

Challenges on data and knowledge integration, analysis, life cycle and reproducibility in life sciences

Olivier Dameron

<https://orcid.org/0000-0001-8959-7189>

Université de Rennes

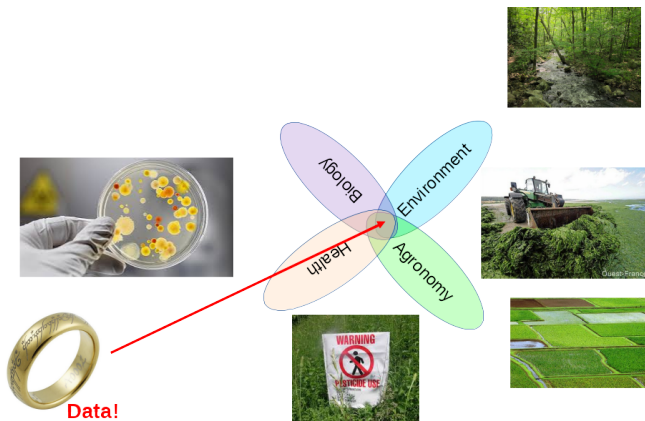
2025-10-15



CC BY-NC-SA

Life science data

Life science : a large, complex, inter-dependent domain...



... that stands out among other experimental sciences

“Biology has become an information science” [T. Lenoir, 1998 Stanford]

What to expect for 2025 ?

Our estimation is that genomics is a “four-headed beast” – it is either **on par with or the most demanding domain** [...] in terms of

- data acquisition
- data storage
- data distribution
- **data analysis**

Big Data: Astronomical or Genomical?

Zachary D. Stephens¹, Skylar Y. Lee¹, Faraz Faghri², Roy H. Campbell², Chengxiang Zhai³, Miles J. Efron⁴, Ravishankar Iyer¹, Michael C. Schatz^{5*}, Saurabh Sinha^{3*}, Gene E. Robinson^{6*}

PLOS Biology | DOI:10.1371/journal.pbio.1002195 July 7, 2015

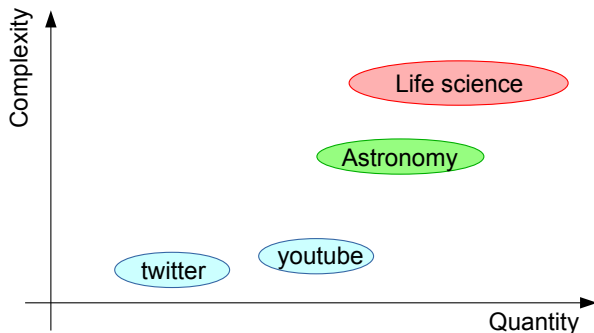
Table 1. Four domains of Big Data in 2025. In each of the four domains, the projected annual storage and computing needs are presented across the data lifecycle.

Data Phase	Astronomy	Twitter	YouTube	Genomics
Acquisition	25 zetta-bytes/year	0.5–15 billion tweets/year	500–900 million hours/year	1 zetta-bases/year
Storage	1 EB/year	1–17 PB/year	1–2 EB/year	2–40 EB/year
Analysis	In situ data reduction	Topic and sentiment mining	Limited requirements	Heterogeneous data and analysis
	Real-time processing	Metadata analysis		Variant calling, ~2 trillion central processing unit (CPU) hours
	Massive volumes			All-pairs genome alignments, ~10,000 trillion CPU hours
Distribution	Dedicated lines from antennae to server (600 TB/s)	Small units of distribution	Major component of modern user's bandwidth (10 MB/s)	Many small (10 MB/s) and fewer massive (10 TB/s) data movement

doi:10.1371/journal.pbio.1002195.t001

Degrees of data complexity

- multiple scales (heterogeneity)
- (highly) interdependent at each scale
- interdependent between scales
- variability
- incompleteness
- distributed (and lack of interoperability)
- evolving



Life science : beyond “data science” ...

- we have accumulated a trove of data
- we store and share these data
(in 2.000+ reference databases [Rigden2025])

> [Nucleic Acids Res.](#) 2025 Jan 6;53(D1):D1-D9. doi: 10.1093/nar/gkae1220.

The 2025 Nucleic Acids Research database issue and the online molecular biology database collection

Daniel J Rigden ¹, Xosé M Fernández ²

Affiliations + expand

PMID: 39658041 PMCID: [PMC11701706](#) DOI: [10.1093/nar/gkae1220](#)

Life science : a domain that stands out

- by its complexity
- by the scarcity of its unifying laws
- by its long history of knowledge description and formalization

... and toward “knowledge(-based) science” ?

We need a framework to support this transition

Attempt at defining some underlying notions

Database (e.g. Gene Expression Omnibus)

structured description of sets of homogeneous (as in “of the same class”) instances and the relations between them.

Ontology (e.g. GeneOntology, ChEBI, HPO)

formal description of the general concepts (as in “the classes of things”) of a domain and the relations between them.
(support general inferences).

Knowledge base (e.g. UniProt, Reactome, Rhea)

combines elements of databases and ontologies
(supports domain-based inferences about instances).

Knowledge graph (e.g. Reactome)

knowledge base that emphasizes graph-based capabilities.

Requirements for coping with life science data complexity

- **Requirement 1 : identify** resources with interoperable identifiers
- **Requirement 2 : describe** resources
 - ▶ their characteristics (e.g. start and end position of a gene,...)
 - ▶ their relations to other entities (e.g. the transcripts associated to a gene, the transcription factors that regulate it,...)
 - ▶ the categories they belong to
- **Requirement 3 : combine** descriptions from different origins, different points of view, different granularity levels
- **Requirement 4 : query** these descriptions
- **Requirement 5 : support semantically-rich** querying and reasoning (because of the inner complexity) using domain knowledge (this is required for capturing *expertise*)
- **Requirement 6 : cover the whole data life cycle**
- **Requirement 7 : enforce reproducibility**

If only the solutions to all these requirements were compatible !

The Semantic Web : a framework for
integrating seamlessly (data,) metadata and
knowledge, and querying and reasoning over it

Semantic Web and Linked (Open) Data

Semantic Web offers a unified framework to Linked Data

- **URIs** for identifying entities
- **RDF** for representing and aggregating entities descriptions
- **RDFS+OWL** for representing domain knowledge (and combine it with data descriptions)
- **SPARQL** for querying everything (possibly from multiple repositories)

SPARQL endpoints offer unified query access to RDF repositories
ex : Fuseki, Virtuoso, QLever...

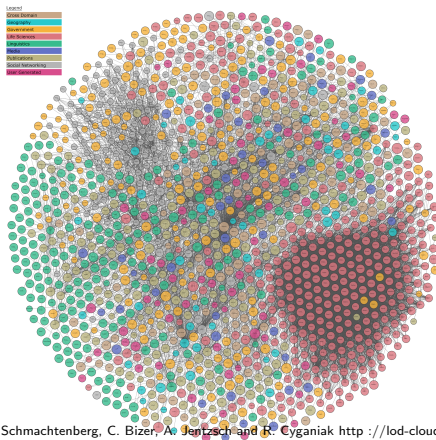
Linked Open Data : a federation of RDF repositories

LODStats (<http://lodstats.aksw.org/>) [Ermilov2016]

- 9960 datasets ; 149.10^9 triples
- general scope ; Life sciences = major field (size+density)

Linked open data (in 2025-09-02)

- RDF repositories can be queried in SPARQL via endpoints
- data from one endpoint can make references to data from another endpoint



Linked open data cloud, by M. Schmachtenberg, C. Bizer, A. Jentzsch and R. Cyganiak <http://lod-cloud.net/>

So... problem solved ?

Some successes

- Semantic Web technologies have been around for 20+ years
- Most reference knowledge bases are part of the LOD cloud

[Review](#) > [NPJ Digit Med.](#) 2019 Sep 10;2:90. doi: 10.1038/s41746-019-0162-5. eCollection 2019.

Enabling Web-scale data integration in biomedicine through Linked Open Data

Maulik R Kamdar ¹, Javier D Fernández ², Axel Polleres ², Tania Tudorache ¹, Mark A Musen ¹

Affiliations + expand

PMID: 31531395 PMID: PMC6736878 DOI: 10.1038/s41746-019-0162-5

> [Sci Data.](#) 2021 Jan 21;8(1):24. doi: 10.1038/s41597-021-00797-y.

An empirical meta-analysis of the life sciences linked open data on the web

Maulik R Kamdar ¹, Mark A Musen ²

Affiliations + expand

PMID: 33479214 PMID: PMC7819992 DOI: 10.1038/s41597-021-00797-y

So... problem solved ?

Some successes

- Semantic Web technologies have been around for 20+ years
- Most reference knowledge bases are part of the LOD cloud

[Review](#) > [NPJ Digit Med.](#) 2019 Sep 10;2:90. doi: 10.1038/s41746-019-0162-5. eCollection 2019.

Enabling Web-scale data integration in biomedicine through Linked Open Data

Maulik R Kamdar ¹, Javier D Fernández ², Axel Polleres ², Tania Tudorache ¹, Mark A Musen ¹

Affiliations + expand

PMID: 31531395 PMID: PMC6736878 DOI: 10.1038/s41746-019-0162-5

> [Sci Data.](#) 2021 Jan 21;8(1):24. doi: 10.1038/s41597-021-00797-y.

An empirical meta-analysis of the life sciences linked open data on the web

Maulik R Kamdar ¹, Mark A Musen ²

Affiliations + expand

PMID: 33479214 PMID: PMC7819992 DOI: 10.1038/s41597-021-00797-y

But world domination is not there yet

- Most reference knowledge bases... **... but not all of them**
 - ▶ no Gene Ontology SPARQL endpoint
 - ▶ The EBI SPARQL endpoints have closed
 - ▶ SIB (Switzerland) and NBDC+DBCLS (Japan) are the new hotspots
- Adoption by both users and developers is still far away
- Some are even abandoning SW for more fancy solutions that address immediate needs but miss the big picture
 - ▶ integration : OmniPath
 - ▶ querying : BioCypher

Challenges

Challenge : data and knowledge integration

- 2.000+ reference knowledge bases
 - countless project-specific databases and knowledge base
- each with its own, complex data schema
 - they are here and yet, we seldom (re-)use them
 - the reflex still too often is to integrate snapshots, expose it and let it become soon out of date
 - the decentralized web is still a (far away) vision

> *Sci Data*. 2024 Dec 18;11(1):1338. doi: 10.1038/s41597-024-04070-w.

Petagraph: A large-scale unifying knowledge graph framework for integrating biomolecular and biomedical data

Benjamin J Stear ^{✉ 1}, Taha Mohseni Ahooei ^{✉ 1}, J Alan Simmons ², Charles Kollar ², Lance Hartman ³, Katherine Beigel ³, Aditya Lehari ³, Shubha Vasisht ³, Tiffany J Callahan ³, Christopher M Nemerich ¹, Jonathan C Silverstein ², Deanne M Taylor ^{4 5}

Affiliations + expand

PMID: 39695169 PMCID: PMC11655564 DOI: 10.1038/s41597-024-04070-w

> *Sci Data*. 2022 Jul 11;9(1):393. doi: 10.1038/s41597-022-01510-3.

The heterogeneous pharmacological medical biochemical network PharMeBInet

Cassandra Königs ¹, Marcel Friedrichs ², Theresa Dietrich ²

Affiliations + expand

PMID: 35821017 PMCID: PMC9276653 DOI: 10.1038/s41597-022-01510-3

Challenge : data and knowledge analysis

Beyond the reductionist approach

- not quite sure what questions to ask
- not even in the mindset of (most) end-users
- how to do it ? before, after, before and after integration ?
- where to run the analysis ? Infrastructure like clusters, Galaxy are available but their user base is far from spanning the life science community
- what to do with the results ? there are so many of them, we still need methods for making sense out of them

Challenge : data and knowledge life cycle

- the FAIR principles have been around for 10 years [Wilkinson2016]

1. Solution: $\frac{1}{2} \times 100 = 50$ percent

The FAIR Guiding Principles for scientific data management and stewardship

[illegible]

diffusion: σ expand

5. *Front Plant Sci.* 2016 May 12;7:641. doi: 10.3389/fpls.2016.00641. eCollection 2016.

Publishing FAIR Data: An Exemplar Methodology Utilizing PHI-Base

Alejandro Rodríguez-Iglesias¹, Alejandro Rodríguez-González², Alistair G. Irvine³, Ane Sesma¹,
Martin Lirio⁴, Kim E. Hammond-Kroack⁴, Mark D. Wilkinson¹

PMID: 27433158. PLoS ONE: [10.1371/journal.pone.0160645](https://doi.org/10.1371/journal.pone.0160645)

[5. #1000866](#), 2007 Jan 31 4:54:18, doi: 10.13648/1000866arch.10044.5, eCollection 2007

Best practice data life cycle approaches for the life sciences

Philippe C. Goffin^{1,2}, Jyoti Khosla³, Kate S. Lohme⁴, Suzanne E. Lemaire⁵, Sarah Orbach⁶, Andrew Peak⁷, Bernard Pope⁸, Keith Reusser⁹, Keith Russell¹⁰, Tordur Sæmundsson¹¹, Andrew Tindler¹², Sonika Tripathi¹³, Jeffrey M. Christensen¹⁴, Sorenson-Stapleton¹⁵, Simon Gladwin¹⁶, Sander B. Hangeroth¹⁷, Helen¹⁸, Kyeles¹⁹, William W. Haines²⁰, Gabriel Kozlovsky-Ginsburg²¹, Peter K. Rothman²², Peter Natch²³, Phyllis R. Presson²⁴, Mark F. Richardson²⁵, Nathan S. Wilson-Hughes²⁶, Kelly L. Womers²⁷, Neil S. Young²⁸, Marko Vukobratovic²⁹ & 30

Affiliations + contact

- Few datasets actually follow them (but support is developping)
- <https://fair-checker.france-bioinformatique.fr>

► J Biomed Semantics. 2023 Jul 1;14(1):7. doi: 10.1186/s13326-023-00289-5

FAIR-Checker: supporting digital resource findability and reuse with Knowledge Graphs and Semantic Web standards

Alban Gaignard¹, Thomas Rosnet^{2,3}, Frédéric De Lamoignon⁴, Vincent Lefort^{5,3},
Marie-Dominique Devignes⁶

PMID: 37303266 PMCID: PMC10315041 DOI: 10.1186/s13326-023-00289-5

N. Sci. Data, 2015 Feb 24;13(1):1329. doi: 10.1038/nrd1587-025-04453-4

Applying the FAIR Principles to computational workflows

Sean R Wilkinson ¹, Mahmud Aliqasbi ², Khalid Bhatnagar ³, Michael R Cranoe ⁴,
Bruno de Paula Kinoshita ⁵, Luiz Godinho ⁶, Daniel Garjo ⁷, Ove Johan Ragnar Gestafsson ⁸,
Mark Juty ⁹, Peter Kawai ¹⁰, Farah Zaid Khan ¹¹, Johannes Köster ¹²,
Kristen Selensky von Gehlen ¹³, Luke Pouchard ¹⁴, Randy K Rannow ¹⁵, Sten Seefeld-Ries ^{2, 16},
Nicola Sonzogni ¹⁷, Shaohui Su ², Ziheng Sun ¹⁶, Balázs Vile ¹⁵, Merrilee A Wouters ¹⁸,
David Vaux ¹⁹, Christine Voigt ²⁰

[Affiliations + expand](#)

FAIRify

- data
- software
- workflows + workflow runs

Will generate even more data

Challenge : reproducibility

How to ensure the reproducibility ?

- of the data we produce
- of the software we develop
- of the analyses we perform
- of the conclusion we draw

Will generate even more data

		Data	
		Same	Different
Code	Same	Reproducible	Replicable
	Different	Robust	Generalisable

[Whitaker2016]

Review > Patterns (N Y). 2021 Sep 10;2(9):100322. doi: 10.1016/j.patter.2021.100322.

The role of metadata in reproducible computational research

Jeremy Leipzig ¹, Daniel Nüst ², Charles Tapley Hoyt ³, Karthik Ram ⁴, Jane Greenberg ¹

Affiliations + expand

PMID: 34553169 PMID: PMC8441584 DOI: 10.1016/j.patter.2021.100322



> PLoS Comput Biol. 2022 Mar 9;18(3):e1009935. doi: 10.1371/journal.pcbi.1009935. eCollection 2022 Mar.

Urgent need for consistent standards in functional enrichment analysis

Kaumadi Wijesooriya ¹, Sameer A Jadaan ², Kaushalya L Perera ¹, Tanuveer Kaur ¹, Mark Ziemann ¹

Affiliations + expand

PMID: 35263338 PMID: PMC8936487 DOI: 10.1371/journal.pcbi.1009935



A road map to improve software sharing practices

Roberto Di Cosmo, Sabrina Granger, Konrad Hinsén, Nicolas Jullien, Daniel Le Berre, Violaine Louvet, Camille Maumet, Clémentine Maurice, Raphaël Monat & Nicolas P. Rougier

The road ahead

Open questions

How to promote usage by users ?

- How to make data + knowledge + softwares + workflows available ?
- How to shield them from the technical difficulties ?
- How to handle the complexity of data schema that is a consequence of Life Science intrinsic complexity ?
- How to accomodate the various kinds of reasoning that users (legitimately) want to perform but are beyond the scope of Semantic Web (boolean networks, path finding, temporal reasoning,...) ?

Open questions

How to promote usage by users ?

- How to make data + knowledge + softwares + workflows available ?
- How to shield them from the technical difficulties ?
- How to handle the complexity of data schema that is a consequence of Life Science intrinsic complexity ?
- How to accomodate the various kinds of reasoning that users (legitimately) want to perform but are beyond the scope of Semantic Web (boolean networks, path finding, temporal reasoning,...) ?

How to promote contribution by users ?

- How to facilitate integration of project-specific datasets with SW ?
- How to facilitate exposing their datasets ?
(hopefully in a FAIR perspective and workflow-compatible way)

Open questions

How to promote usage by users ?

- How to make data + knowledge + softwares + workflows available ?
- How to shield them from the technical difficulties ?
- How to handle the complexity of data schema that is a consequence of Life Science intrinsic complexity ?
- How to accomodate the various kinds of reasoning that users (legitimately) want to perform but are beyond the scope of Semantic Web (boolean networks, path finding, temporal reasoning,...) ?

How to promote contribution by users ?

- How to facilitate integration of project-specific datasets with SW ?
- How to facilitate exposing their datasets ?
(hopefully in a FAIR perspective and workflow-compatible way)

Beyond *Homo sapiens* : how to connect with ML ?