# CORBEL

shared services for life-science

4$^{th}$ AGM, 02 March 2020, Brussels, Belgium

# WP6 – Data access, management and integration

**Carole Goble** on behalf of **Helen Parkinson,** Morris Swertz, Jason Swedlow, Thomas Keane, Ilkka Lappalainen and WP6 Partners

# WP6  ACTIVITIES



Year 1      Year 2      Year 3      Year 4

Identifier best practice  and checklists

Semantic infrastructure and ontologies

Secure data access technology

Pan infrastructure deployment / enhancement -> EOSC Life

A **checklist based framework** for systematic documentation, gap analysis, recommendations and actions developed through **6 case studies**



**BBMRI-ERIC** ®

1 Rare Disease

2 Biobanking

EURO-BIOIMAGING

3 Imaging data

EMBRC
EUROPEAN MARINE BIOLOGICAL RESOURCE CENTRE

4 Marine metazoan models

5 Ocean sampling

Open PHACTS
Open Pharmacological Space

ISBE Infrastructure for Systems Biology Europe

6 Genes, proteins and drugs

**Checklist**
**Identifier Strategy**
- ✓ what is being identified?
- ✓ what is the granularity of the entity?
- ✓ data & identifier life cycles?
- ✓ entity visibility outside the RI?
- ✓ relationships are there between identifiers?
- ✓ how names and ontology terms mapped to names?
- ✓ identifier properties, policies and practices?
- ✓ lookup and resolve the identifier?
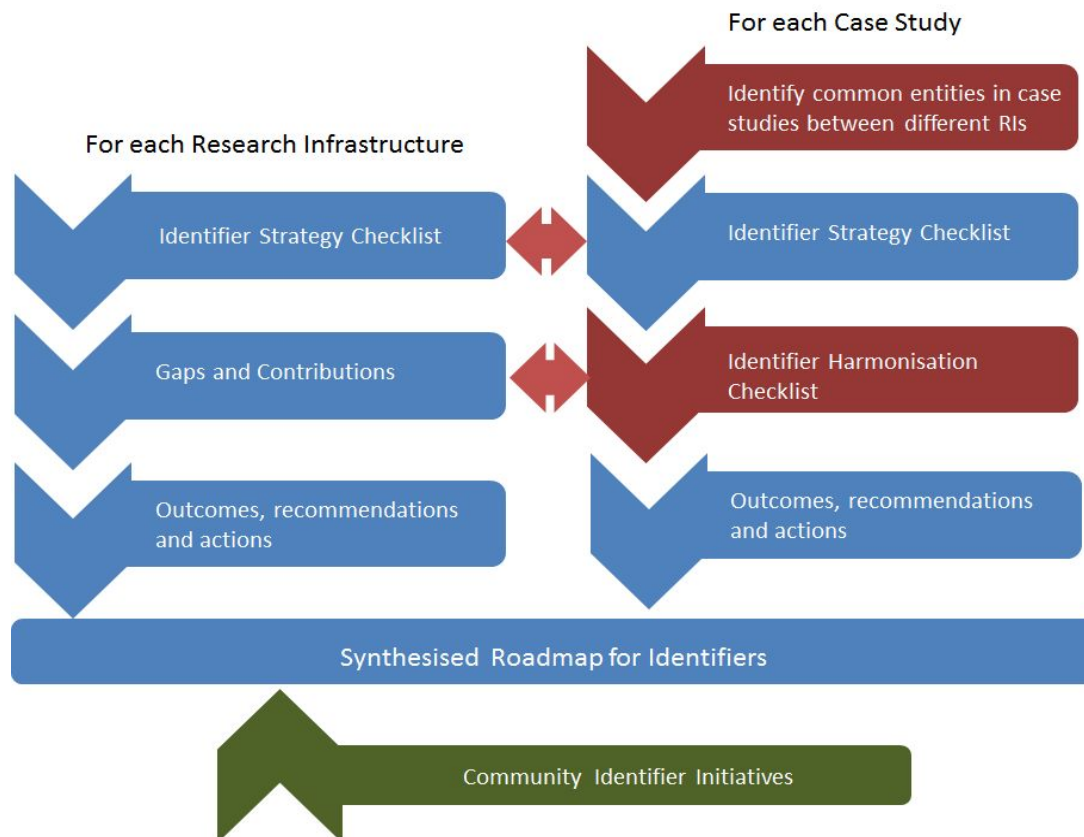- ✓ identifier metadata & data citation?

5

A **checklist based framework** for systematic documentation, gap analysis, recommendations and actions



For each Case Study

Identify common entities in case studies between different RIs

For each Research Infrastructure

Identifier Strategy Checklist

Identifier Strategy Checklist

Gaps and Contributions

Identifier Harmonisation Checklist

Outcomes, recommendations and actions

Outcomes, recommendations and actions

Synthesised Roadmap for Identifiers

Community Identifier Initiatives

**Checklist**
**Identifier Strategy**

- what is being identified?
- what is the granularity of the entity?
- data & identifier life cycles?
- entity visibility outside the RI?
- relationships are there between identifiers?
- how names and ontology terms mapped to names?
- identifier properties, policies and practices?
- lookup and resolve the identifier?
- identifier metadata & data citation?

Aligned ontology ids and Identifier.org
-> Ontology Lookup Service supports CURIEs

Provides unique stable, resolvable and location-independent compact URIs to identify and locate scientific data

Harmonisation with USA CDL name2thing (n2t.net) resolution service

- Standards compliant web and programmatic based access to ontologies and linked open data from an ontology access service

- Validated data-ontology maps with provenance between data and ontologies

- Ontology-ontology mappings with provenance supporting data integration across infrastructures

- Semantic infrastructure

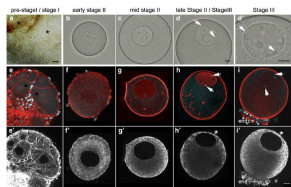| | | |
|---|---|---|
| How do I access ontologies? | OLS | View on GitHub |
| How do I map data to ontologies? | ZOOMA | View on GitHub |
| How do I translate from one ontology to another? | OxO | View on GitHub |
| How can I extend an ontology? | WEBULOUS | View on GitHub |
| How do I build "ontology aware" search applications? | BioSolr | View on GitHub |
| How do I publish this data? | RDF Platform | |

https://ebispot.github.io

# TASK 6.2 A SEMANTIC TOOLKIT

Build/extend an ontology
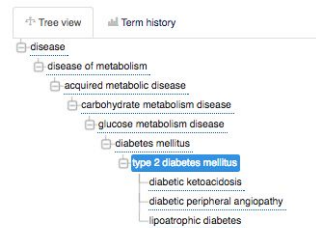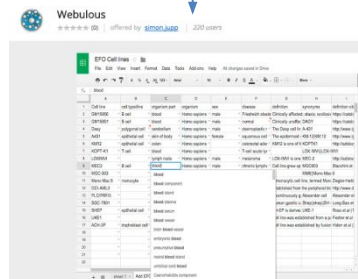
Access and visualise an ontology

Map between ontologies 'xref'

Map text to ontologies

"Euro-BioImaging has used identifiers for defining genes, proteins, antibodies, drugs, species, phenotypes, and organ systems and pathologies to publish ~170 TBytes of original image datasets

Identifiers are key for IDR to fulfil its function as an added value knowledgebase, making critical reference datasets well-annotated and ultimately linked, searchable and reusable.
Everything IDR has done has followed guidance and used tools developed within or related to CORBEL."



https://idr.openmicroscopy.org

Gene Product Targeting

Genetic

Geographic

Chemical

Histopathology

3D-

Super-resolution

Integrated studies

Thumbnails (of 5D Images)

Experimental metadata

Biomolecular annotations

Analysis results

Cross-data browsing

Cloud analysis

Download (local analysis)

Jason Swedlow

# RI DEPLOYMENT

# IMPROVED METADATA

Improving the handling of image annotations

Using the Cellular Microscopy Phenotype Ontology and the OLS can group attributes

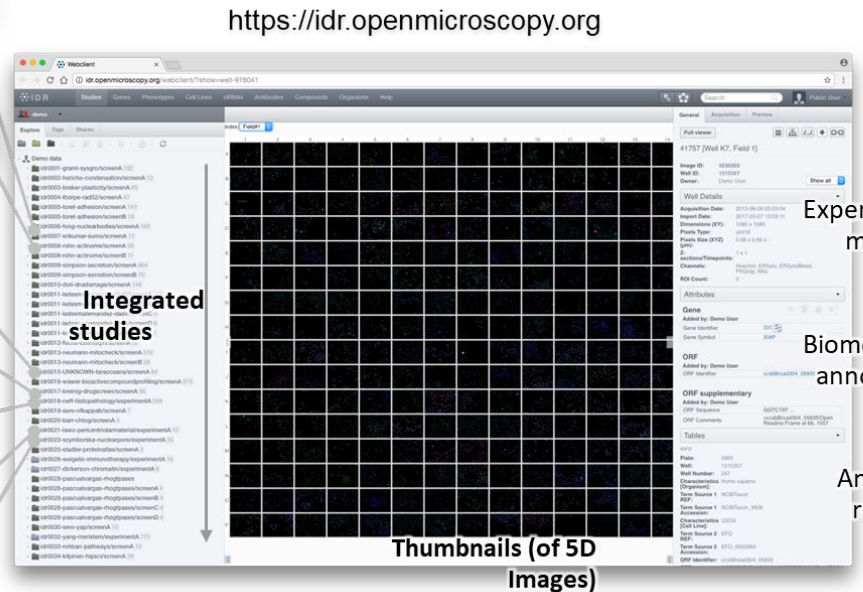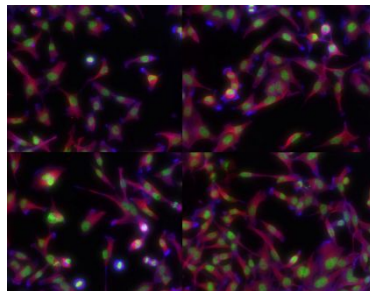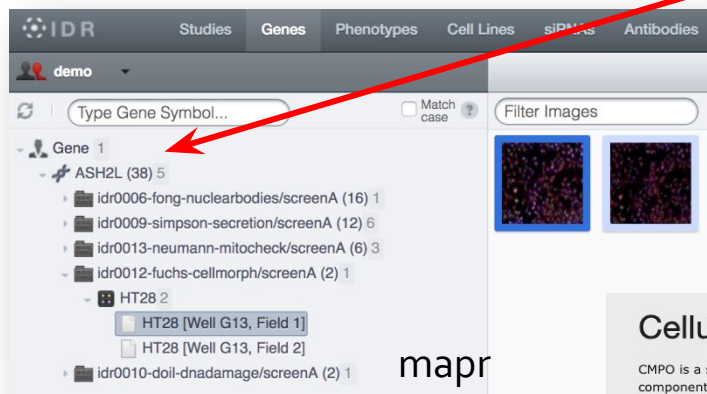Now can link between studies through the ontology

idr0012



**Attributes** 8

**Cell Lines**
Added by: Public data
Cell Line — HeLa

**Gene**
Added by: Public data
Gene Identifier — 9070
Gene Symbol — ASH2L

**Phenotype**
Added by: Public data
Phenotype — elongated cells
Phenotype Term Name — elongated cell phenotype
Phenotype Term Accession — CMPO_0000077

**OLS** ONTOLOGY SEARCH

**Cellular Microscopy Phenotype Ontology**

CMPO is a species neutral ontology for describing general phenotypic observations relating to the whole cell, cellular components, cellular processes and cell populations.

Search CMPO

mapr

# TASK 6.3 SECURE ACCESS TO SENSITIVE DATA

Federated authentication and access solution to data service provides selected by the project who support BMS RI data management, analyses, deposition and distribution

Key components

Authentication and authorisation infrastructure (AAI)

Secure data streaming

Metadata standardisation and synchronisation

Policy components linking to country level/Infrastructure level best practice

Improve interoperability with European e-infrastructures and leverage existing investments in these capacities within the biomedical and life science domain

CORBEL driver projects

Delegated access to digitalized biobank samples

GoNL - Federated AAI + secure streaming

BBMRI-NL - Bioschemas and Beacons

BBMRI-ERIC/RD-connect - harmonization and 'matchmaking' service
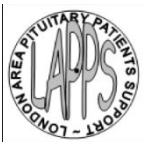
**molgenis**

Hosts many databases, amongst which many patient and mutation registries.

**deb register**
the international database of dystrophic epidermolysis bullosa patients and *COL7A1* mutations
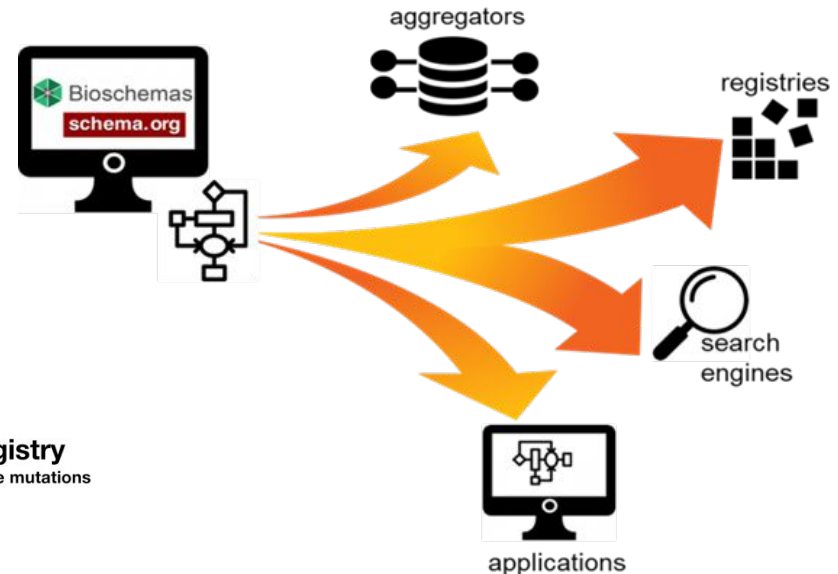
**International Microvillus Inclusion Disease (MVID) Patient Registry**
International registry of patients with microvillus inclusion disease and database of associated gene mutations

AIP Mutation Database

CHD7 Database

Added structured machine readable metadata descriptions to multiple patient registries using BioSchemas.
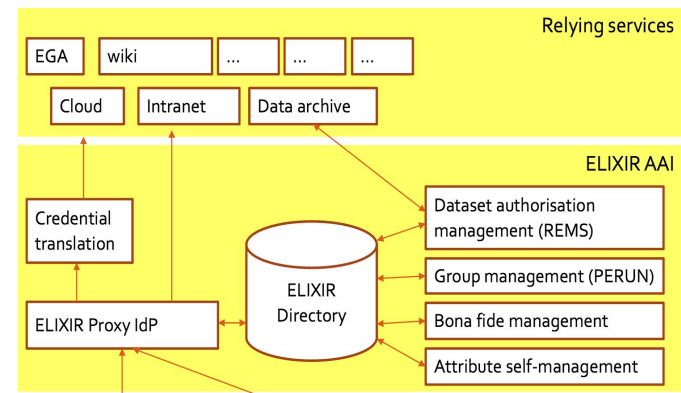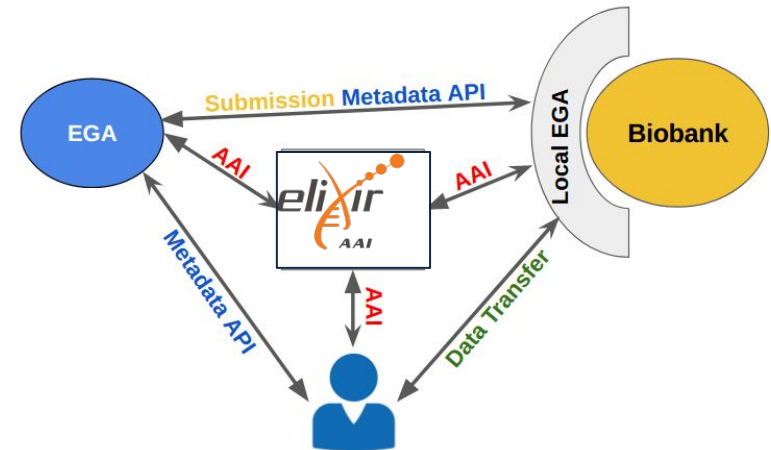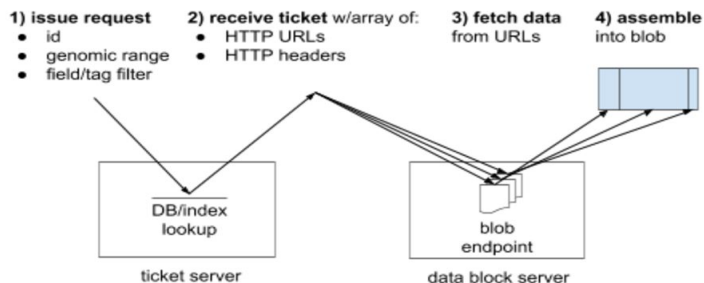
Increase the findability of multiple patient and mutation registries, making it easier to find and reuse critical information for rare disease patient care and research.

https://f1000research.com/posters/7-1228

- User authentication and authorisation
  - Align with ELIXIR AAI
- Dataset authorisation
  - REMS: Electronic tool for the management of access rights to controlled access research data
  - Beacon: automated process to all data in "registered" data access layer
- Secure data delivery
  - Align with GA4GH
  - htsget: remote streaming protocol







htsget (http://samtools.github.io/hts-specs/htsget.html)

# TASK 6.3 DATA MATCHMAKING & HARMONIZATION TOOLKIT FOR POOLED ANALYSIS OF SENSITIVE DATA

## Heterogeneous data sharing across 20+ child cohorts - the LifeCycle Project

| Find data collections matching your needs | Match data items to research parameters and create algorithms | Apply harmonization algorithms and then automated (meta) analysis |
|---|---|---|
| **new** | **FAIRifier** | |
| **BiobankUniverse** | **BiobankConnect** | **Secure Digital Research Environment** |
| *Find 'similar' collections based on data item metadata only, using ontology based lexical and semantic annotation and matching (6.2)* | *Shortlist attributes matching research needs and auto-generate ETL algorithms; provide sharing of harmonization rules* | *Apply ETL algorithms on the data in secure analysis environment following DAC, e.g. analysis of 250.000 children (LifeCycle)* |

All available as open source http://molgenis.githubio. BiobankUniverse: manuscript in prep; beta@ http://biobankuniverse.org. BiobankConnect:PMID:27153686 ; demo @ http://biobankconnect.org .

**BBMRI**
Biobanking and
Biomolecular
Resources Research
Infrastructure

elixir

Virtual HPC
cloned across
sites

Deliver
reproducibility

**WHAT**

>30 million
PEOPLE IN EU are affected
by rare diseases

Rare disease
projects

GDPR compliance
use cases
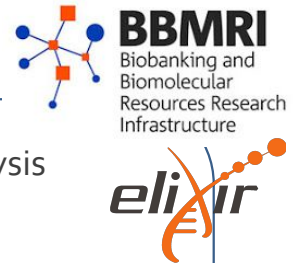
Secure Federated
analysis solution

**WHY**

Openstack API

Ansible installer
Plug and play workflow

Embassy cloud
deployment –
proximity to EGA

**HOW**

Federated analysis
enabled

Close to data archives at
EBI, UMCG, NIKHEF

19K+165k individual's
data

Broaden access to tools
and compute

**IMPACT**

"We developed a Bioinformaticians Sandbox - a fully virtual HPC cluster for bioinformaticians which can be automatically and reproducibly cloned on different sites, ensuring reproducibility when analyses are done on multiple sites. This was motivated by a large multi-center study to have a controlled data access site with analysis capability that is GDPR-Compliant"

Morris Swertz

# CONTINUITY FOR EOSC LIFE

Identifier best practice and checklists

Semantic infrastructure and ontologies
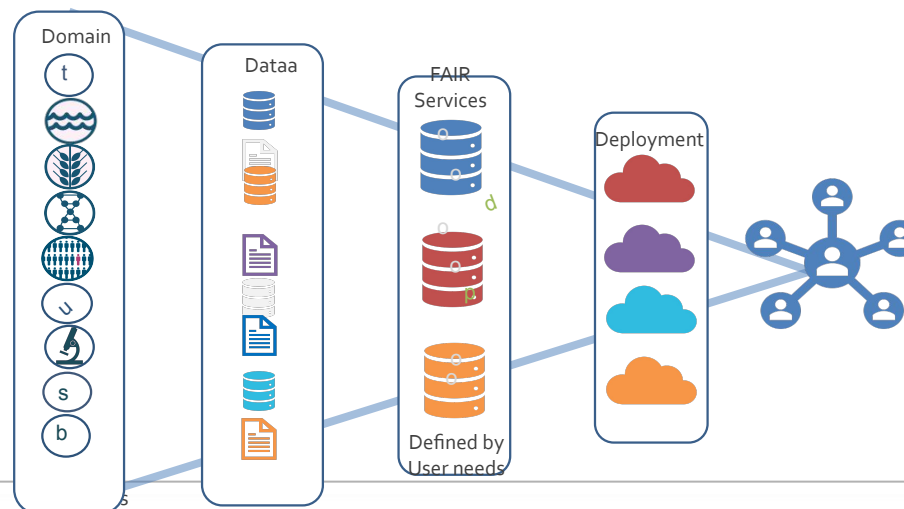
Secure data access technology

Pan infrastructure deployment / enhancement

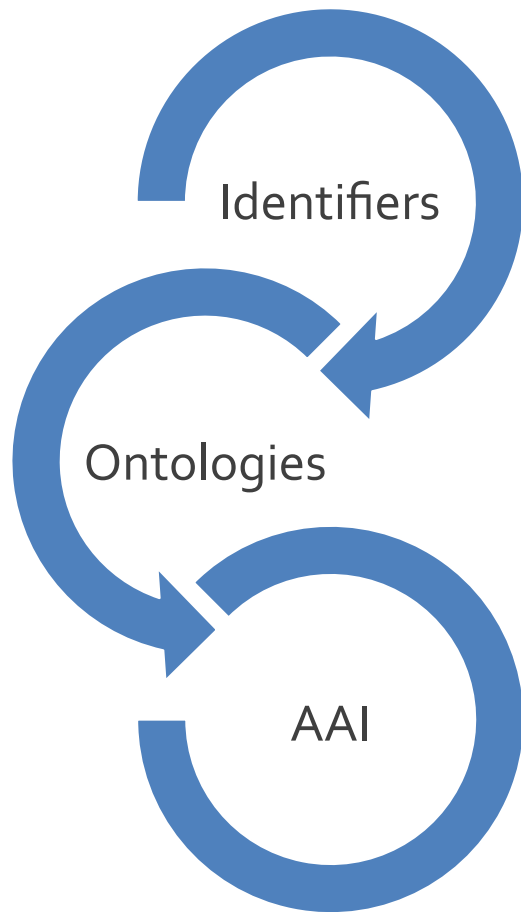**Standards driving FAIR tool suite**

**FAIR Tool suite in the cloud**

**Enabling AI on controlled access data**

**Deployment via Community Funding Calls**

Identifiers

Ontologies

AAI

Within and cross BMS RI – open projects, standards and implementations, and better understanding.

ELIXIR, BBMRI, EuroBioimaging, EMBRC, ISBE

Broader utility than implementations within CORBEL – now seeing use in many projects

Many components now being taken forward to EOSC Life where they will form tech stack for the next phase of projects

Getting infrastructure out and across BMS RI a challenge.

Benefited from cooperation with ELIXIR Excelerate