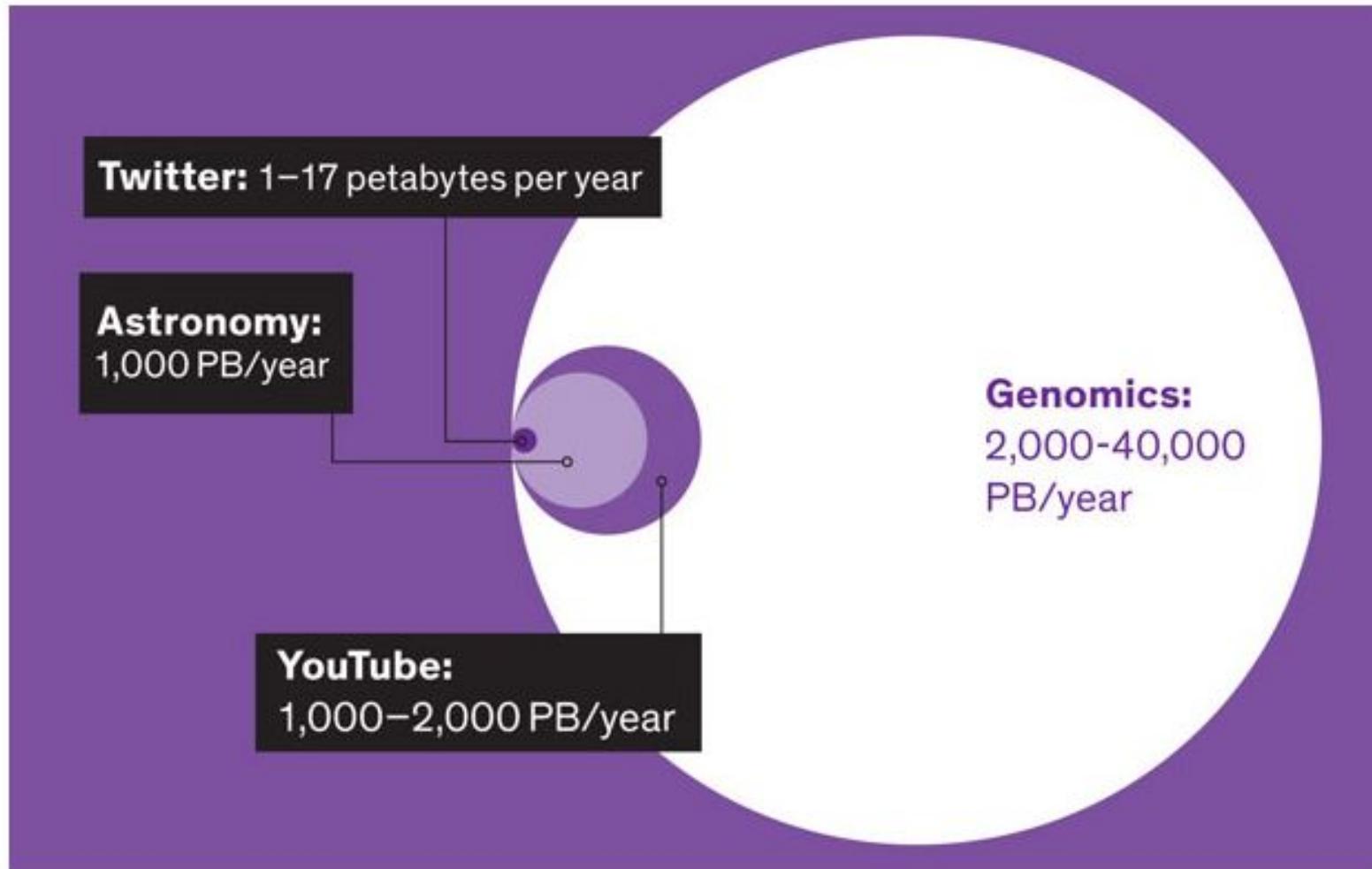


# APPLICATIONS OF LINKED DATA IN THE LIFE SCIENCES

DR. JASPER J. KOEHORST - LABORATORY OF SYSTEMS AND SYNTHETIC BIOLOGY

JASPER.KOEHORST@WUR.NL

## Projected annual storage in 2025



Source: ["Big Data: Astronomical or Genomical?" PLoS Biology, 7 July 2015.](#)

# GENOMICAL AMOUNT OF DATA

- Not only a lot of data
- Extremely diverse
- Interoperability problems: Inconsistencies in describing genomic information
- Current data standards are not sufficient.
- The process of FAIRification is complex and costly
- 99.99% of information obtained is computationally predicted

## The Crazy-Ambitious Effort to Catalogue Every Microbe on Earth

It's the most macro study of the microscopic world ever published.

MENU ▾

nature  
International journal of science

Subscribe

NEWS · 02 NOVEMBER 2018

'Why not sequence everything?' A plan to decode every complex species on Earth

The Earth BioGenome Project aims to sequence 1.5 million genomes and will probably cost US\$4.7 billion.

Science

Contents ▾

News ▾

Careers ▾

Journals ▾

Researchers launch plan to sequence 66,000 species in the United Kingdom. But that's just a start

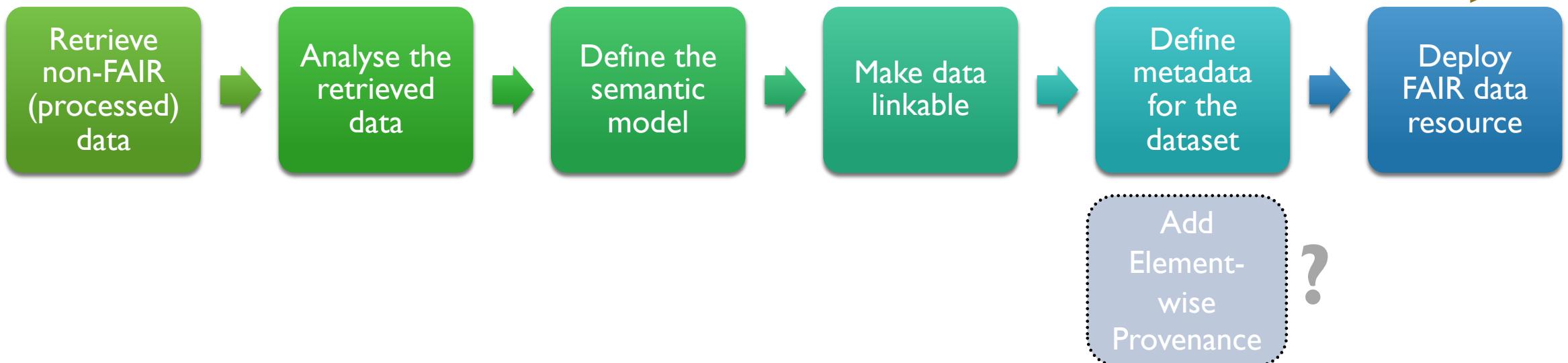
By Erik Stokstad | Nov. 1, 2018, 3:15 PM

## A CHINESE PROVINCE IS SEQUENCING ONE MILLION OF ITS RESIDENTS' GENOMES

THE PROJECT IS BEING ADVISED BY HARVARD GENETICIST GEORGE CHURCH.

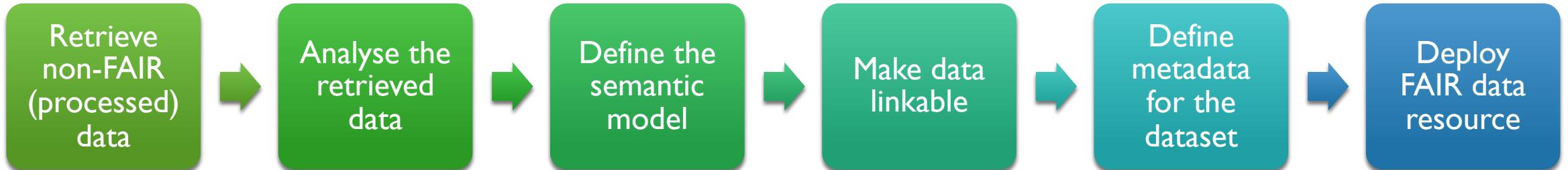
# FAIRIFICATION

> Degree of interoperability >

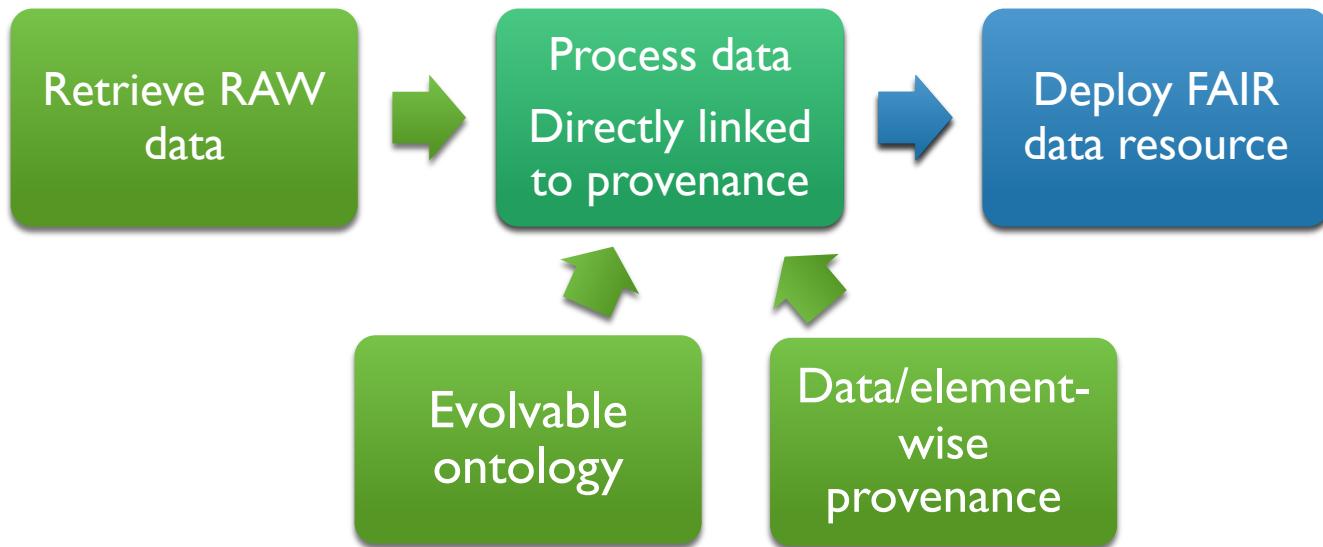


# FAIRIFICATION

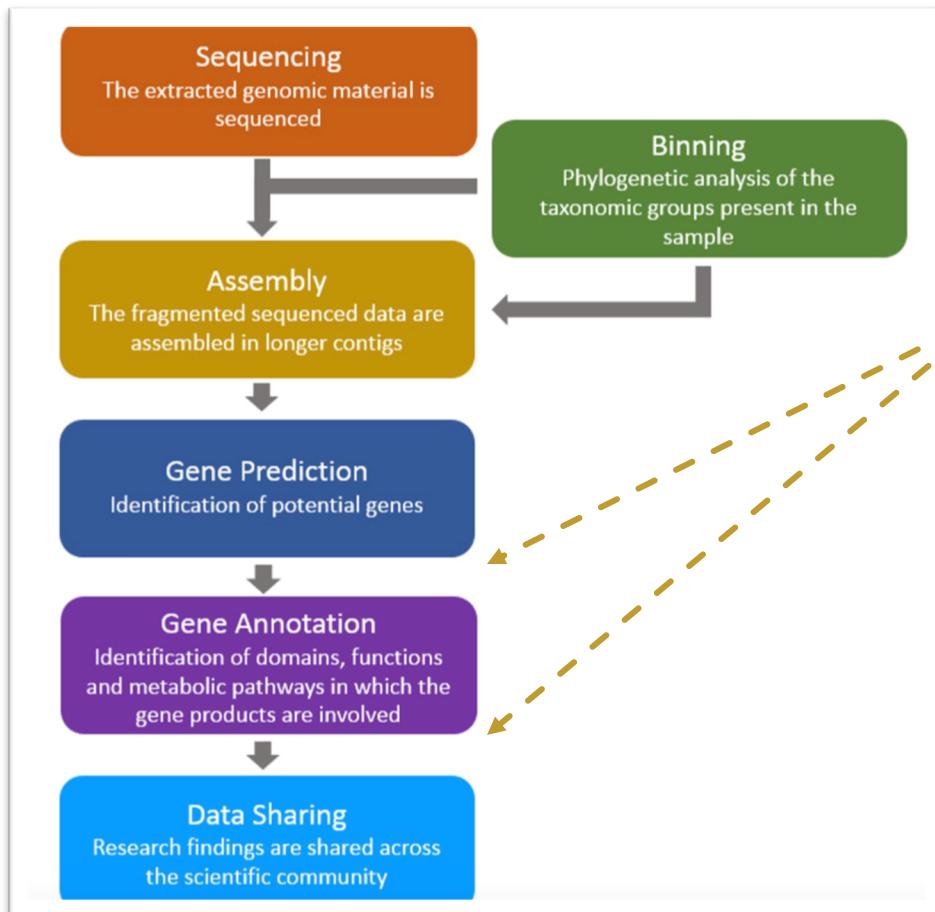
> Degree of interoperability >



# FAIR BY DESIGN



# CLASSICAL ANALYSIS WORKFLOW: DATA MEETS THE TOOL

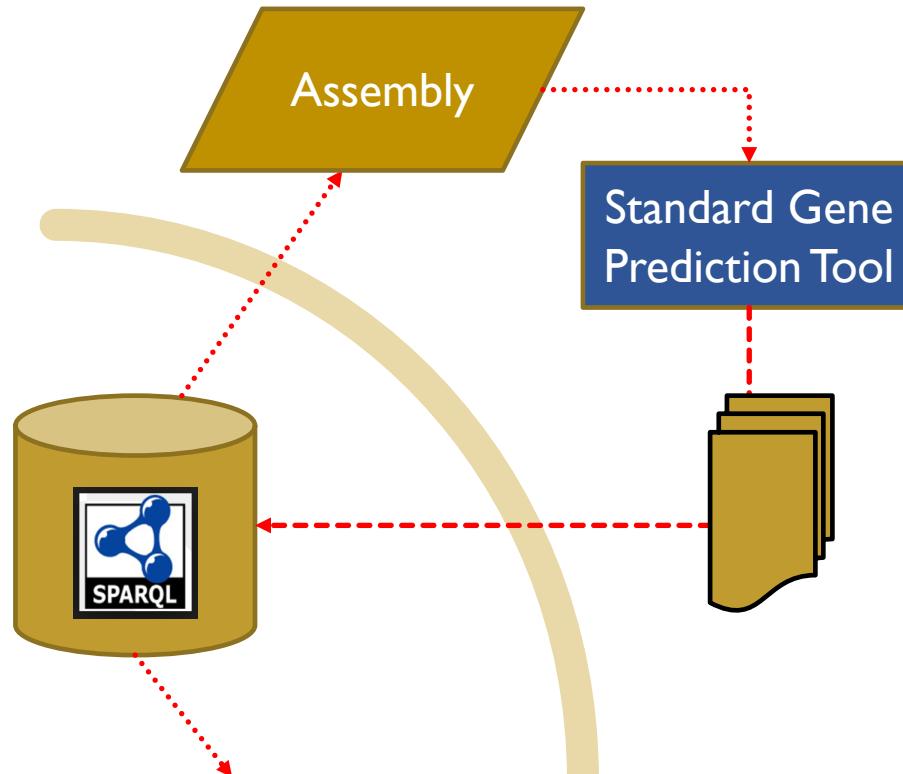


IN EACH STEP A DECISION IS MADE BASED ON A PRE-DEFINED THRESHOLD

## THE ACTUAL ELEMENT-WISE PROVENANCE

- E.G CONFIDENCE SCORES - IS LOST
- NO STATISTICAL ANALYSIS ON METADATA POSSIBLE
- CHANGING THE THRESHOLD VALUE = RECOMPUTE

# FAIR BY DESIGN SEMANTIC ANALYSIS WORKFLOW: TOOL MEETS THE DATA



FAIR Data Sharing – Linked Data  
Triple store can be queried with  
SPARQL

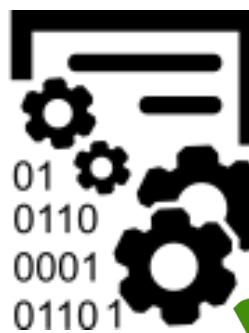
**DATA/ELEMENT-WISE PROVENANCE  
AUTOMATICALLY AND DIRECTLY  
LINKED TO THE PREDICTION**

Algorithm = ..  
Version = ..  
Run date = ..  
.....  
Genomic location  
Confidence score  
.....

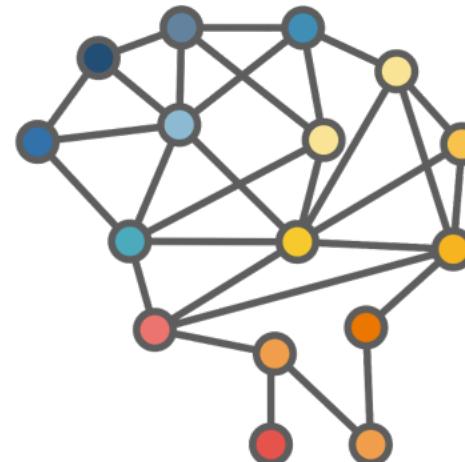
GBOL ontology for  
Interoperable  
genome annotation

**DYNAMIC THRESHOLD SELECTION (AFTERWARDS)  
MULTIPLE (SIMILAIR) TOOLS CAN BE COMPARED**

# FROM FAIRIFICATION TO FAIR BY DESIGN



RAW DATA



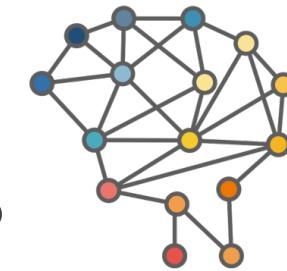
ONTOLOGIES



APPLICATION  
PROGRAMMING  
INTERFACE

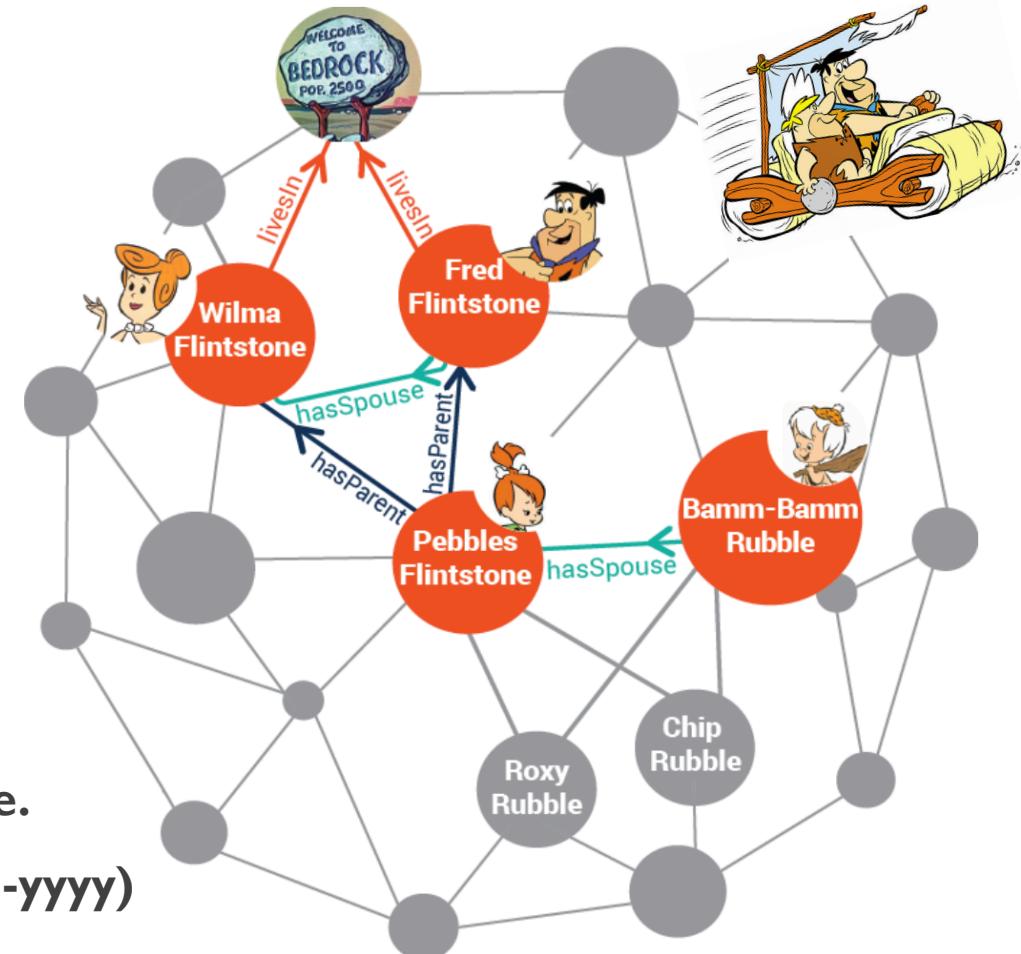
SEMANTIC FRAMEWORK & LINKED (RDF) DATA MODEL

# HOW TO ENSURE A HIGH DEGREE OF INTEROPERABILITY USING A GENOMICAL AMOUNTS OF DATA: **ONTOLOGIES**



- DEFINITIONS, DEFINITIONS AND MOST IMPORTANTLY DEFINITIONS
  - MINIMAL INFORMATION MODEL REQUIRED FOR INTEROPERABILITY
  - ONTOLOGY AND CONSTRAINED RELATIONS
  - SHEx DEFINITION (a language for validating and describing the RDF data according to the ontology)

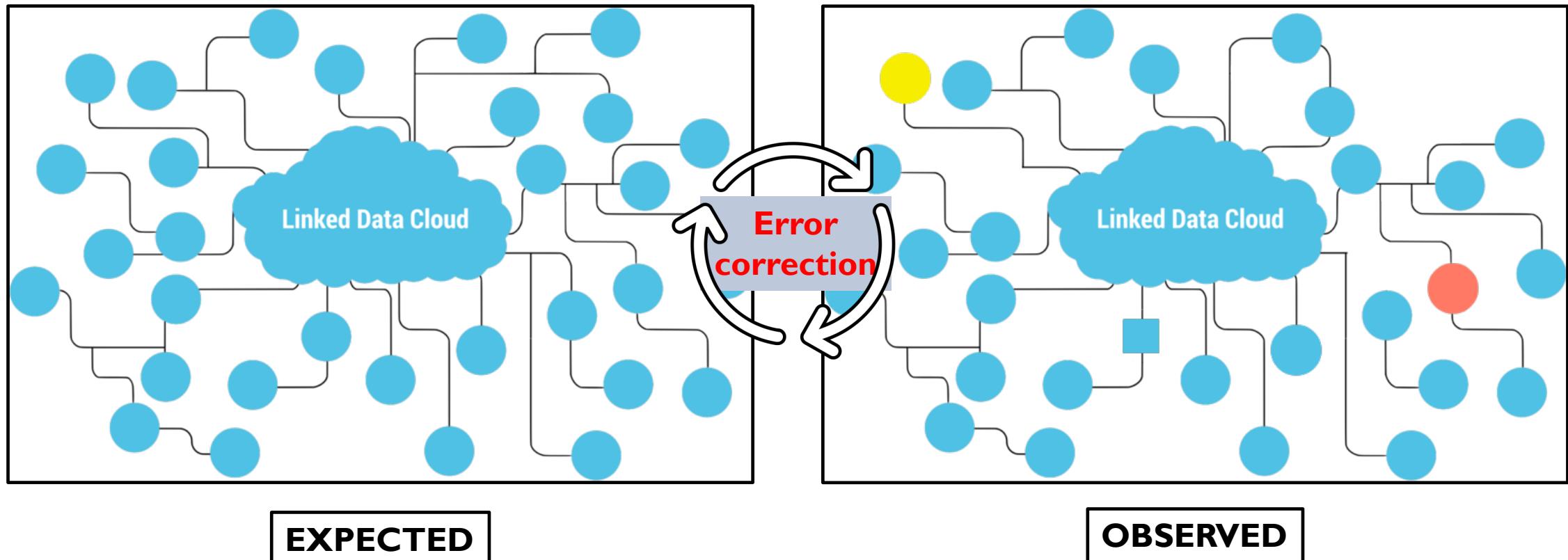
```
<cartoonfigureShape> {      #A cartoonfigure (e.g Fred) has:  
:name xsd:string+;          # one or more names  
:birthdate xsd:date?        # and optionally a (single) birthdate.  
}  
                                         in a specific format (e.g. dd-mm-yyyy)
```



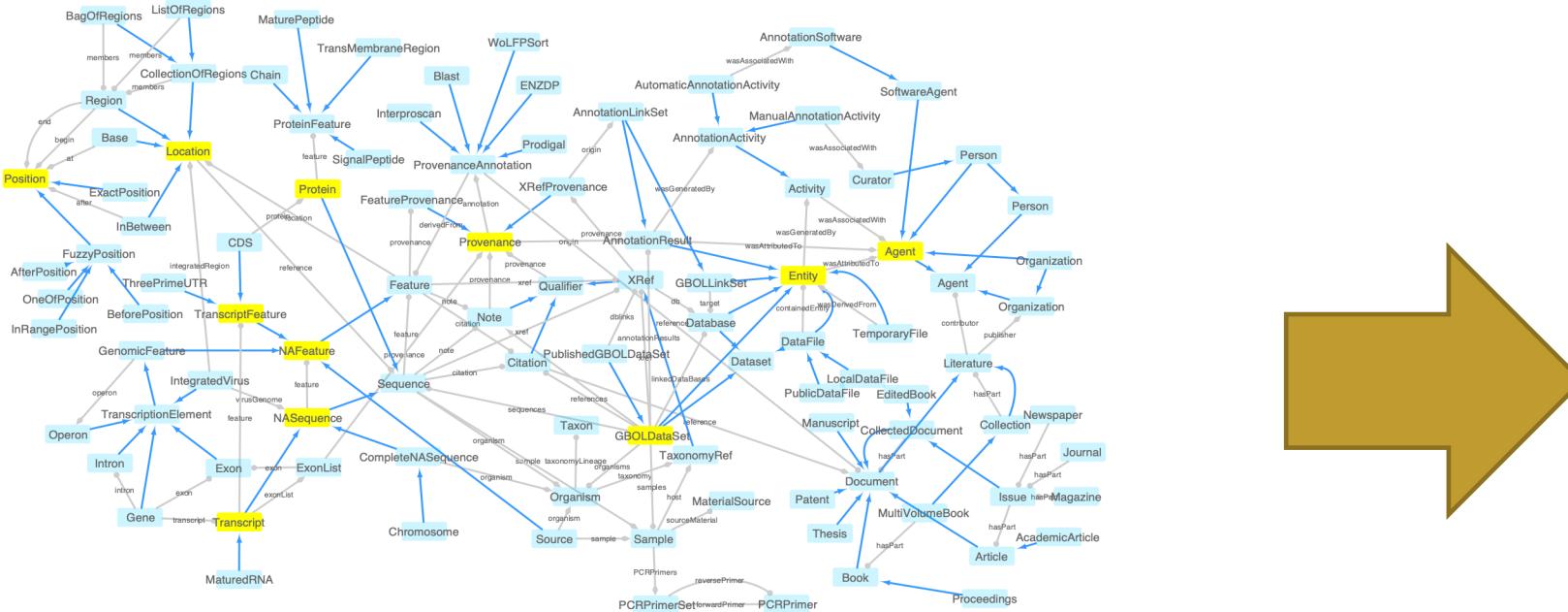




# THE DYNAMIC NATURE OF SEMANTIC WEB



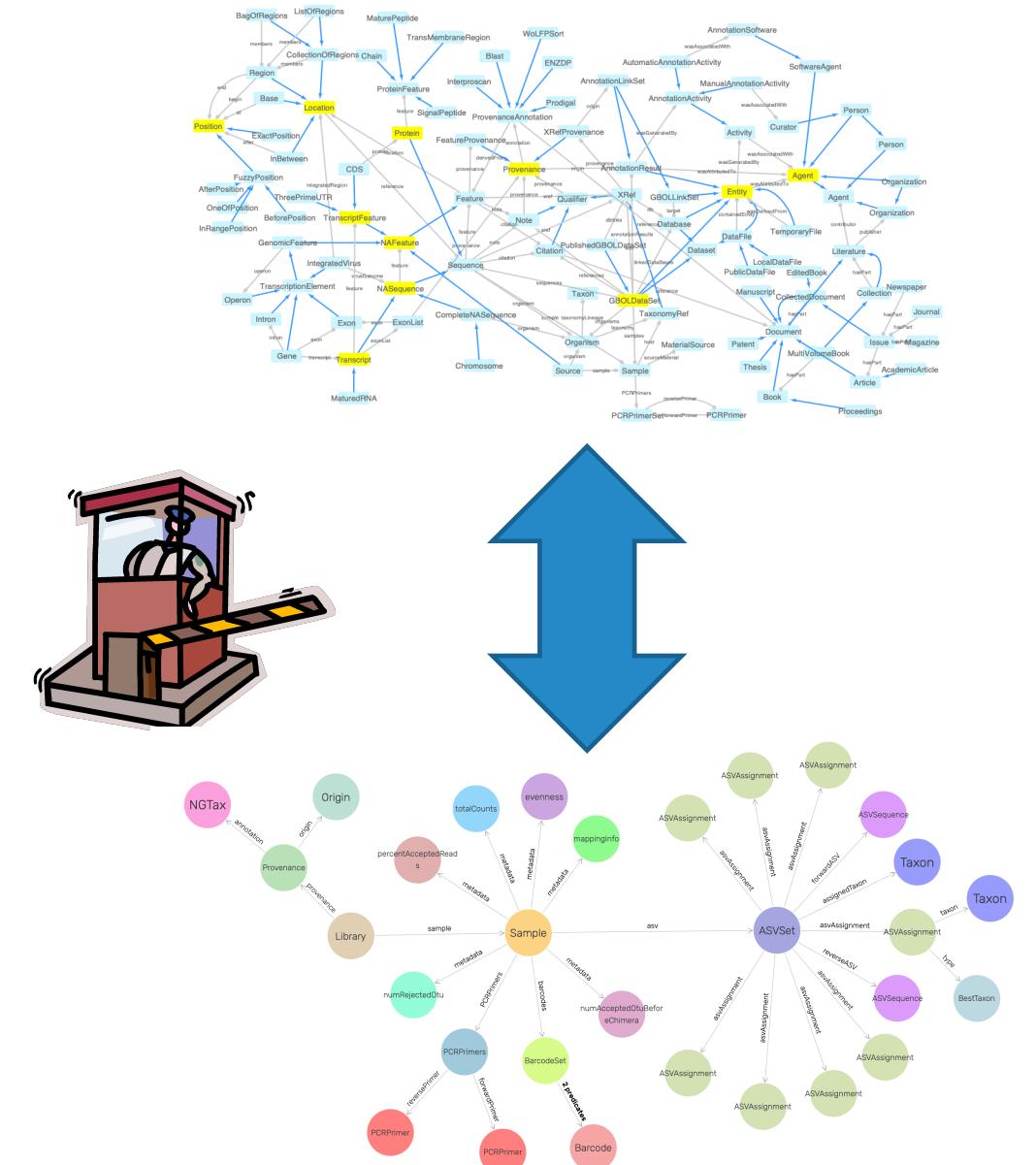
# EMPUSA: ONTOLOGY VALIDATOR & CODE GENERATOR



APPLICATION  
PROGRAMMING  
INTERFACE

EMPUSA GENERATES THE CODE THAT CAN ACT AS A GATEKEEPER FOR FAIR LINKED DATA GENERATION

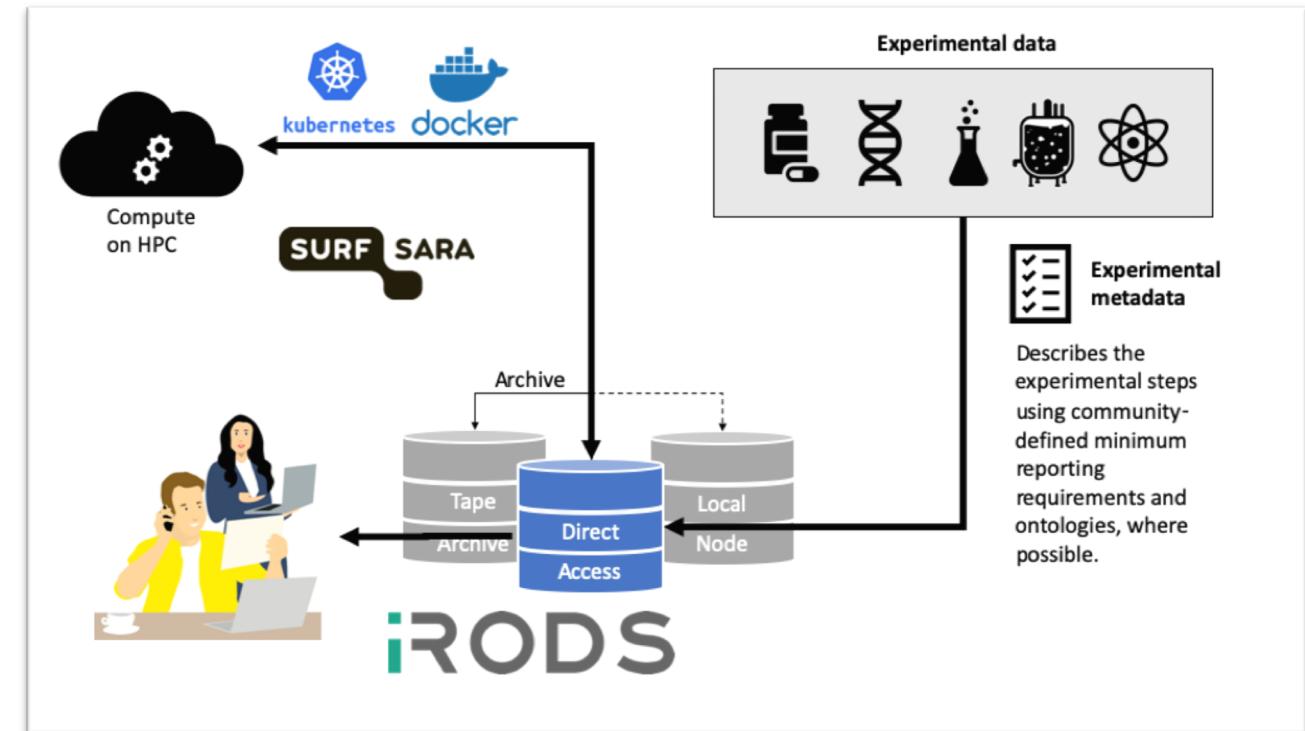
# FAIR BY DESIGN FROM DATA TO LINKED DATA



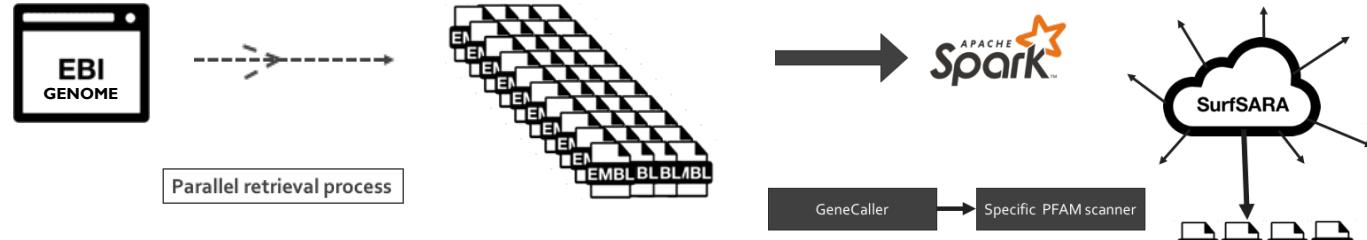
Validated linked (meta-)dataset

# APPLICATION: PROJECT MANAGEMENT

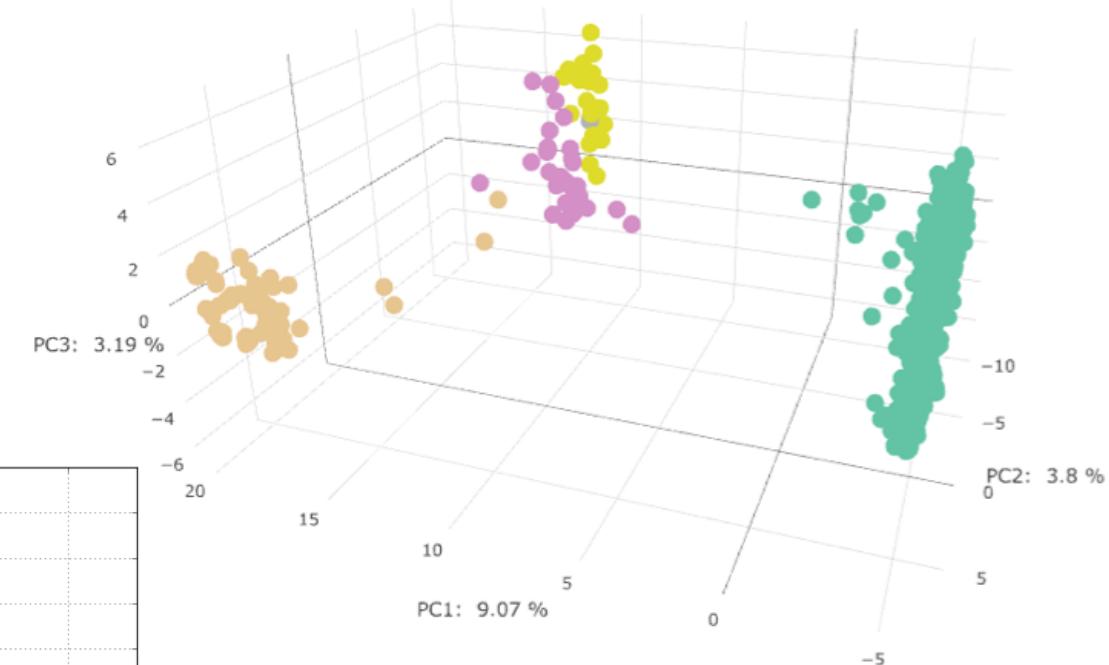
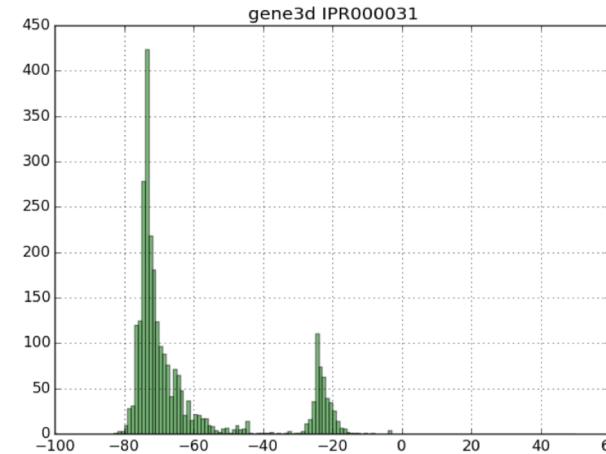
- According to the P-ISA principles
  - Project, Investigation, Study, Assay
  - All (meta-)data is available as linked data
  - Assay: compute results and provenance translated to RDF and embedded in the system



# APPLICATIONS: SAPP, FUNCTIONAL GENOME ANNOTATION & COMPARATIVE FUNCTIONAL GENOMICS OF THOUSANDS OF SPECIES



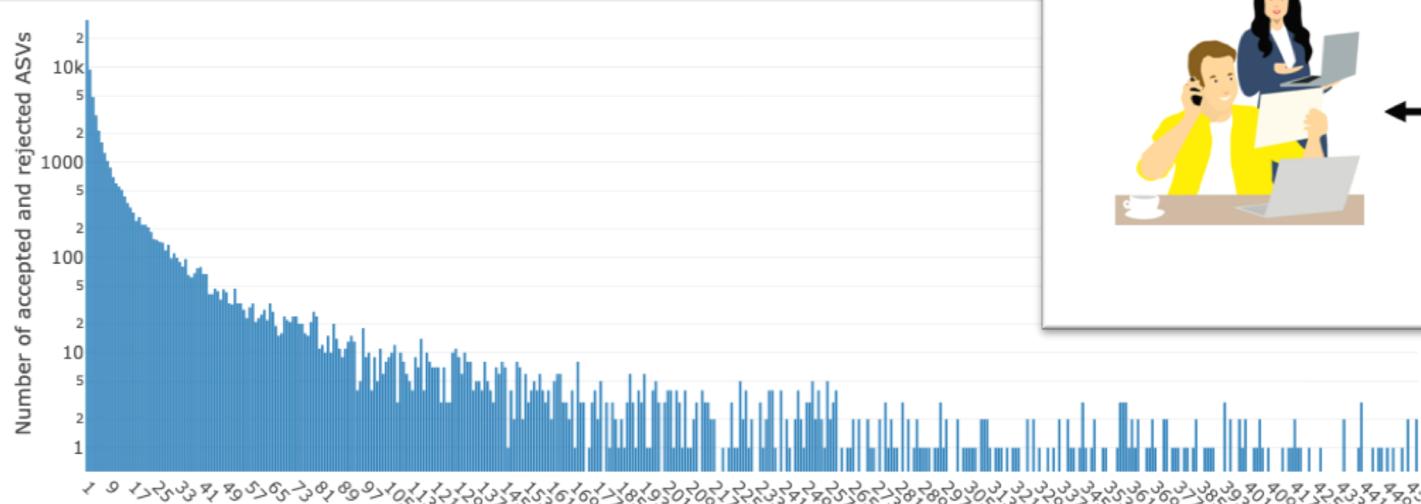
```
PREFIX gbol: <http://gbol.life/0.1/>
SELECT ?sample ?accession ?db
WHERE {
  ?sample a gbol:Sample .
  ?dnaobject gbol:sample ?sample .
  ?dnaobject gbol:feature ?gene .
  ?gene gbol:transcript ?transcript .
  ?transcript gbol:feature ?cds .
  ?cds gbol:protein ?protein .
  ?protein gbol:feature ?domain .
  ?domain gbol:xref ?xref .
  ?xref gbol:db <http://gbol.life/0.1/db/pfam> .
  ?xref gbol:accession ?accession .
} LIMIT 10
```



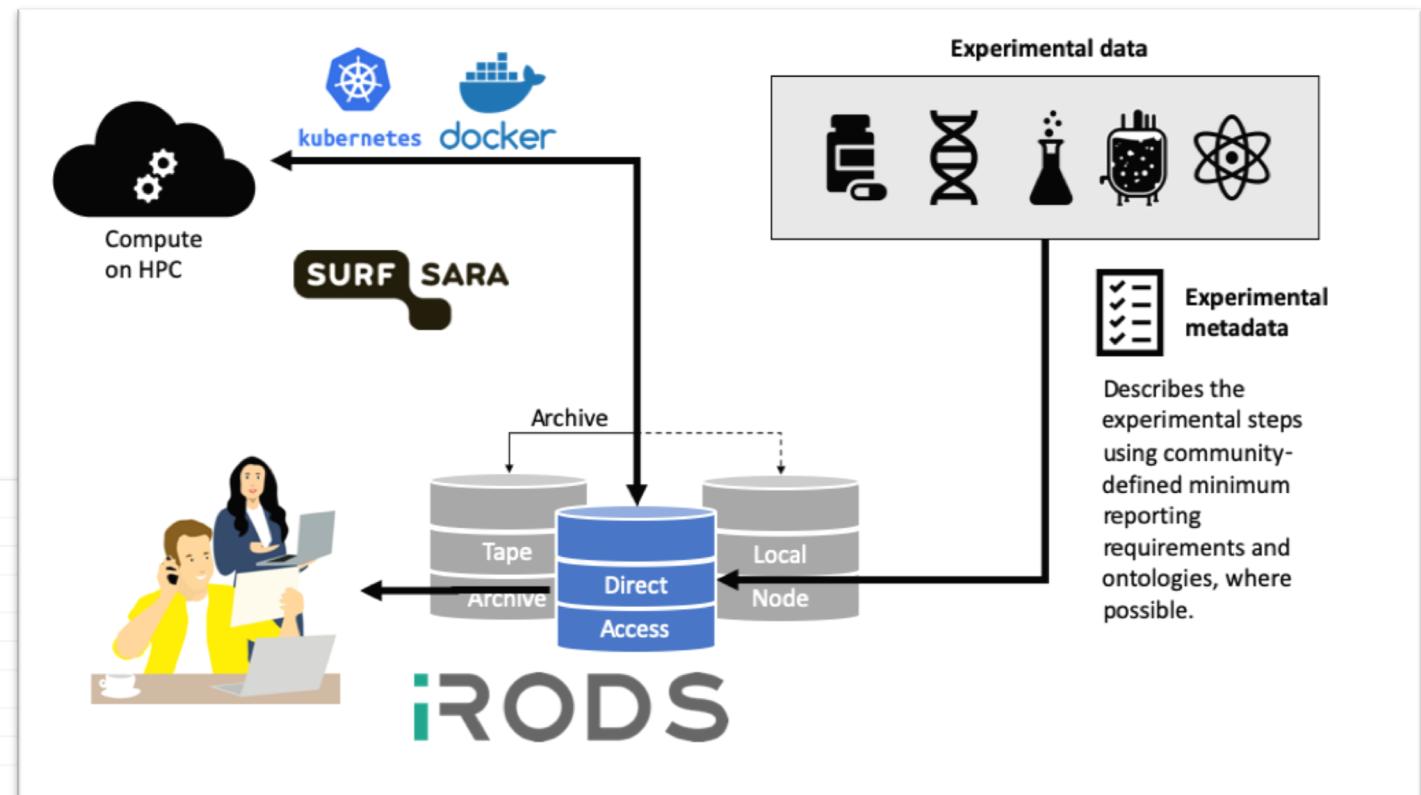
- Pseudomonas aeruginosa
- Pseudomonas fluorescens
- Pseudomonas putida
- Pseudomonas syringae

# APPLICATIONS: NGTAX 2.0 HIGH THROUGHPUT AMPLICON ANALYSIS

- TRACK & TRACE EACH SPECIES
- SEQUENCE MAPPING ACROSS SAMPLES
- LINKED METADATA
  - AGE, GENDER, DIET, ...



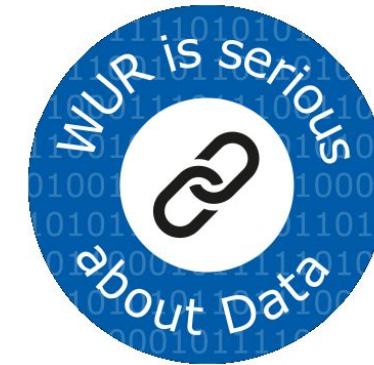
NUMBER OF ASVs SHARED BETWEEN ~1800 SAMPLES



ANY OMICS CAN BE ANALYSED  
USING THE SAME PRINCIPLES

## SOFTWARE - AVAILABILITY

- GENOME BIOLOGY ONTOLOGY LANGUAGE (GBOL)  
[HTTP://GBOL.LIFE](http://gbol.life)
- EMPUSA, THE CODE GENERATOR  
[HTTP://EMPUSA.ORG](http://empusa.org)
- SAPP, SEMANTIC GENOME ANNOTATION  
[HTTP://SAPP.GITLAB.IO](http://sapp.gitlab.io)
- PROJECT MANAGEMENT  
[HTTP://M-UNLOCK.NL](http://m-unlock.nl)
- NGTAX 2.0, AMPLICON ANALYSIS  
[HTTP://WURSSB.GITLAB.IO/NGTAX/](http://wurssb.gitlab.io/ngtax/)



WDCC

