# CINECA

Common Infrastructure for National Cohorts in Europe, Canada, and Africa

# Bringing it all together: human cohort standards, tools and applications

Presenter: Dr. Melanie Courtot, Ontario Institute for Cancer Research (OICR)
Host: Vera Matser (EMBL-EBI)

# About this webinar

This webinar is being recorded and will be disseminated afterwards

After the presentation we will address the questions posted by the audience using the Q&A function

# Common Infrastructure for National Cohorts in Europe, Canada and Africa

**The vision:**

Accelerating disease research and improving health by facilitating transcontinental human data exchange

**Stay informed**

@CinecaProject

www.cineca-project.eu

**The challenges:**

1. Federated Discovery

REGISTER & LOGIN
2. AAI

3. Harmonised Metadata

4. Federated Analysis in research/healthcare
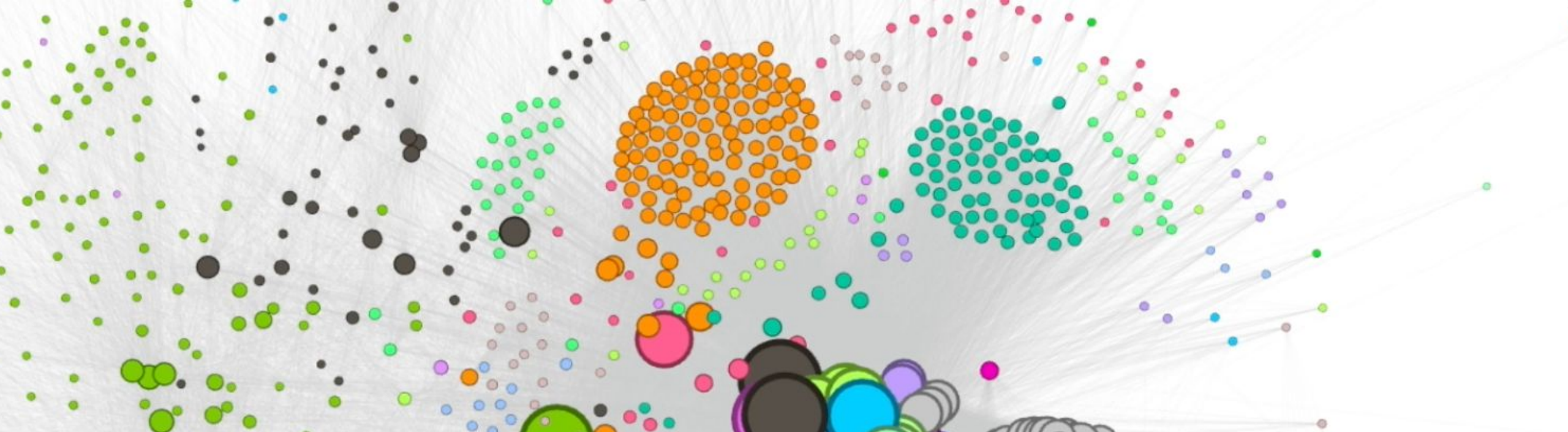
5. ELSI Framework

# Today's presenter

**Dr. Melanie Courtot** is Director of Genome Informatics and incoming Principal Investigator at the Ontario Institute for Cancer Research (OICR). Her team develops new software, databases and other necessary components to store, organize and compute over the large and complex datasets being generated by OICR's cancer research programs.

Dr Courtot is passionate about translational informatics - building intelligent systems to gain new insights and impact human health.

She co-leads the Data Use and Cohort representation groups for the Global Alliance for Genomics and Health (GA4GH), as well as cohort harmonization efforts for Common Infrastructure for National Cohorts in Europe, Canada, and Africa (CINECA), the International HundredK+ Cohorts Consortium (IHCC) and the Davos Alzheimer's Collaborative.

Melanie can be found twitter, @mcourtot, where she often posts about science, equity and diversity, food and silly things she or her children do.

# Bringing it all together: human cohort standards, tools and applications

**Mélanie Courtot, PhD**
Ontario Institute for Cancer Research
mcourtot@oicr.on.ca
@mcourtot

*CINECA webinar, March 31$^{st}$ 2022*

**Human cohorts for disease research**



**FAIR data management**

**Human cohorts for disease research**

**FAIR data management**

# International HundredK+ Cohort consortium



- Canada (835,000)
- Iceland (250,000)
- Norway (5,750,000)
- Finland (950,000)
- Estonia (52,000)
- United Kingdom (2,341,000)
- Sweden (998,000)
- Denmark (5,400,000)
- China (2,513,000)
- United States (7,787,000)
- France (162,000)
- Israel (140,000*)
- Iran (230,000)
- India (21,000)
- Japan (540,000)
- Korea (742,000)
- Mexico (160,000)
- Saudi Arabia (102,000)
- Taiwan (92,000)
- Malaysia (107,000)
- Multi-Country (1,051,000)
- Brazil (15,000)
- Singapore (10,000)
- Australia (267,000)
- Chile (7,000)

Legend: <100K | 100K - 249K | 250K - 999K | 1M and up

~60 cohorts, ~30M participants

# Global Alliance for Genomics and Health (GA4GH)

Common Infrastructure for National Cohorts in Europe, Canada, and Africa (CINECA)

Use cases for care providers and patients

# Building the IHCC cohort atlas



Cohort presentation and display

Intuitive filtering by cohort metadata & data dictionary attributes

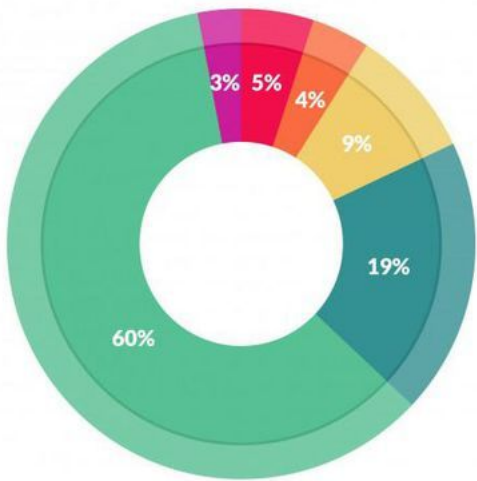Reference to external cohort sites

# Human cohorts for disease research



# FAIR data management

What data scientists spend the most time doing

- Building training sets: 3%
- Cleaning and organizing data: 60%
- Collecting data sets; 19%
- Mining data for patterns: 9%
- Refining algorithms: 4%
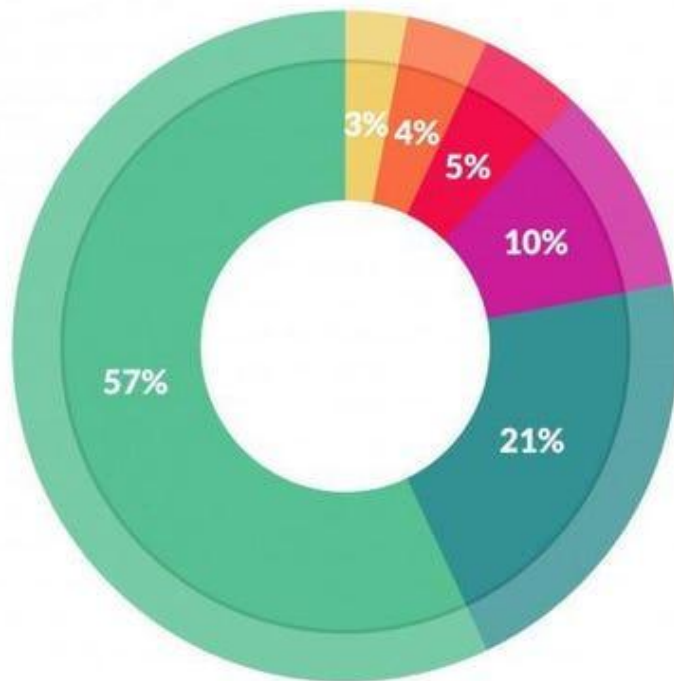- Other: 5%

Forbes article on 2016 Data Scientist Report
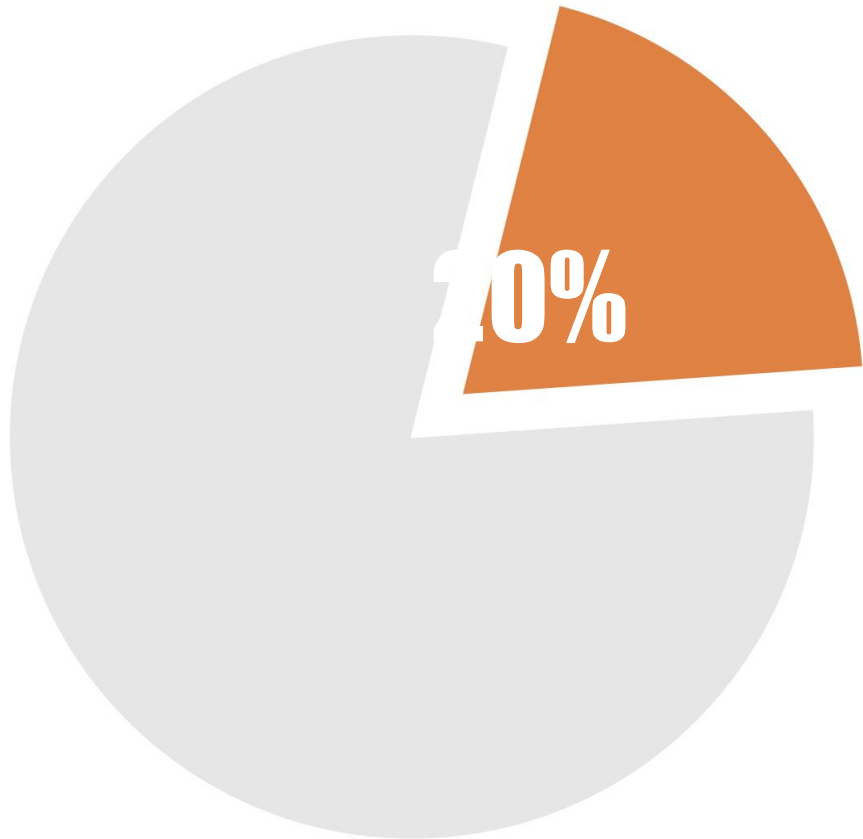
What data scientists spend the most time doing

- Building training sets: 3%
- Cleaning and organizing data: 60%
- Collecting data sets; 19%
- Mining data for patterns: 9%
- Refining algorithms: 4%
- Other: 5%



What's the least enjoyable part of data science?

- Building training sets: 10%
- Cleaning and organizing data: 57%
- Collecting data sets: 21%
- Mining data for patterns: 3%
- Refining algorithms: 4%
- Other: 5%

16

Forbes article on 2016 Data Scientist Report

# "Sometimes it's easier to rewrite genetics than update Excel"

20%

**Symbols that affect data handling and retrieval.** For example, all symbols that autoconverted to dates in Microsoft Excel have been changed (for example, *SEPT1* is now *SEPTIN1*; *MARCH1* is now *MARCHF1*); tRNA synthetase symbols that were also common words have been changed (for example, *WARS* is now *WARS1*; *CARS* is now *CARS1*).

*Ziemann, M., Eren, Y. & El-Osta. Genome Biol 17, 177 (2016). https://doi.org/10.1186/s13059-016-1044-7*
*Bruford, E.A., Braschi, B., Denny, P. et al. Nat Genet52, 754–758 (2020). https://doi.org/10.1038/s41588-020-0669-3*
*https://www.theverge.com/2020/8/6/21355674/human-genes-rename-microsoft-excel-misreading-dates*
*Abeysooriya M, Soria M, Kasu MS, Ziemann M (2021) PLOS Comp Bio 17(7): e1008984. https://doi.org/10.1371/journal.pcbi.1008984*

> **Everyone wants to do the model work, not the data work.**
>
> Sambavisan et al.

# The importance of FAIR data



F indable  A ccessible  I nteroperable  R eusable
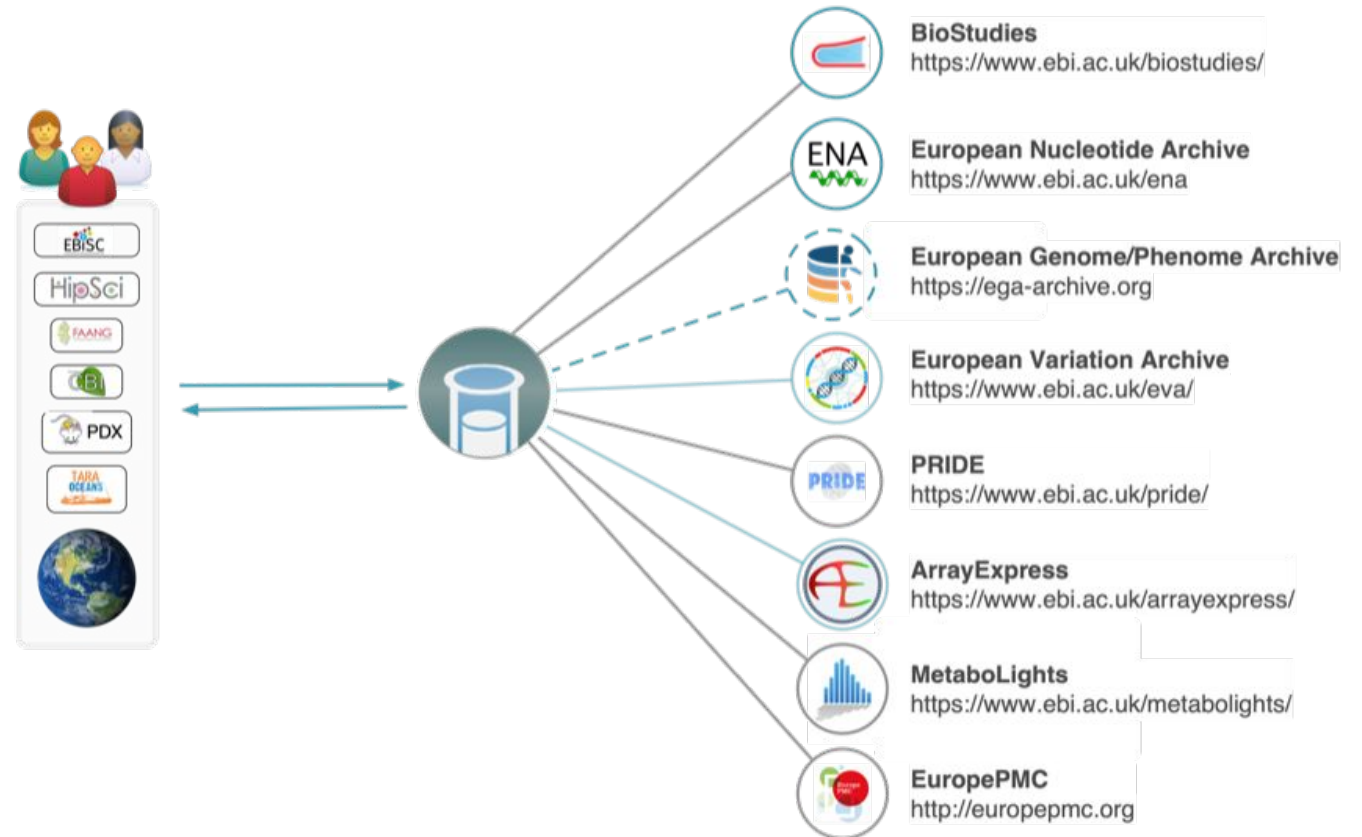
Image credit: Sungya Pundir, Wikimedia Commons CC BY-SA 4.0

- More data resources than ever before
- Combining data from "federated" data repositories for analysis is becoming increasingly common
- The impact of AI and machine learning is rapidly growing, and require data that computers can understand
- Current data generation efforts do not support this level of reuse

# EMBL-EBI BioSamples : a metadata hub

- Archive of information about biological materials

- From internal and external sources, and directly from submitters.

- Enables technology independent linking between assay data and sample metadata



BioStudies
https://www.ebi.ac.uk/biostudies/

ENA — European Nucleotide Archive
https://www.ebi.ac.uk/ena

European Genome/Phenome Archive
https://ega-archive.org

European Variation Archive
https://www.ebi.ac.uk/eva/

PRIDE
https://www.ebi.ac.uk/pride/

ArrayExpress
https://www.ebi.ac.uk/arrayexpress/

MetaboLights
https://www.ebi.ac.uk/metabolights/

EuropePMC
http://europepmc.org

https://www.ebi.ac.uk/biosamples/

# Linking Samples through relationships

| Relationship types | Reverse relationships | Description |
| --- | --- | --- |
| derived from | derived from (reverse) | Sample A is derived from Sample B. E.g. - Tissue samples derived from donor samples - Cell line samples derived from tissue samples - Viral samples separated from saliva samples - Organoid samples cultured from tissue samples |
| same as | same as | Sample A is the same as Sample B. This can be used to link duplicated samples |
| has member | has member (reverse) | Sample A is a member of Sample group G. BioSamples create a sample group for each sampleTab submission*. It's also possible to put patient samples as a sample group. |
| child of | child of (reverse) | Sample A is the child of Sample B. E.g - Patient A is the child of Patient B |



*Related patient-derived xenografts(PDX) samples*

https://www.ebi.ac.uk/biosamples/docs/guides/relationships

# Finding the right data is hard without good metadata



GA4GH
use cases

Find all biosamples
for male individuals

Find all individuals with an
AML diagnosis
AND Down's Syndrome (trisomy 21)

Find all individuals originating from
a specified geographical location
by latitude/longitude

Find all datasets for which
individuals have consented
to code X

...

# BioSamples dataset size

**55,000**
unique attributes

**18 million**
samples

**375 million**
key:value data points



CS0760_2                                                    SAMN15744766

Updated on: 07-10-2020 18:21

External Id   SAMN15744766     INSDC center name   CSIR-Institute of Genomics and Integrative Biology

INSDC first public   2020-08-10T00:00:00Z     INSDC last update   2020-08-11T06:20:02.500Z

INSDC secondary accession   SRS7179058     INSDC status   live     NCBI submission model   Pathogen.cl

NCBI submission package   Pathogen.cl.1.0     SRA accession   SRS7179058     host   Homo sapiens

organism   Severe acute respiratory syndrome coronavirus 2     replicate   Biological Replicate 760

strain   SARS-CoV-2     title   Negative Control 8

# Metadata curation



Redundancy and inconsistency in real life data

**Researchers**

I need to find all COVID related samples

**COVID-19 Data Portal**

**Covid 19-related attributes in BioSamples:**
- *severe acute respiratory syndrome*
- *COVID19*
- *novel coronavirus pneumonia*
- *nCoV pneumonia*
- *COVID-19*
- *Coronavirus infected disease-19 (COVID-19)*

**Metadata curation**
- Text curation
- Semantic annotation

**Common challenges in sample metadata**
- Special characters (*COVID19* vs *COVID-19*)
- Acronyms (*T2D for diabetes*)
- Typo
- Synonyms

# Text curation and semantic annotation

- Automatic curation by pipelines
- Manual curation by experts
- Curation tool based on manual curation and machine learning



Manual curation

Semantic annotation

Text curation

Raw metadata

# The **OLS** to access and visualize ontologies

The Ontology Lookup Service (OLS) is repository for biomedical ontologies providing access to up-to-date ontology resources (UI + API)



http://purl.obolibrary.org/obo/MONDO_0100096

# ZOOMA **for automated ontology annotation**

- An annotation service mapping ontology terms to free text
- Stores known rules (e.g., manual curation) to guide future annotations

**SARS-CoV-2 related values**
- SARS-CoV2
- Wuhan coronavirus
- Human coronavirus 2019
- SARS-CoV-2
- 2019-nCoV
- COVID-19 virus

OLS / NCBI organismal classification   NCBITAXON  /  NCBITaxon:2697049   Copy

## Severe acute respiratory syndrome coronavirus 2

http://purl.obolibrary.org/obo/NCBITaxon_2697049   Copy

Tree view | Term mappings | Term history

root
 Viruses
  Riboviria
   Orthornavirae
    Pisuviricota
     Pisoniviricetes
      Nidovirales
       Cornidovirineae
        Coronaviridae
         Orthocoronavirinae
          Betacoronavirus
           Sarbecovirus
            Severe acute respiratory syndrome-related coronavirus
             Severe acute respiratory syndrome coronavirus 2

Graph view
Reset tree
Show all siblings

https://www.ebi.ac.uk/spot/zooma/

# ~5 Million samples in the EMBL-EBI COVID-19 Data Portal (March 2022)

# Graph search across archives

Researchers

As a researcher, I want to find the immunotyping data of all lung samples from **COVID19 patients** and corresponding genome sequencing data of the **viral isolate**, to study how the immune systems response to viral infection.

# Connecting attributes for recommendation

# Connecting attributes for recommendation

## Semantic based validation

https://www.npmjs.com/package/elixir-jsonschema-validator

## Structured schema annotation

https://bioschemas.org/types/BioSample/0.1-RELEASE-2019_06_19/

## Structured phenotype exchange

https://phenopacket-schema.readthedocs.io/en/latest/biosample.html

**Human cohorts for disease research**



**FAIR data management**

# Use case : As a researcher, I am looking for cohorts with XXX data

# Bringing cohort data together

1. **Data models** to represent both access conditions and cohort data
2. **Tools** and processes for implementations
3. Deployment over **clinical cohorts**

# GA4GH Data Use Ontology

- Vocabulary describing permitted data uses and modifiers

- "General research use", "disease-specific research", "not for profit only"...

https://www.ebi.ac.uk/ols/ontologies/duo

https://github.com/EBISPOT/DUO



https://ega-archive.org/datasets/EGAD00010001859

# Semantic harmonisation

**To promote and publish it, the CINECA model was formalised as an ontology**

- using World Wide Web Consortium standards
- Adopting OBO Foundry best practices
- Leverages (and contributes) to existing resources for maximal interoperability
- Available publicly (CC-BY)



https://www.ebi.ac.uk/ols/ontologies/gecko

# Genomics Cohorts Knowledge Ontology

- **Commonly used attributes to describe cohort metadata**
- **"Medication", "sample type", "genomics datatypes"...**

https://www.ebi.ac.uk/ols/ontologies/gecko
https://github.com/IHCC-cohorts/GECKO

**Harmonization process**

# Data harmonisation process:

1. Data collection
2. Metadata model design
3. Harmonisation

# Harmonisation example



43

# Harmonisation example



**Clinical Measurement Ontology**

GECKO classes

KoGES classes

# Harmonisation example

**Applying these techniques to clinical cohorts...**

# Use case : As a researcher, I am looking for cohorts with XXX data

I am looking for cohorts with 'blood measurement' data

IHCC cohort atlas

KoGes:Glucose (fasting) is a subtype of Blood glucose level value: defined in OWL as 'blood measurement value' and ('has target' some glucose)

GCS:bpleft1dbp is a subtype of diastolic blood pressure value, subtype of blood pressure measurement value

....

# Automating the process



Refine import pipeline to (semi) automatically ingest data dictionaries provided as CSV files

Develop automated mapping using existing tooling and text-mining processes

# Automated mapping pipeline for cohort owners



**Cohort registry**

# Cohort registry

## Deploy cohort registry

- Provide dereferencing for human readability, e.g., http://purl.obolibrary.org/obo/GECKO_0000068 returns a human readable page describing tobacco history in GECKO
- Include versioning and change detection for re upload

## Built-in interoperability with mapping/curation tools

EMBL-EBI OLS



https://registry.ihccglobal.app/index

# IHCC cohort mappings

- Stores mapping between GECKO and cohort terms

- Accessible through APIs

- Parameter to bridge between mappings: If A ⟺ B and B ⟺ C then can infer A⟺ C



https://mapping.ihccglobal.app

EMBL-EBI

# IHCC Cohort Atlas



- User interface and search API leverage the reusable modules from overture

- 12 cohorts deployed

- In the process of adding more cohorts

Cohort presentation and display

Reference to external cohort sites

Intuitive filtering by cohort metadata & data dictionary attributes

https://atlas.ihccglobal.org

# Summary: bringing it together

# EMBL-EBI semantic toolkit



**Search/visualize ontologies**

**Ontology cross mapping**

**Annotate data**

# The Overture suite



- Flagship product powering Software Engineering team projects

- Interoperability with community standards e.g., Global Alliance for Genomics and Health (GA4GH)

*Next steps:*
- *Provide simple self-installer to make the product more accessible*
- *Enable further front-end customization by 3rd party developers*

https://www.overture.bio/

https://www.davosalzheimerscollaborative.org
Video demo: https://vimeo.com/505253841

# VirusSeq data portal

# VirusSeq data portal



1 month from project start to delivery!

# VirusSeq data portal



1 month from project start to delivery!

Over 250,000 viral genomes as of March 2022

# Community engagement is key

60

# >200,000 datasets annotated with the Data Use Ontology

# Cohort Representation - Use Cases

**Global Alliance**
for Genomics & Health

Pain Points:

Lack of interoperability

Demographic variables not standardized

Hard to perform analyses when there are different baseline measures/descriptions

Values for each individual often not available/shareable

Requirements:

Have an interoperable computable cohort definition standard that goes across OMOP / FHIR / HL7 CQL etc.

Provide a "computational" description of the cohort, e.g., one that could be run as a query on the baseline population

Allow for model representation of different sources, including annotation linking to semantic definitions (data elements, vocabularies, ontologies). Ability to create, share and use mappings.
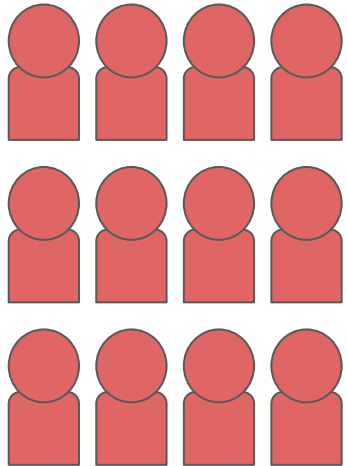
# Cohort Registry Discovery
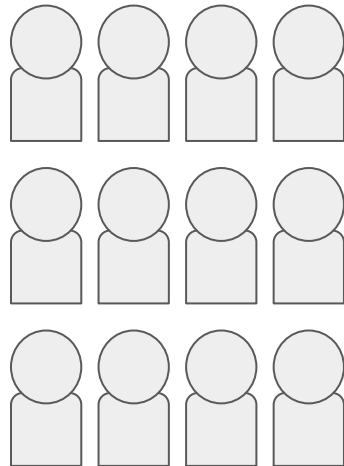
I am looking for cohorts with '**blood measurement**' data"

| Cohort Name | Country | Bio Specimens | Sample Type | ... | Enrollment |
|---|---|---|---|---|---|
| Cohort 1 | South Africa | Yes | Blood (fasting) | ... | 11,681 |
| Cohort 2 | England | No | Urine | ... | 350,000 |
| Cohort 3 | Northern Ireland | Yes | Creatinine | | 2,856 |
| ... | | | | | |

# The IHCC cohort atlas for cohort discovery



Cohort presentation and display

Intuitive filtering by cohort metadata & data dictionary attributes

Reference to external cohort sites

# Computable Cohort Discovery

I am looking for **female** Asthmatic patients with **creatinine** > **0.9**

| Patient ID | Age | Gender | ... | Lab Result (Creatinine) | ... | Diagnosis |
|------------|-----|--------|-----|-------------------------|-----|-----------|
| 1 | 34 | M | ... | 0.72 | ... | SNCT(26929004) |
| 2 | 45 | F | ... | 1.13 | ... | ICD10CM(J45) |
| 3 | 61 | F | ... | 0.81 | | |
| ... | | | | | | |

# Computable Cohort Discovery

# Computable Cohort Representation



- Minimum Information About Computable Cohorts (MIACC)
- Reuse Existing Standards:
  - GECKO / BBMRI / - Registry Alignment
  - FHIR CQL / OMOP Cohort - Query Alignment
  - Phenopackets - Payload alignment
  - Beacons - API Alignment

PID: hjmg45-2344-jnm2b34-2w34@0.0.3
Name: Asthma & Diabetic Patients
Description:
Authors: Susheel Varma, …
Version: 0.0.3
Revisions: 0.0.1, 0.0.2
Changelog: []
Parents []: kj67-2321-3452mn-243234
Created: 2021-04-14 12:01:00
Updated: 2021-04-14 13:01:00
Cohort Type: Study | User
Projects: hg345-234j-2343,…
----
Entry_date:
Exit_date:
Interval Type: Closed | Open
Criteria_Groups: [
  {query_criteria_inclusion_exclusion}
]
Collection_Events: []
  - Event Type: Incident | Prevalent | Other
  - Entry Type: Single | Multiple
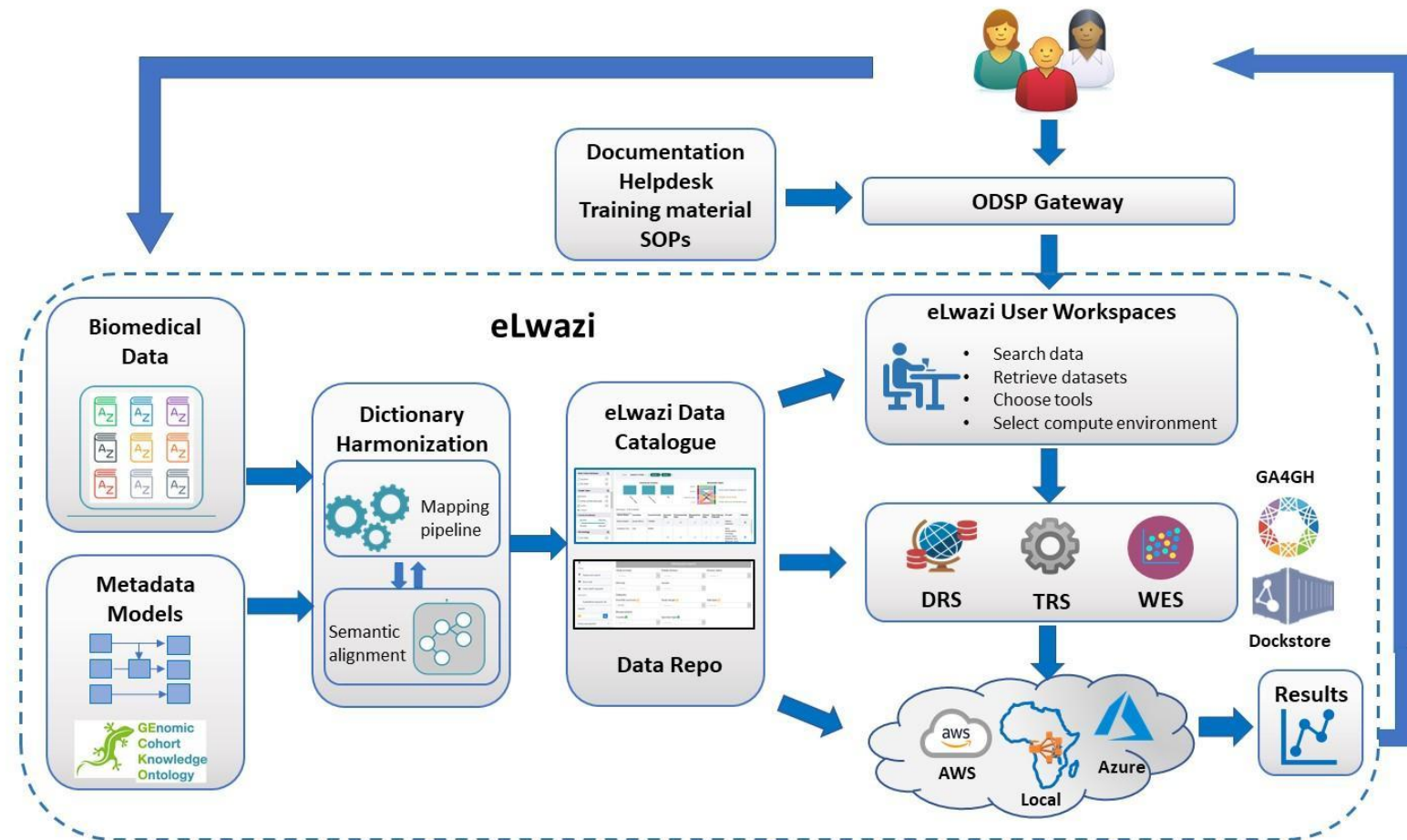----
Datasets: []
  - Name
  - ID
  - Count:
- DUO Code (Restrictions & Limitations)
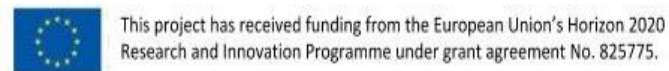
# NOA September 2021

NIH RFA-RM-20-018 on Harnessing Data Science for Health Discovery and Innovation in Africa (DS-I Africa) Open Data Science Platform and Coordinating Center
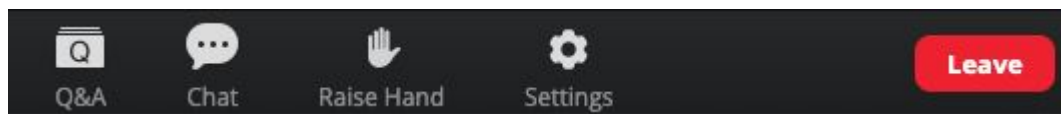
Nicky Mulder

# Acknowledgements

# Questions?

Please write your questions using the Q&A button



Bringing it all together: human cohort standards, tools and applications

Presenter: Dr. Melanie Courtot (OICR)