



PSDI

PHYSICAL SCIENCES
DATA INFRASTRUCTURE

We don't talk about Semantic Web Technologies

ACS Fall 2023

13th August 2023

Dr Samantha Kanza

University of Southampton

<https://www.psd.ac.uk/>

About Me



- ▶ Senior Enterprise Fellow at University of Southampton
- ▶ Pathfinder Lead & Researcher for PSDI Project: Process Recording
- ▶ Coordinates AI4SD & Future Blood Testing Network
- ▶ Research Interests: Semantic Web Technologies, IoT, Research Data Management, Digitisation, Lab of the Future, Paperless Labs, Re-use of Technology
- ▶ @SamiKanza

Lets talk about the Semantic Web

Conclusions



- **The Semantic Web is coming!**
 - ◆ Joint development between DARPA/EU/and W3C communities
 - ◆ Languages and tools are available to play with
 - [Http://www.daml.org/](http://www.daml.org/)
 - ◆ W3C interest group available for those wishing to join the discussion
 - Www-rdf-logic@w3c.org (live or archived)
 - ◆ Ongoing DoD and commercial projects
- **Come join us**
 - ◆ Submit ontologies/marked up pages
 - ◆ Develop tools or help test ours
- **Get in on the next big thing early!**

www.daml.org

<https://www.wi-consortium.org/wicweb/pdf/wi-hendler.pdf>

Common Misconceptions

The concept of machine-understandable documents does not imply some **magical artificial intelligence** allowing machines to comprehend human mumblings. It relies solely on the machine's ability to solve well-defined problems by performing well-defined operations on well-defined data. So, instead of asking machines to understand people's language, the new technology, like the old, involves asking people to make some extra effort, in repayment for which they get major new functionality – just as the extra effort of producing HTML mark-up is outweighed by the benefit of having content searchable on the web.

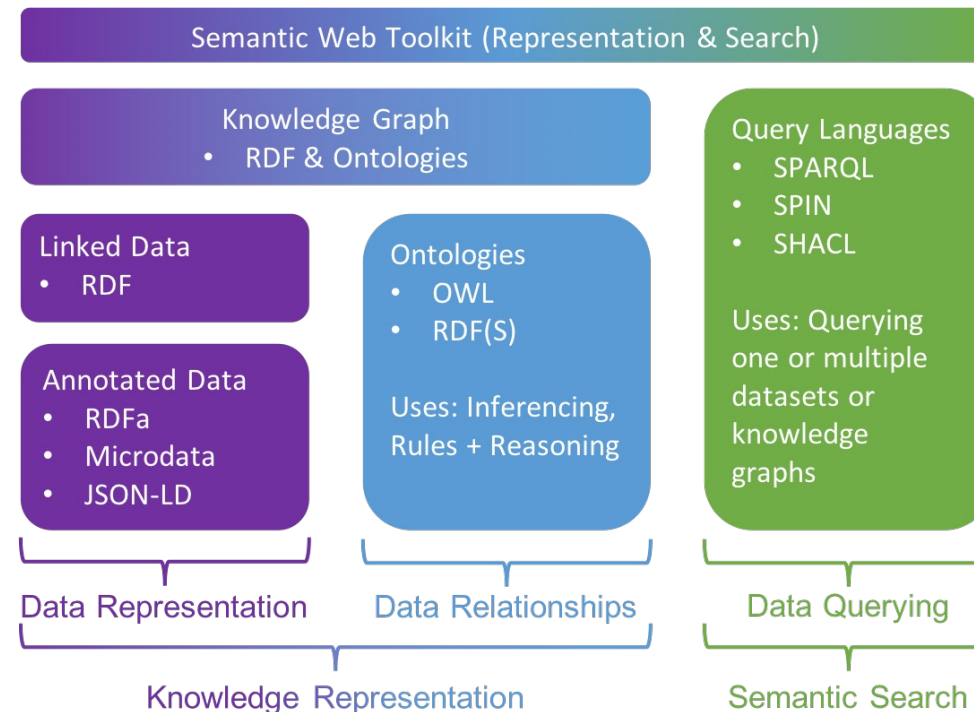
Berners-Lee, T. and Hendler, J., 2001. Publishing on the semantic web. *Nature*, 410(6832), pp.1023-1024.

Berners-Lee, T., Hendler, J. and Lassila, O., 2001. The semantic web. *Scientific american*, 284(5), pp.34-43.

The Semantic Web will enable machines to COMPREHEND semantic documents and data, not human speech and writings.

So what is the Semantic Web?

- ▶ The Web of Linked Data
- ▶ Creating machine-readable/understandable data
- ▶ A way to bring context and meaning to data
- ▶ A set of common standards for data representation, integration, and search



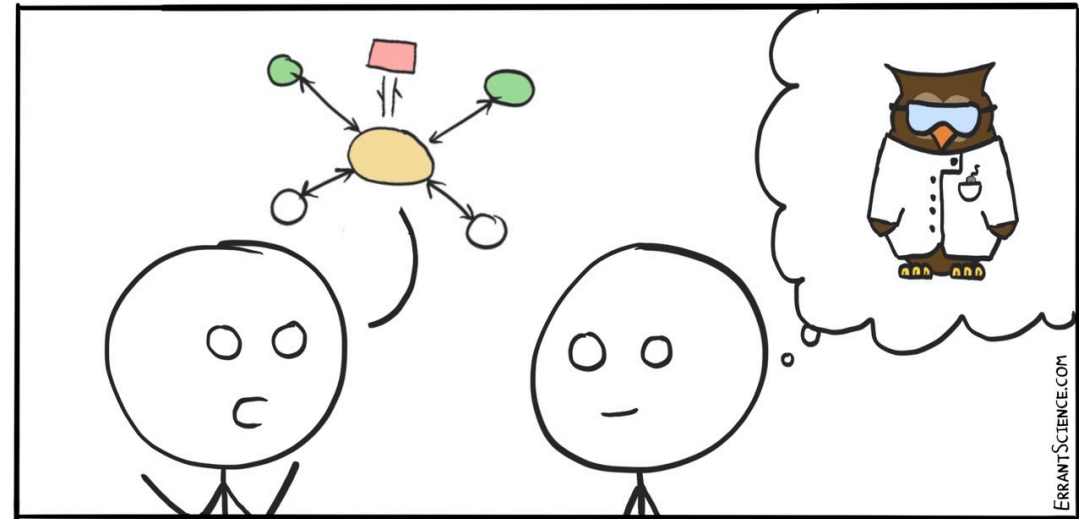
So where did this all go wrong?



<https://www.pinterest.co.uk/pin/539306124105554567/>

Ontologies

- ▶ Too many ontologies (and yet not enough?)
- ▶ “Lets just make a new one” mentality
- ▶ Lack of standards for Ontologies (aside from the biology sphere)
- ▶ Ontology projects frequently don't consider the full data lifecycle
- ▶ Ontologies aren't always well maintained
- ▶ Ontologies aren't always FAIR
- ▶ Ontology tools are poor or expensive



Cartoon created by ErrantScience.com for AI4SD: licensed under [CC-BY-NC](https://creativecommons.org/licenses/by-nc/4.0/)

Linked Data Creation/Conversion

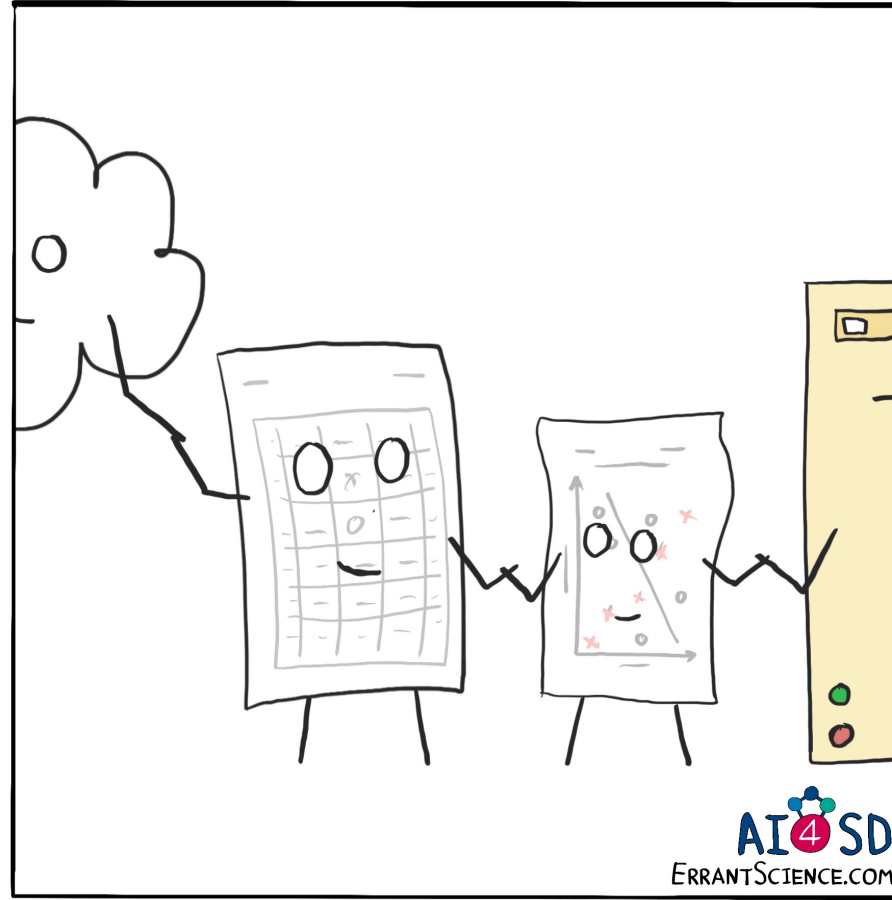
- ▶ Tools to convert structured data into linked data are still lacking
- ▶ R2RML even with additional libraries is repetitive
- ▶ Existing tools (OpenRefine or Karma) are very clunky and require a lot of human editing
- ▶ Writing custom scripts seems to be the best current method but isn't very user friendly



<https://imgflip.com/memegenerator/>

Data Storage & Data Types

- ▶ Some data doesn't lend itself well to being represented semantically (e.g. numerical classifications)
- ▶ Large datasets lead to even larger knowledge graphs
- ▶ Issues with reading, writing and deleting data with certain (free) versions of triple stores



Cartoon created by ErrantScience.com for AI4SD: licensed under [CC-BY-NC](https://creativecommons.org/licenses/by-nc/4.0/)

Given these issues, why should we care?

- ▶ Defining Common Shared Vocabularies
- ▶ Machine Readable Data
- ▶ Interoperable Metadata
- ▶ Linking Datasets
- ▶ Inferencing
- ▶ Semantic Search
- ▶ Unlocking the potential of AI/ML



<https://imgflip.com/i/4x69cr>

Common Shared Vocabularies

- ▶ Define a common set of terms to describe scientific data concepts consistently
- ▶ Markup your datasets to talk about the **SAME** concepts



Inconsistent
data
terminology

Using
an ontology

<https://imgflip.com/i/7uq27y>

Machine Readable Data

- How can computers be expected to analyse and make decisions on scientific data they don't understand?



"My Computer doesn't understand me" by Etta Hulme is licensed under [CC-BY-NC](#)

Rich Interoperable Metadata

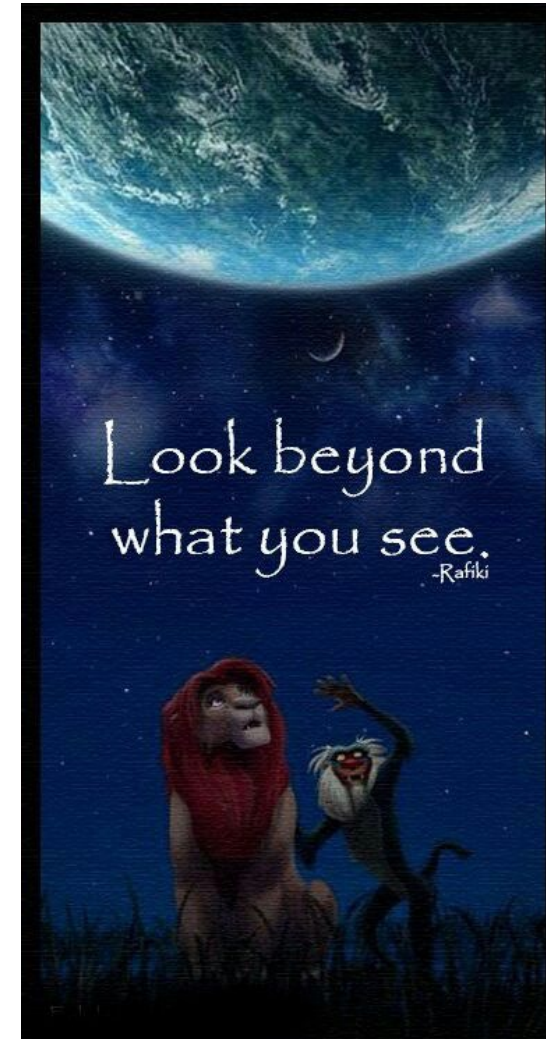
- ▶ Adding rich interoperable semantic metadata to your documents / webpages adds a whole new dimension to your data



https://www.pinterest.co.uk/jaci_mize/metadata/

Semantic Knowledgebases

- ▶ Scientific Research can be significantly enhanced by linking datasets together to create knowledgebases:
 - ▶ Find undiscovered links
 - ▶ Answer questions that cannot be addressed with a single data source
 - ▶ Very useful for domains like drug discovery which require the integration of disparate datasets



<https://www.pinterest.co.uk/pin/2674081014251214/>

Inferencing

- ▶ Description Logic can be embedded into Ontologies
- ▶ This enables machines to “infer” additional information that is not explicitly defined in the data

Example: If we know that:

`Sami isAllergicTo Juniper`
`Gin hasBotanical Junpier`



<https://quoteswell.com/cheshire-cat-quotes/>

Then we can infer: `Sami isAllergicTo Gin`

Semantic Search

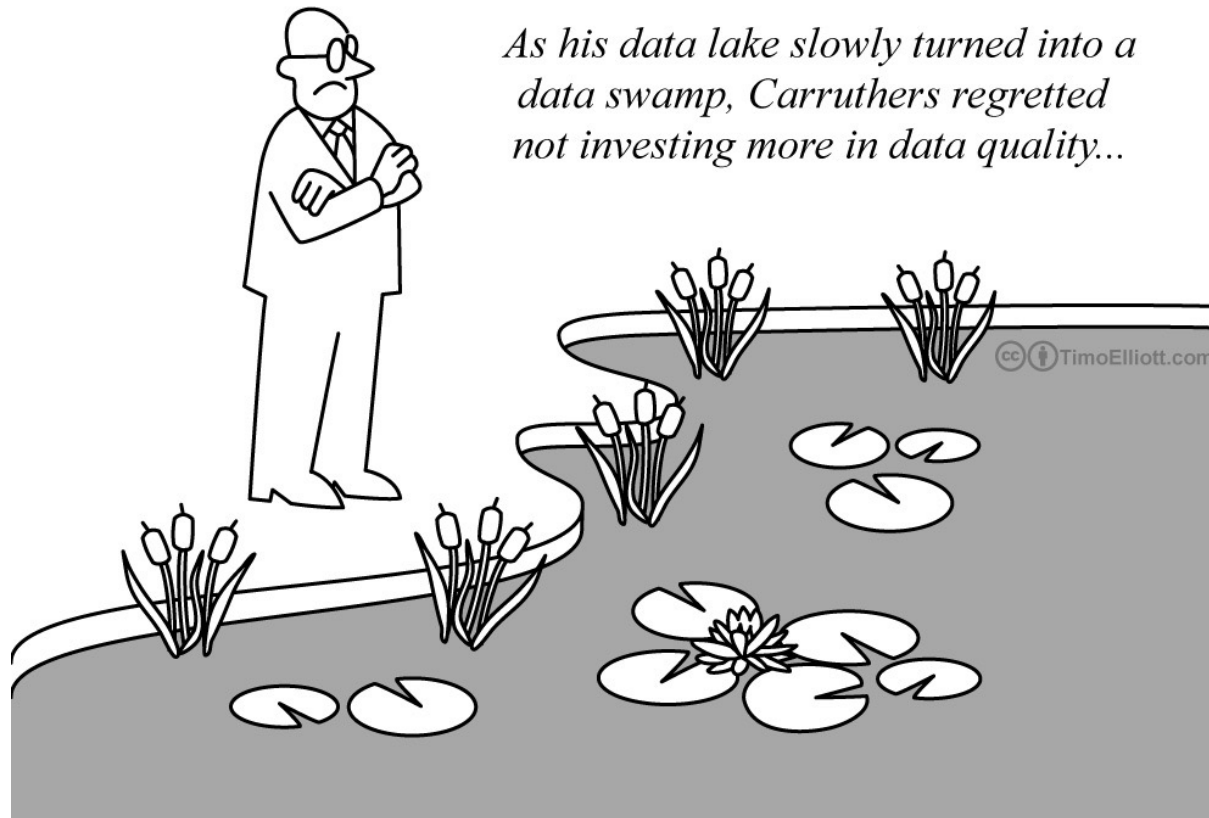
► Search on CONCEPTS across LINKED DATASETS



<https://garfield.com/comic?keywords=Jon&page=1113>

Unlock the Potential of AI & ML

- ▶ Don't waste your algorithms time with poorly formed inconsistent data



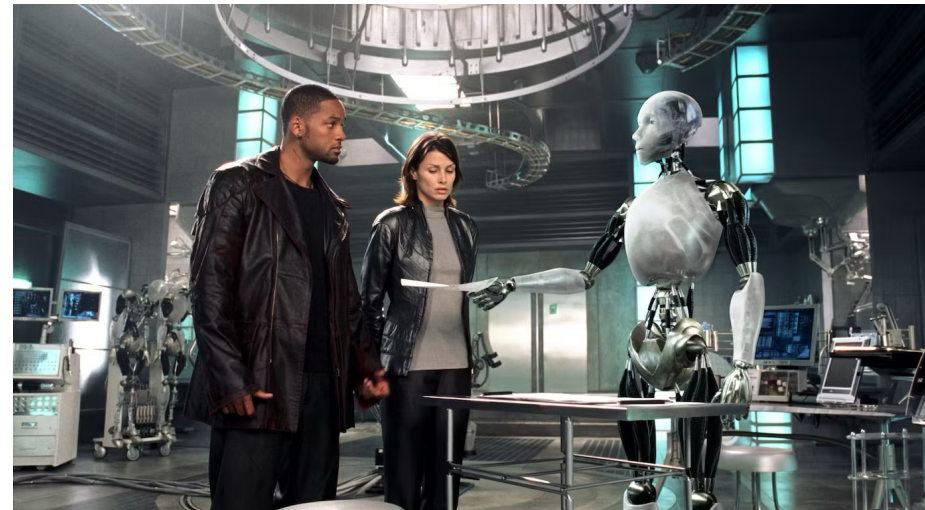
So what can we do?



<https://memegenerator.net/instance/68158697/hermione-leviosa-swish-flick>

You need to ask the right questions

- ▶ What is your use case for semantics?
- ▶ Do you need an ontology?
- ▶ Do you need a knowledge graph?
- ▶ How much data do you need to represent semantically?
- ▶ What data would benefit from being in a semantic format?
- ▶ What do you want your data to look like and how are you going to use it?



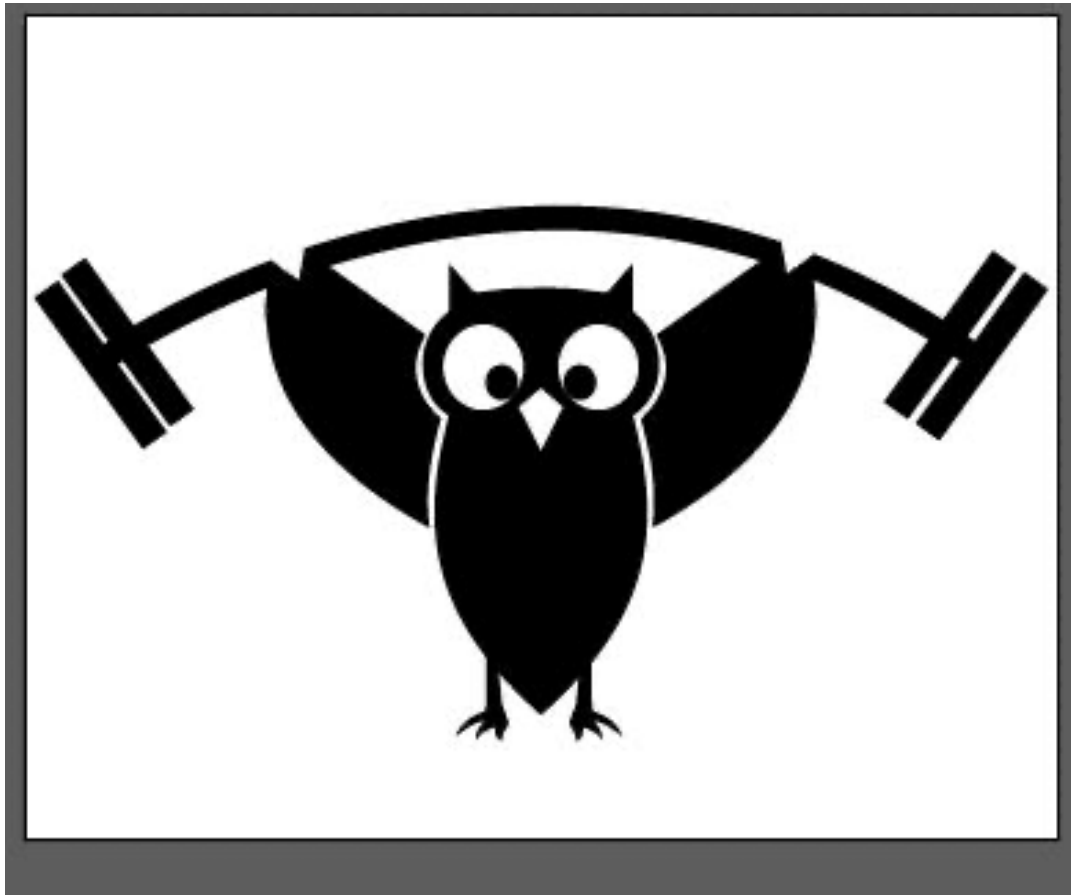
<https://www.inverse.com/culture/sci-fi-movies-january-2023-hbo-max-i-robot>

Consider your Conversion in Advance

- ▶ Converting data into linked data is non-trivial
- ▶ Conceptualising an Ontology vs using it in a dataset are two different things
- ▶ You need to consider the entire semantic data lifecycle of a project



Think about the WHY and the WHAT



<https://cuk00jan.wordpress.com/2014/02/14/who-is-alpha-owl/>

- ▶ Think about **WHY** you are making Ontologies and **WHAT** they are for
- ▶ Play to OWLS Strengths

Re-use and extend Ontologies where possible

- ▶ There are lots of ontologies out there! See if there's any you can re-use before making a new one (Search Papers/OLS/Google)
- ▶ Re-use doesn't just mean ontologies! You can re-use design patterns as well!
- ▶ Consider trying to create an extension to an existing ontology if appropriate



Modularize your Ontologies

- Break your ontologies into smaller related modules

Relation to time Granularity	Continuant				Occurrent	
	Independent		Dependent			
Complex of organisms	Family, community, deme, population		Environment	Organ function (FMP, CPRO)	Population phenotype	Population process
Organ and organism	Organism (NCBI Taxonomy)	(FMA, CARO)		Phenotypic Quality (PaTO)	Biological process (GO)	
Cell and cellular component	Cell (CL)	Cell component (FMA, GO)				
Molecule	Molecule (ChEBI, SO, RnaO, PrO)			Molecular function (GO)		Molecular process (GO)

Don't forget your standards

- ▶ Even your standards should have standards
- ▶ Consider:
 - ▶ Design patterns
 - ▶ Upper Level Ontologies
 - ▶ Consistency



FAIR is a Four Letter Word

"ALL RESEARCH SHOULD AIM
TO BE F.A.I.R."

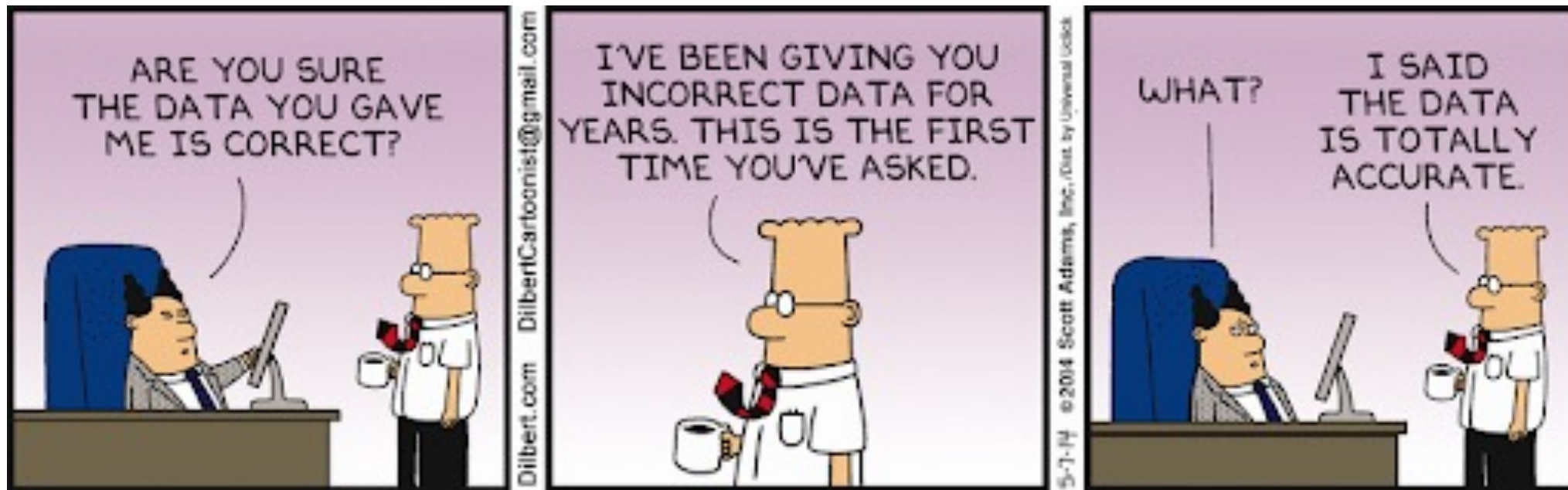
#FIGSHAREFEST

	GOOD	BAD
FINDABLE	ONLINE DATABASE	FILING CABINET IN A BATH IN THE BASEMENT UNDER A LEAKING PIPE
ACCESSABLE	OPEN ACCESS FOR EVERYONE (NO LOGIN)	THE FILING CABINET ALSO IS HOME TO A NEST OF WILD BADGERS
INTEROPERABLE	ALL DATA IS IN OPEN FORMATS	ALL DOCUMENTS ARE PRINTED IN COMIC SANS AND WRITTEN IN ESPERANTO
REUSEABLE	GOOD META DATA AND SECURELY STORED FOR 10 YEARS	THE PAPER EXPLODES IF IT'S READ

ERRANTSCIENCE.COM

Garbage In = Garbage Out

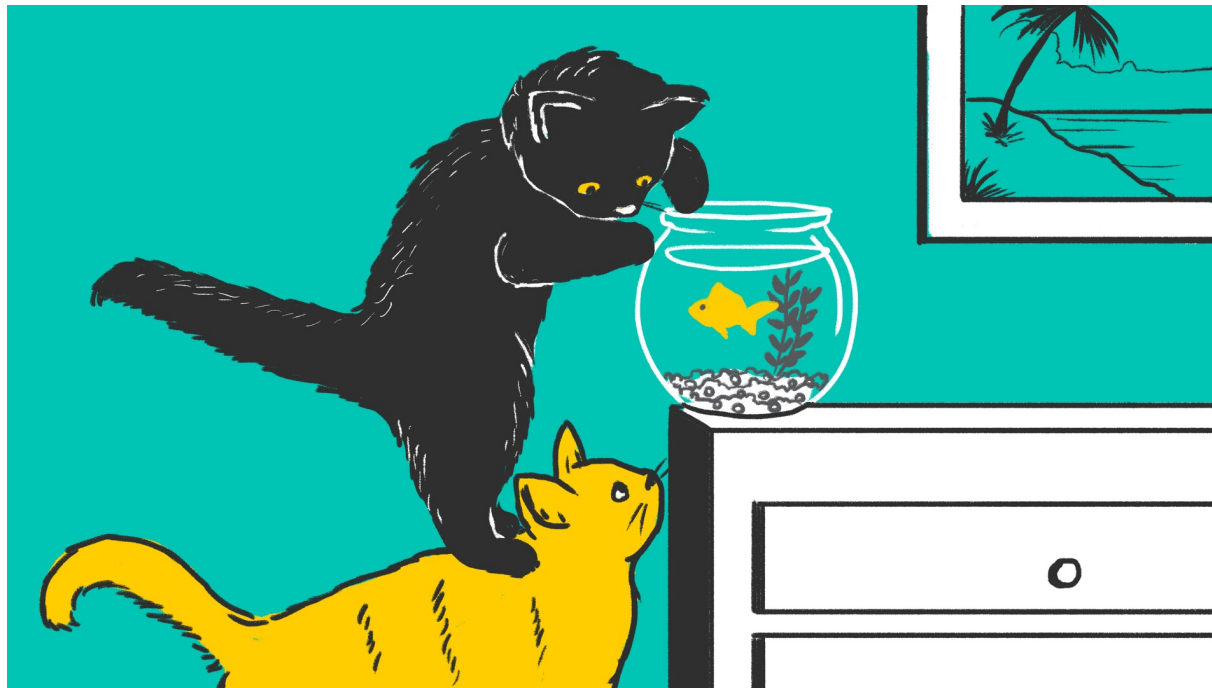
- GIGO can still apply even in the world of semantics



https://ffineis.github.io/blog/2017/09/13/case_for_big_govt.html

Conclusions

- ▶ We're all in this together and we need to work together!
- ▶ A lot of the technology is there, we just need to use it properly!
- ▶ This is as much a human endeavor as a technical one!



Relevant Publications & Talks

- ▶ Kanza, S., Willoughby, C., Gibbins, N., Whitby, R., Frey, J.G., Erjavec, J., Zupančič, K., Hren, M. and Kovač, K., 2017. Electronic lab notebooks: can they replace paper?. *Journal of cheminformatics*, 9(1), pp.1-15. <https://doi.org/10.1186/s13321-017-0221-3>
- ▶ Kanza, S., Stolz, A., Hepp, M. and Simperl, E., 2018. What does an ontology engineering community look like? A systematic analysis of the schema.org community. In *The Semantic Web: 15th International Conference, ESWC 2018, Heraklion, Crete, Greece, June 3–7, 2018, Proceedings 15* (pp. 335-350). Springer International Publishing. https://doi.org/10.1007/978-3-319-93417-4_22
- ▶ Kanza, S., 2018. What influence would a cloud based semantic laboratory notebook have on the digitisation and management of scientific research? (Doctoral dissertation, University of Southampton). <https://eprints.soton.ac.uk/421045/>
- ▶ Kanza, S., Gibbins, N. and Frey, J.G., 2019. Too many tags spoil the metadata: investigating the knowledge management of scientific research with semantic web technologies. *Journal of cheminformatics*, 11(1), p.23. <https://doi.org/10.1186/s13321-019-0345-8>
- ▶ Kanza, S. and Frey, J.G., 2019. A new wave of innovation in Semantic web tools for drug discovery. *Expert Opinion on Drug Discovery*, 14(5), pp.433-444. <https://doi.org/10.1186/s13321-019-0345-8>
- ▶ Kanza, S. and Frey, J.G., 2020. Semantic technologies in drug discovery—potential, practical, possibilities. <http://dx.doi.org/10.1016/B978-0-12-801238-3.11520-X>
- ▶ Kanza, S., 2021. AI3SD Video: Semantic Web in Scientific Research—Possibilities & Practices. <http://dx.doi.org/10.5258/SOTON/P0078>
- ▶ Kanza, S., Willoughby, C., Bird, C.L. and Frey, J.G., 2022. eScience Infrastructures in Physical Chemistry. *Annual Review of Physical Chemistry*, 73, pp.97-116. <https://doi.org/10.1146/annurev-physchem-082120-041521>
- ▶ Kanza, S., Willoughby, C., Knight, N.J., Bird, C.L., Frey, J.G. and Coles, S.J., 2023. Digital research environments: a requirements analysis. *Digital Discovery*. <https://doi.org/10.1039/D2DD00121G>

Acknowledgements

- **PhD Supervisors:** Jeremy Frey & Nicholas Gibbins (*University of Southampton*)
- **Electronic Lab Notebook Research:** Cerys Willoughby & Nicholas Gibbins & Richard Whitby & Jeremy Frey (*University of Southampton*) Jana Erjavec, Klemen Zupančič, Matjaz Hren & Kristina Kovač (*SciNote*)
- **PSDI Team:** Simon Coles, Jeremy Frey, Nicola Knight, Cerys Willoughby & Colin Bird (*University of Southampton*), Juan Bicarregui, Barbara Montanari, Brian Matthews, Vasily Bunakov (*Science & Technologies Facilities Council*)
- **Semantic Web Research Teams**
 - **SEED:** Pistoia Alliance, Pfizer, Biovia, Dotmatics, Sanofi, BenchSci, AstraZeneca, Elsevier, GSK, Bayer, Merck, SciByte, University of Southampton, Bristol-Myers Squibb, ChemAxon, IDBS, Linguamatics, Arxspan
 - **BSP:** Hugo Mills & Hannah Wickenden & Derek Scuffell (Syngenta), Jeremy Frey & Nicholas Gibbins (*University of Southampton*)
 - **S3W:** Susan Halford & Faranak Hardcastle (*University of Bristol*), Nicholas Gibbins & Mark Weal (*University of Southampton*), Cathy Pope (*University of Oxford*)

PSDI & Personal Details - Questions



www.psdi.ac.uk



@PSDI_UK



@PSDI_UK



[linkedin.com/company/psdiuk](https://www.linkedin.com/company/psdiuk)



linktr.ee/samanthakanza

Mailing List: <https://www.jiscmail.ac.uk/PSDI>