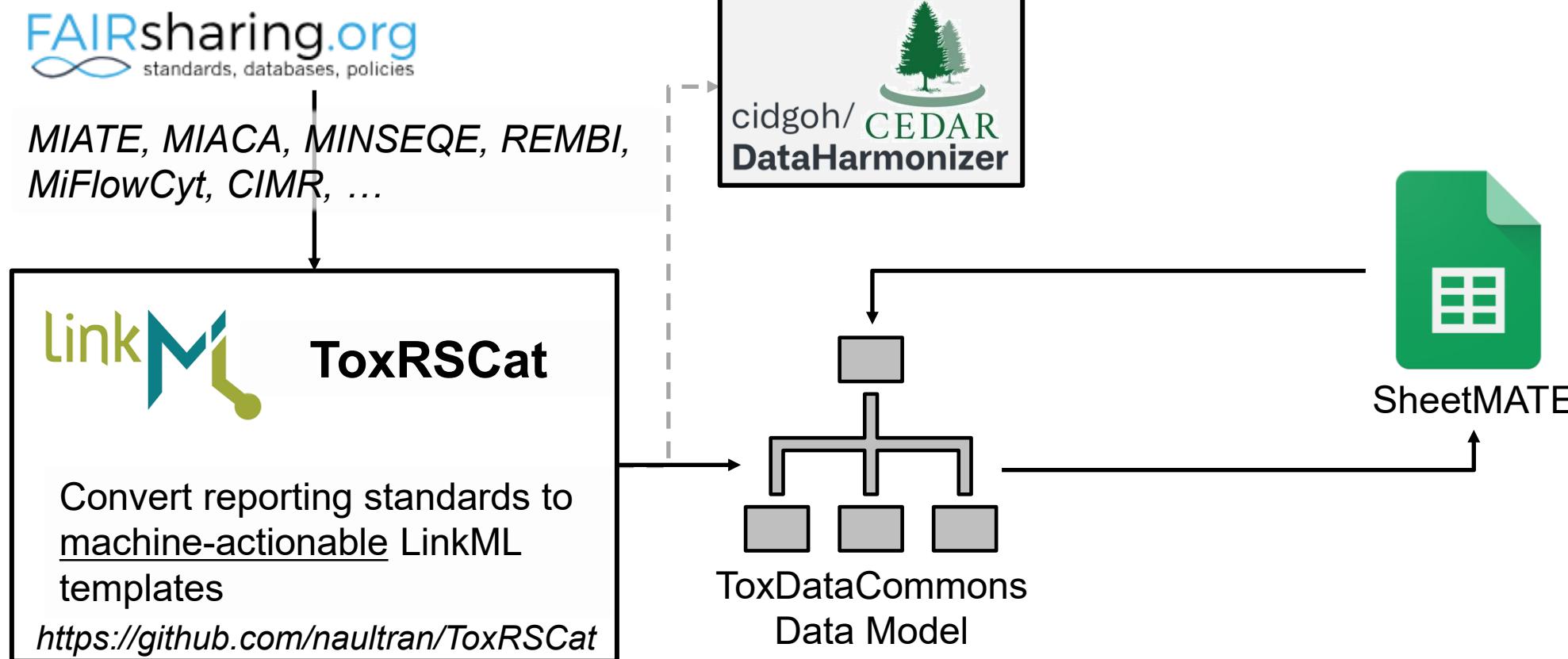


# MAKING METADATA COLLECTION ACCESSIBLE

The research community has aimed to make metadata collection more accessible for years through ***Templatization***:

- Who creates the templates?
- Who has access to the templates?
- Do we use web forms or spreadsheets?
- ...
- **Domain experts**
- **Everyone/Anyone**
- **Spreadsheets**
- **Machine-actionable**

# BUILDING A COMMUNITY RESOURCE

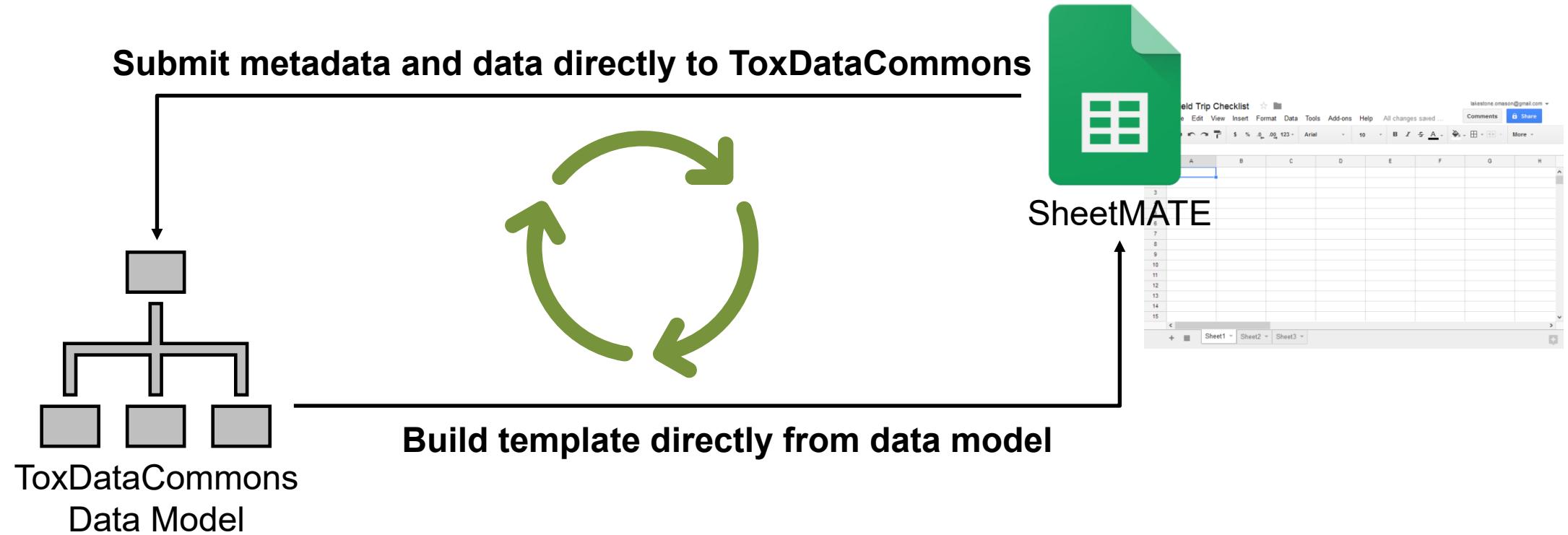


# WHY SHEETMATE?

SheetMATE was primarily driven by user experience:

- **ISACreator:** Tool was retired and no longer maintained. Template format was machine-actionable but not widely used.
- **DataHarmonizer:** Builds templates from LinkML but the spreadsheet functionality often felt limiting compared to Excel/Sheets.
- **CEDAR Workbench:** Web form format. Templates built and shared by users (can vary a lot). Difficult to do batch collections.
- None of the tools could be directly submitted to ToxDataCommons

# SHEETMATE: CLOSING THE LOOP





# ToxDataCommons IN ACTION: DATA UPLOAD

ToxDataCommons holds all the metadata. Can we use it to improve data submissions to established repositories?

## Without using ToxDataCommons

Treatment ID:	TR003988
Treatment Summary:	TCDD elicits similar effects as TCDD in humans (heightened sensitivity to right-shifted). To observe this effect, day (PND) 25 C57BL/6N mice were administered TCDD (Innovive Innocages (San Diego, CA) and a 12 h/12 light/dark cycle (Sigma-Aldrich, St. Louis, MO). A total of 7 mice were administered TCDD at different doses and study duration as opposite to the human studies. The doses and study duration are as follows: 0, 1, 3, 10, or 30 micromolar TCDD in sesame oil (0.1 milliliters) for 28 days.
Treatment Compound:	2,3,7,8-tetrachlorodibenzo-p-dioxin
Treatment Route:	Oral gavage
Treatment Dose:	0, 1, 3, 10, or 30 micromolar
Treatment Dosevolume:	0.1 milliliters
Treatment Doseduration:	28 days
Treatment Vehicle:	sesame oil

## Using ToxDataCommons

Treatment ID:	TR003667
Treatment Summary:	Mice were orally gavaged with TCDD at different doses and study duration as follows: 0, 0.01, 0.03, or 0.1 mL TCDD in sesame oil (ad libitum) for 28 days.
Treatment Compound:	DTXSID2021315:1
Treatment Route:	Oral Gavage Route
Treatment Dose:	[0.0, '0.01', '0.03', 0.1]
Treatment Dosevolume:	0.1 mL
Treatment Vehicle:	DTXSID9033971:3
Animal Fasting:	ad libitum
Animal Endp Euthanasia:	Carbon dioxide asphyxiation
Animal Endp Tissue Coll List:	Liver
Treatment ID:	TR002872
Treatment Summary:	Mice were orally gavaged with TCDD at different doses and study duration as follows: 0, 0.03, 0.1, or 0.3 mL TCDD in sesame oil (ad libitum) for 28 days.
Treatment Compound:	DTXSID4051378:1
Treatment Route:	Oral Gavage Route
Treatment Dose:	[0.0, '0.03', '0.1', '0.3']
Treatment Dosevolume:	0.1 mL
Treatment Vehicle:	DTXSID9033971:3
Animal Fasting:	ad libitum
Animal Endp Euthanasia:	Carbon dioxide asphyxiation
Animal Endp Tissue Coll List:	Liver

# ToxDATACOMMONS IN ACTION: INTEGRATION

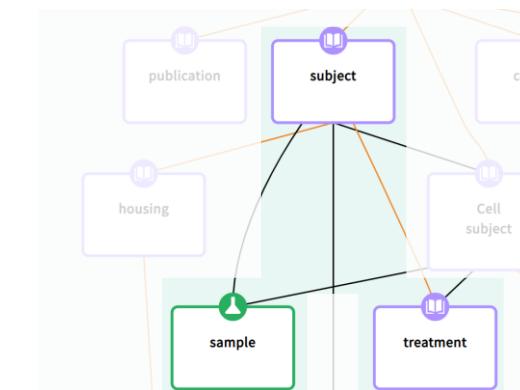
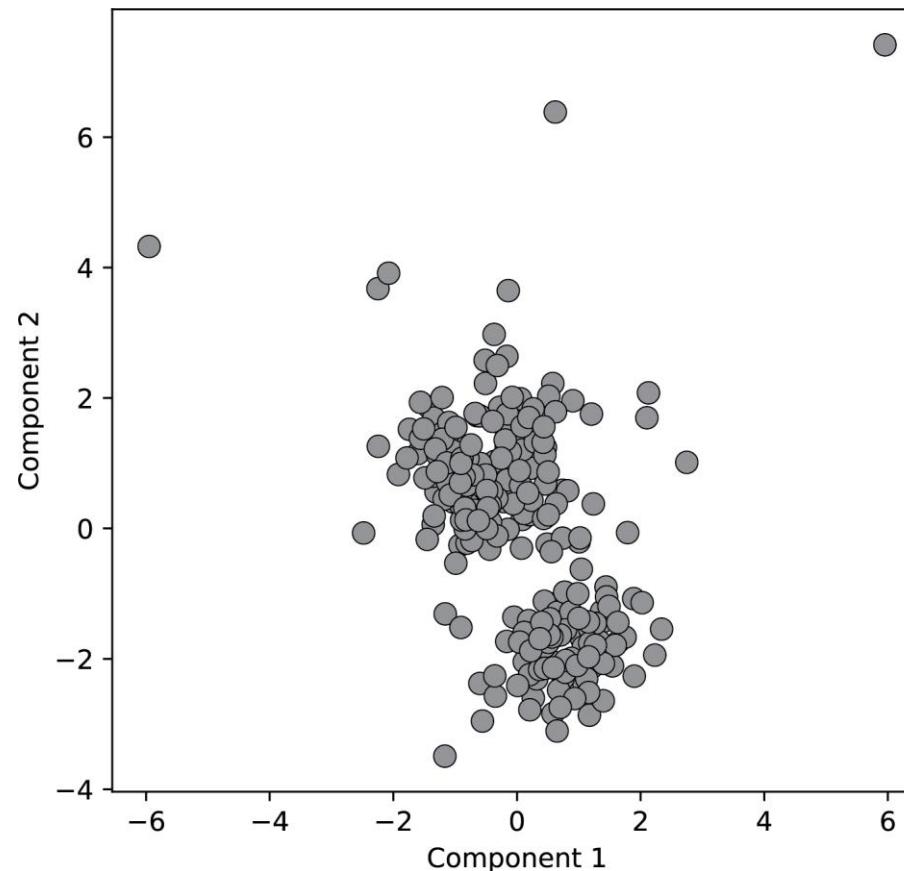
Datasets were generated in at least 3 independent studies with several measured endpoints (weights, pathology, gene expression, ...).

	Time									
Dose (Chemical X)	1	2	3	4	5	6	7	8	9	10
1	■									
2		■								
3			■							
4				■						
5					■					
6						■				
7							■			
8								■		
9									■	
10										■

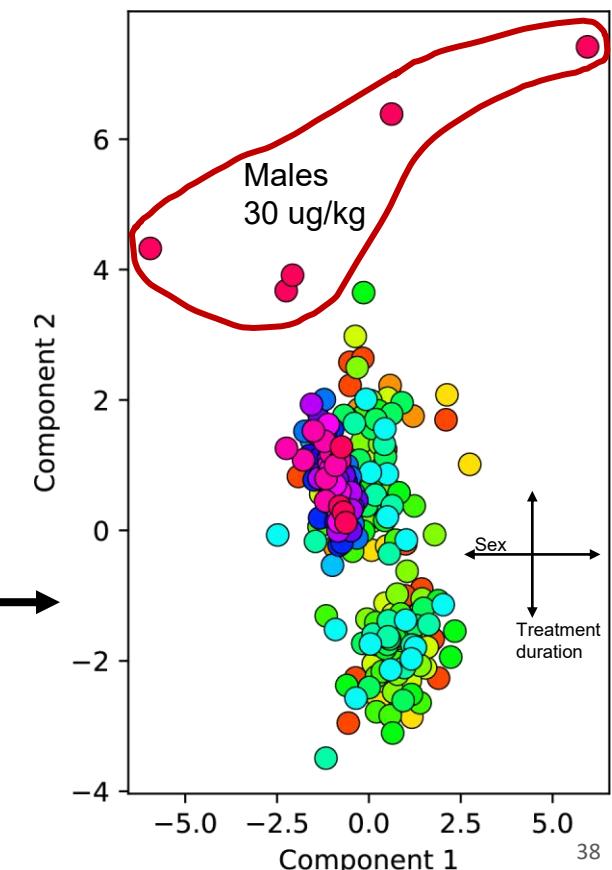
**Can we identify informative patterns about chemical toxicity by integrating these studies?**

# ToxDATACOMMONS IN ACTION: INTEGRATION

Dynamic time warping analysis of daily body weights of each *subject*.

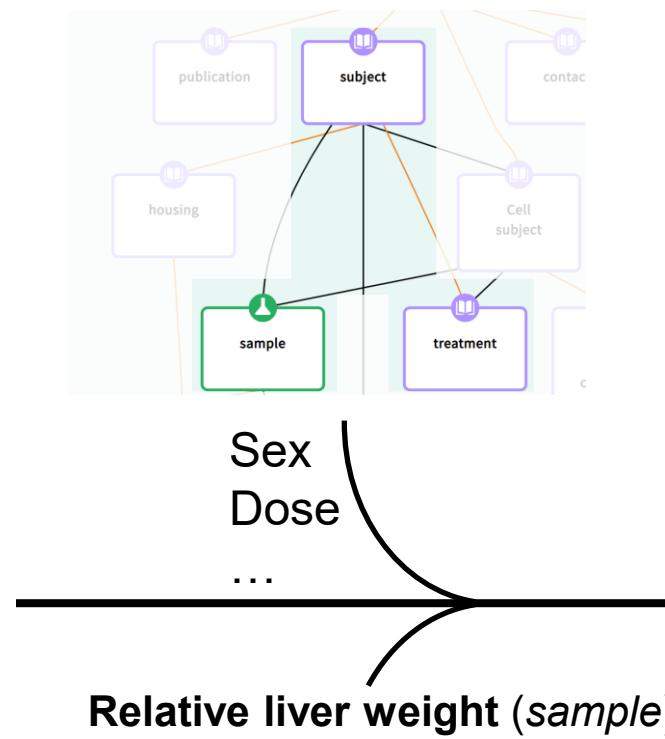
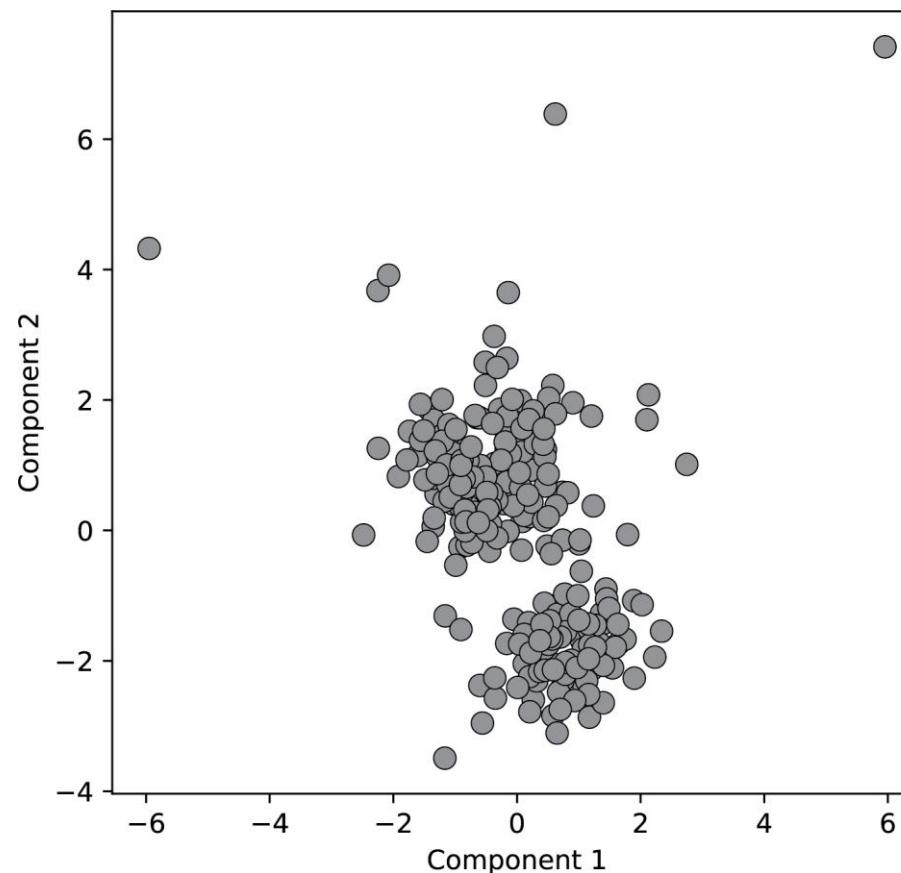


Metadata annotation reveals clustering

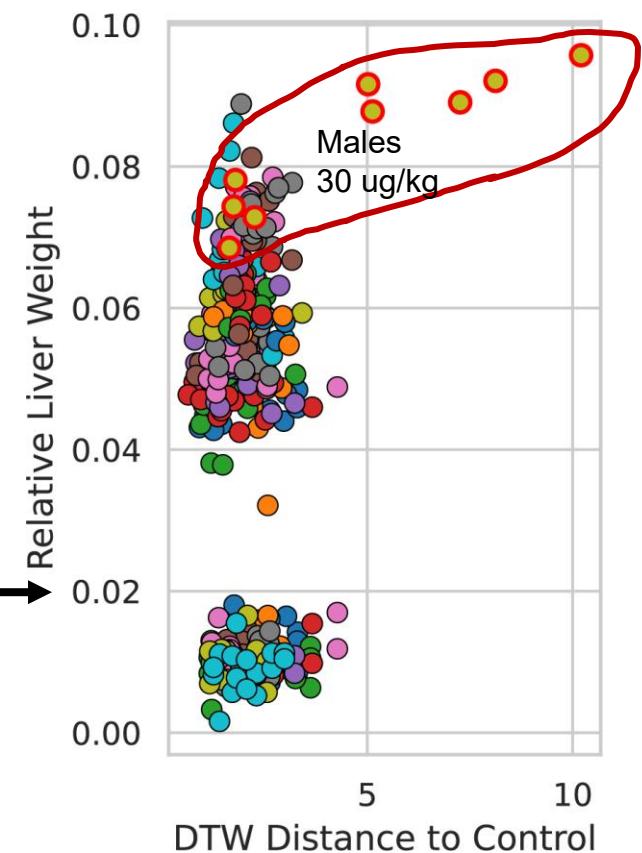


# ToxDATACOMMONS IN ACTION: INTEGRATION

Dynamic Time Warping analysis of daily body weights of each *subject*.



Metadata annotation reveals clustering



# SUMMARY

- Improved data sharing has the potential to benefit the entire environmental health research community
- The MSU SRC has tackled this by building infrastructure and resources. *It is only a single piece of the puzzle!*
- **Collaboration**, feedback, ideas, and transparency regarding challenges will help us move forward more effectively.



# ToxDATACOMMONS VISION

## MSU Superfund Research Program

30+ Investigators across  
many departments and  
colleges

## MSU EHS Researchers

500+ Investigators across  
many departments and  
colleges

## External Data

30+ Superfund Research  
Centers and even more  
EHS researchers

2025

2026

2028

# ACKNOWLEDGEMENTS

## MSU Superfund Research Center

\*Keji Yuan

Tim Zacharewski

Giovan Cholico

Eric Kasten

Jonathan Babbage

Todd Hall

## Funding:

NIEHS SRP P42 ES004911

## CTDS:

Chris Meyer

Ed Malinowski

Jawad Qureshi



National Institute of  
Environmental Health Sciences  
*Superfund Research Program*

