ecetoc

October 16 2024 https://www.ecetoc.org/event/ ai-workshop/

# FAIR Data and data that are Fully AI Ready

Erik Schultes http://orcid.org/0000-0001-8888-635X eriks@gofair.foundation These slides: <u>https://osf.io/ueqvw</u>





ecetoc

# Opportunities

### **Al-assisted Next Generation Risk Assessment** enabled by FAIR Data and Knowledge Graphs

Barry Hardyl, Asmaa Alil, Tomaz Mohoric1, Daniel Burgwinkel1, Jeff Wiseman1, Thomas Luechtefeld2, Zaki Mughal2 [1] EdelweissConnect GmbH, Switzerland, [2] InSilica, Bethesda, MD 20817, US

#### BACKGROUND

#### **KNOWLEDGE GRAPH CREATION**

**AI-ASSISTED APPLICATION** 

As an initial AI-Knowledge Graph integration,

🖋 WikiPathways Query Tool

Configuring BioBricks & OpenA

**DESIGN FOR AI SAFETY (CBRN)** 

We are designing our toxicological safety risk assessment system to

comply with AI safety regulations, such as the NIST AI Risk Management

Framework and the EU AI Act. The providers of Large Language Models

(LLMs) have to prevent the misuse in context of CBRN threats (Chemical,

Biological, Radiological, and Nuclear). We analyse how these AI safety

measures impact our use cases for safety risk assessments.

we used the OpenAI API and WikiPathways'

RDF Knowledge Graph for information

retrieval via natural language queries.

The open knowledge resources offer a harmonized data

focusing on harmonization, referencing, and integrity

best practices, and a standardized data collection approach

iteatosis AOP network

🖋 WikiPathways Query Tool

#### We developed two toxicology knowledge resources for case studies:

- A dataset and knowledge graph from the EU-ToxRisk[2] program on New Approach Methods in Toxicology. - A knowledge graph based on a network model for steatosis, supporting ASPIS[3] case studies and SSbD4CheM.



#### **TESTING AND REFINING**

#### Our AI-enabled risk assessment test cases include:

1. Efficient retrieval of reliable, well-referenced knowledge via a knowledge graph;

2. Enhanced toxicology learning through ML and Generative AI models:

3. Al integration into risk assessment workflows for analysis, evidence interpretation, and decision-making;

4. Generation of auditable reports documenting risk assessment

processes with supporting data and knowledge records.

#### NEXT STEPS AND COLLABORATION

#### Collaboration opportunities:

- sharing open science resources supporting NAMS-based risk

- assessment with the regulatory acceptance purpose.
- Establishing FAIR practices and approaches to sustainability.



#### References

[1] Findable, Accessible, Interoperable, Reusable (FAIR)

[2] EU-ToxRisk, An Integrated European 'Flagship' Programme Driving Mechanism-based Toxicity Testing and Risk Assessment for the 21st century [3] Animal-free Safety assessment of chemicals: Project cluster for Implementation of novel Strategies (ASPIS);



his project has received funding from the European Union's Horizon Europe research and innovation programme under grant agreement n° 101138475. UK participants in SSbD4CheM project are supported by UKRI. CH participants in SSbD4CheM project receive ding from the Swiss State Secretariat for Education, Research and Innovation (SERI).



Insilica is the proud recipient of an NSF SBIR and NIH SBIR grant. ToxTrack, Inc is a subsidiary of Insilica

SUMMARY

#### . onfederaziun svizra

Project funded by

tate Secretariat for Educ search and Innovation SERI

Our AI-assisted knowledge from FAIR data is integral to NGRA workflows, including ASPA Workflow[3], template aligned with FAIR principles, omics data management supporting stages like problem formulation, hazard characterization, mechanistic interpretation, and

Our goal is to develop AI-assisted Next Generation Risk

biological response data, pathways and key events.

Assessment. In our first project phase we are developing processes and applications to prepare FAIR[1] datasets and knowledge graphs connecting compound information,



#### **QUALITY AND TRUST ISSUES**

We design systematic, automated workflows for data processing and utilize the biobricks.ai platform to transform our FAIR datasets into FAIR knowledge graphs using the semantic web's RDF format.







#### NGRA WORKFLOW INTEGRATION

# FAIR Principles & AI?

### **Key Questions:**

- 1. What are the essential elements of data curation and formatting for AI readiness? (What type of data curation and formatting is necessary to ensure that AI is fed with all relevant data?)
- 2. How can we ensure that all relevant data, including minor parameters, are accurately captured and utilized?
- 3. Which standards or ontologies should be adopted to improve data interoperability and machine readability?
- 4. How AI could potentially be used to promote FAIR principles in chemical safety assessment?
- 5. What measures can be taken to identify and mitigate biases in AI models resulting from data selection or curation?
- 6. How can we manage the lifecycle of AI models, including updates and decommissioning, in response to changes in data?

https://www.ecetoc.org/wp-content/uploads/2024/10/AI-WS\_Programme\_FINAL\_web.pdf



# scientific data

Explore content v About the journal v Publish with us v

nature > scientific data > comment > article

#### 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 Data3, Article number: 160018 (2016)Cite this article488kAccesses4555Citations2031AltmetricMetrics

...the FAIR Principles put specific emphasis on **enhancing the ability of machines to automatically find and use the data**, in addition to supporting its reuse by individuals.



### 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.
- 12. (meta)data use vocabularies that follow FAIR principles
- 13. (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



## **Box 2** | The FAIR Guiding Principles

### Computability

• FAIR data use knowledge representation languages and controlled vocabularies that reduce/eliminate ambiguity.

### Trustworthiness

• FAIR data have (rich) provenance which provides evidence for the source of data. Provenance can include uncertainties and error estimation. Trust also requires large amounts of metadata.

### Equitabilty

• FAIR data make explicit the conditions for reuse. Equitability can be technologically ensured, when data and metadata are FAIR



ecetoc

# Challenges

Key Considerations:

- **Filtering:** The data collected from these sources are heavily filtered to remove noise, low-quality text, and duplicates.
- Bias and Ethical Concerns: Training on public data can introduce biases present in the original content (e.g., social biases, regional biases), which is why researchers try to apply techniques to mitigate these issues.
- **Copyrighted Material:** While efforts are made to use publicly available content, there are ongoing debates and concerns about training on data that may come from copyrighted works, especially when it is scraped without explicit permission.

Content generated using OpenAl's ChatGPT-4

![](_page_10_Picture_5.jpeg)

FAIR Considerations:	Key Considerations:
Computability —	<ul> <li>Filtering: The data collected from these sources are heavily filtered to remove noise, low-quality text, and duplicates.</li> </ul>
Trustworthiness—	• <b>Bias and Ethical Concerns:</b> Training on public data can introduce biases present in the original content (e.g., social biases, regional biases), which is why researchers try to apply techniques to mitigate these issues.
Equitability	<ul> <li>Copyrighted Material: While efforts are made to use publicly available content, there are ongoing debates and concerns about training on data that may come from copyrighted works, especially when it is scraped without explicit permission.</li> </ul>
	Content generated using OpenAI's ChatGPT-4

![](_page_11_Picture_1.jpeg)

ecetoc

# Path forward...

![](_page_13_Figure_0.jpeg)

![](_page_14_Figure_0.jpeg)

![](_page_15_Figure_0.jpeg)

![](_page_16_Figure_0.jpeg)

![](_page_17_Figure_0.jpeg)

![](_page_18_Picture_0.jpeg)

![](_page_18_Picture_2.jpeg)

2024

FAIR as in "AI-Ready"

See also: FAIR in

ML, AI Readiness,

& Reproducibility (FARR) Workshop

October 9-10 https://www.farr-

rcn.org/ workshop24

![](_page_18_Picture_5.jpeg)

![](_page_18_Picture_6.jpeg)

 $\bigcirc$ Amazon Q

Google Al

![](_page_18_Picture_9.jpeg)

CCETOC

October 16 2024 https://www.ecetoc.org/event/ ai-workshop/

# **FAIR well**

![](_page_19_Picture_4.jpeg)

Erik Schultes http://orcid.org/0000-0001-8888-635X eriks@gofair.foundation These slides: <u>https://osf.io/ueqvw</u>

![](_page_19_Picture_6.jpeg)

![](_page_19_Picture_7.jpeg)