# FAIR Principles practices supporting recognition
## (experiences from ELIXIR and FAIRDOM)

## Professor Carole Goble
## The University of Manchester, UK



Head of UK Node
Co-lead Interoperability
Platform

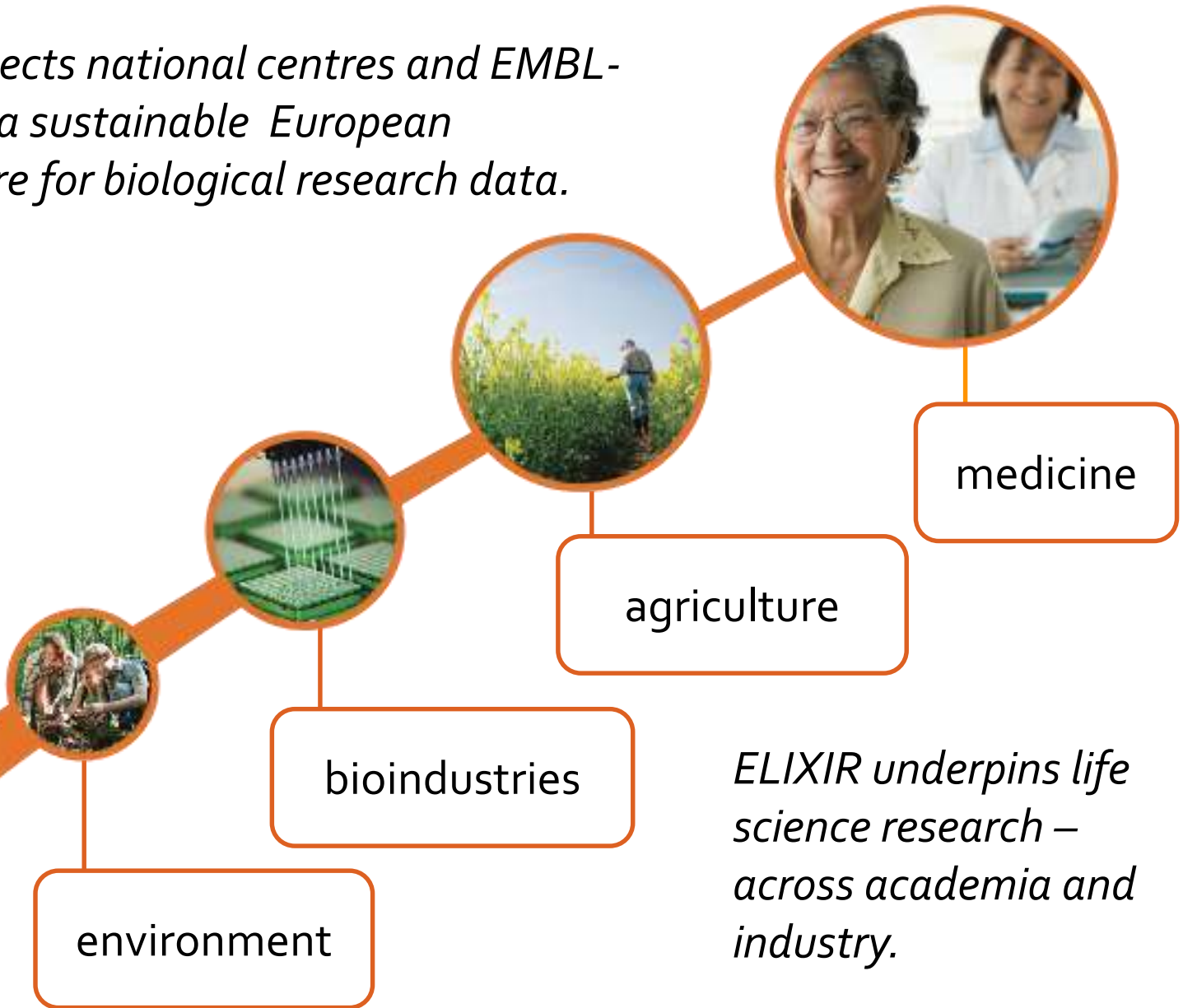Co-ordinator
Chair FAIRDOM
Association e.V.

Leadership team

DCIP Advisor
Advisory Board

UK Biobanking Showcase 2017, London, 18 October 2017

*ELIXIR connects national centres and EMBL-EBI to build a sustainable European infrastructure for biological research data.*

medicine

agriculture

bioindustries

*ELIXIR underpins life science research – across academia and industry.*

environment

## European Research Infrastructure

Life Sciences Data
Standards, Portals, Platforms,
Lobbying, Sustain Core Resources
Archive Level: e.g. BioSamples

22 Countries

## Systems and Synthetic Biology Projects

Management for Research
Data, Operations, Models
Assets at the Project level
e.g. Samples, Strains,
Specimens

80+ Projects
700+ Researchers

# Recognition: being FAIR
## Find – Download – Go



Data, model, SOP, sample… provider

make it easier to be found and to track credit

Data, model, SOP, sample… user

make it easier to find and to action credit

# Getting recognition

## Software in the scientific literature: Problems with seeing, finding, and using software mentioned in the biology literature

James Howison, Julia Bullard

### Abstract

Software is increasingly crucial to scholarship, yet the visibility and usefulness of software in the scientific record are in question. Just as with data, the visibility of software in publications is related to incentives to share software in reusable ways, and so promote efficient science. In this article, we examine software in publications through content analysis of a random sample of 90 biology articles. We develop a coding scheme to identify software "mentions" and classify them according to their characteristics and ability to realize the functions of citations. Overall, we find diverse and problematic practices: Only between 31% and 43% of mentions involve formal citations; informal mentions are very common, even in high impact
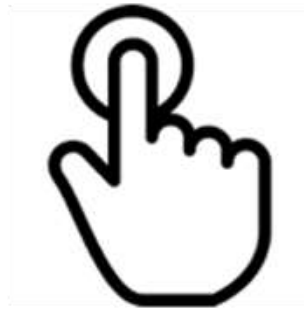
37% formal citations
mostly "mentions"

## Database Citation in Full Text Biomedical Articles

Şenay Kafkas*, Jee-Hyub Kim, Johanna R. McEntyre

European Molecular Biology Laboratory – European Bioinformatics Institute Wellcome Trust Genome Campus, Cambridge, United Kingdom

### Abstract

Molecular biology and literature databases represent essential infrastructure for life science research. Effective integration of these data resources requires that there are structured cross-references at the level of individual articles and biological records. Here, we describe the current patterns of how database entries are cited in research articles, based on analysis of the full text Open Access articles available from Europe PMC. Focusing on citation of entries in the European Nucleotide Archive (ENA), UniProt and Protein Data Bank Europe (PDBe), we demonstrate that text mining doubles the number of structured annotations of database record ci... literature-database relationships are found by t... cited by database records. We recommend th... databases, such as ArrayExpress and Pfam, ent... and high-throughput of this text-mining pipe... allow the development of new integrated dat...

### Review

## Reuse of public genome-wide gene expression data

Johan Rung & Alvis Brazma

Our understanding of gene expression has changed dramatically over the past decade, largely catalysed by technological developments. High-throughput experiments — microarrays and next-generation sequencing — have generated large amounts of genome-wide gene expression data that are collected in public archives. Added-value databases process, analyse and annotate these data further to make them accessible to every biologist. In this Review, we discuss the utility of the gene expression data that are in the public domain and how researchers are making use of these data. Reuse of public data can be very powerful, but there are many obstacles in data preparation and analysis and in the interpretation of the results. We will discuss these challenges and provide recommendations that we believe can improve the utility of such data.

*25% Publications that used the public Arrayexpress Archive cited it*

F <span style="color:#5b7ca3">Credit</span>
indable

A <span style="color:#5b7ca3">Track</span>
ccessible

I <span style="color:#5b7ca3">Understand</span>
nteroperable

R <span style="color:#5b7ca3">Reproduce</span>
eusable

The FAIR Guiding Principles for scientific data management and stewardship
https://www.nature.com/articles/sdata201618 (2016)

# FAIR Data Principles

*Access to public funded research, Reproducible results*
*Value and CREDIT all research outputs*

# SCIENTIFIC DATA

## Comment: The FAIR Guiding Principles for scientific data management and stewardship

Mark D. Wilkinson *et al.*#

There is an urgent need to improve the infrastructure supporting the reuse of scholarly data. A diverse set of stakeholders—representing academia, industry, funding agencies, and scholarly publishers—have come together to design and jointly endorse a concise and measureable set of principles that we refer to as the FAIR Data Principles. The intent is that these may act as a guideline for those wishing to enhance the reusability of their data holdings. Distinct from peer initiatives that focus on the human scholar, the FAIR Principles put specific emphasis on enhancing the ability of machines to automatically find and use the data, in addition to supporting its reuse by individuals. This Comment is the first formal publication of the FAIR Principles, and includes the rationale behind them, and some exemplar implementations in the community.

https://www.nature.com/articles/sdata201618 (2016)

## Supporting discovery through good data management

Good data management is not a goal in itself, but rather is the key conduit leading to knowledge discovery and innovation, and to subsequent data and knowledge integration and reuse by the community after the data publication process. Unfortunately, the existing digital ecosystem surrounding scholarly data publication prevents us from extracting maximum benefit from our

EUROPEAN COMMISSION
Directorate-General for Research & Innovation

## H2020 Programme

Guidelines on

FAIR Data Management in Horizon 2020

Realising
the European
Open Science Cloud

First report and recommendations
of the Commission High Level Expert Group
on the European Open Science Cloud

EUROPEAN COMMISSION
DIRECTORATE-GENERAL FOR RESEARCH & INNOVATION

The Director-General

Brussels, 10 July 2017

### EOSC Declaration

RECOGNISING the challenges of data driven research in pursuing excellent science;

GRANTING that the vision of European Open Science is that of a research data commons, widely inclusive of all disciplines and Member States, sustainable in the long-term,

CONFIRMING that the implementation of the EOSC is a process, not a project, by its nature iterative and based on constant learning and mutual alignment;

UPHOLDING that the EOSC Summit marked the beginning and not the end of this process, one based on continuous engagement with scientific stakeholders, the European Commission,

PROPOSES that all EOSC stakeholders consider sharing the following intents and will actively support their implementation in the respective capacities:

### Data culture and FAIR data

➢ **[Data culture]** European science must be grounded in a common culture of data stewardship, so that research data is recognised as a significant output of research and is appropriately curated throughout and after the period conducting the research. Only a considerable cultural change will enable long-term reuse for science and for innovation of data created by research activities: no disciplines, institutions or countries must be left behind.

➢ **[Open access by-default]** All researchers in Europe must enjoy access to an open-by-default, efficient and cross-disciplinary research data environment supported by FAIR data principles. Open access must be the default setting for all results of publicly funded research in Europe, allowing for proportionate limitations only in duly justified cases of personal data protection, confidentiality, IPR concerns, national security or similar (e.g. 'as open as possible and as closed as necessary').

➢ **[Skills]** The necessary skills and education in research data management, data stewardship and data science should be provided throughout the EU as part of higher education, the training system and on-the-job best practice in the industry. University associations, research organisations, research libraries and other educational brokers play an important role but they need substantial support from the European Commission and the Member States.

➢ **[Data stewardship]** Researchers need the support of adequately trained data stewards. The European Commission and Member States should invest in the education of data stewards via career programmes delivered by universities, research institutions and other trans-European agents.

➢ **[Rewards and incentives]** Rewarding research data sharing is essential. Researchers who make research data open and FAIR for reuse and/or reuse and reproduce data should be rewarded, both

UK Funder Data Policies
http://www.dcc.ac.uk/resources/policy-and-legal/overview-funders-data-policies

# Machine Processability for Distributed Systems

The FAIR Guiding Principles for scientific data management and stewardship
https://www.nature.com/articles/sdata201618 (2016)

search
catalogues

stores

F*Credit* indable A*Track* ccessible I*Understand* nteroperable R*Reproduce* eusable

policy,
authorised
access,
licensing

standards: identifiers, metadata
that is machine processable

# 1. Identifiers and Citation

## Best Practices & Principles

PERSPECTIVE

Identifiers for the 21st century: How to design, provision, and reuse persistent identifiers to maximize utility and impact of life science data

https://doi.org/10.1371/journal.pbio.2001414

#CiteTheData campaigns

$DC^1$
Data Citation Principles

<JATS>

https://www.force11.org/group/joint-declaration-data-citation-principles-final

## Data Citation Implementation

## Identifier Schemes

Resource Identification Initiative

doi

https://www.force11.org/group/resource-identification-initiative

Identifiers.org

DataCite
FIND, ACCESS, AND REUSE DATA

Impactstory

## Services

# Identifiers for 21st Century

1. Credit any derived content using its original identifier
2. Help local IDs travel well: Document prefix and patterns
3. Opt for simple, durable web resolution
4. Avoid embedding meaning or relying on it for uniqueness
5. Design new identifiers for diverse uses by others
6. Implement a version-management policy
7. Do not reassign or delete identifiers
8. **Make URIs clear and findable**
9. Document the identifiers you issue and use
10. Reference and display responsibly

# Data Citation Principles

- 1. Importance
- **2. Credit and Attribution** Data citations should facilitate giving scholarly credit and normative and legal attribution to all contributors to the data
- 3. Evidence
- **4. Unique Identification** A data citation should include a persistent method for identification that is machine actionable, globally unique, and widely used by a community.
- **5. Access** Data citations should facilitate access to the data themselves and to such associated metadata, documentation, code, and other materials, as are necessary for both humans and machines to make informed use of the referenced data.
- 6. Persistence
- 7. Specificity and Verifiability
- 8. Interoperability and Flexibility

# 2. Find and Access



Schema.org adds simple **structured metadata markup** to web pages & sitemaps for harvesting, search and summary snippet making.

Search engines often highlight websites containing Schema.org

Widespread commercial and open source infrastructure -> low barrier to adoption.

No need for an API or special feeds.

schema.org tailored to the Biosciences

# From Potato Salad to Protein Annotation

```
{
"@context": "http://schema.org",
"@type": "BiologicalEntity",
"@id":
"http://www.identifiers.org/uniprot/P00519
",

"biologicalType": "protein",
"isMentionedIn": {
   "@type": "Dataset",
   "@id":
"http://www.uniprot.org/news/2017/03/15/re
lease"
},

"associatedDisease": {
   "@type": "MedicalCondition",
   "@id":
"http://www.omim.org/entry/608232",
   "name": "Leukemia, chronic myeloid
(CML)",
   "code": {
      "@type": "MedicalCode",
      "code": "608232",
      "codingSystem": "OMIM"
   },
   "sameAs":
"http://www.uniprot.org/diseases/DI-03735"
},
"biocoordinates": {
   "@type": "QuantitativeValue",
   "value": "1130"
},
```



```
"taxon": "http://www.uniprot.org/taxonomy/9606",

"alternateName": [
   {
      "@language": "en",
      "@value": "ABL1_HUMAN"
   },
   {
      "@language": "en",
      "@value": "Proto-oncogene c-Abl"
   }
],
"description": {
   "@language": "en",
   "@value": "Non-receptor tyrosine-protein kinase that plays a role..."
},
"identifier": "http://www.identifiers.org/uniprot/P00519",
"image":
"http://www.identifiers.org/uniprot/P00519#showFeaturesViewer",
```
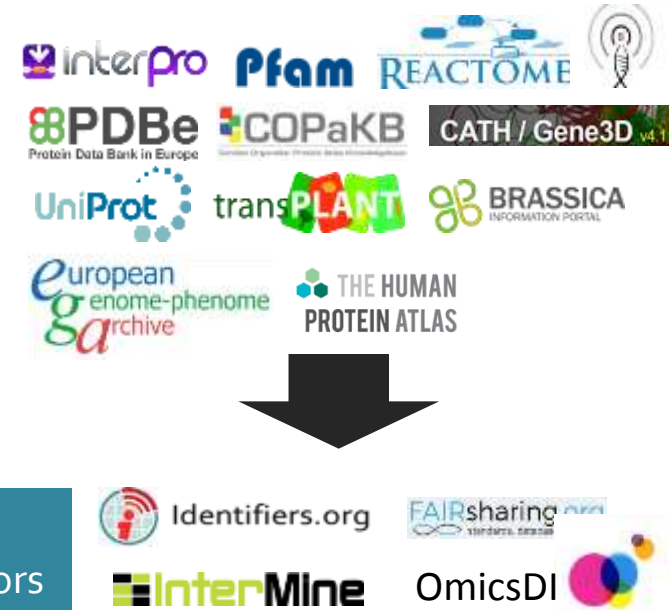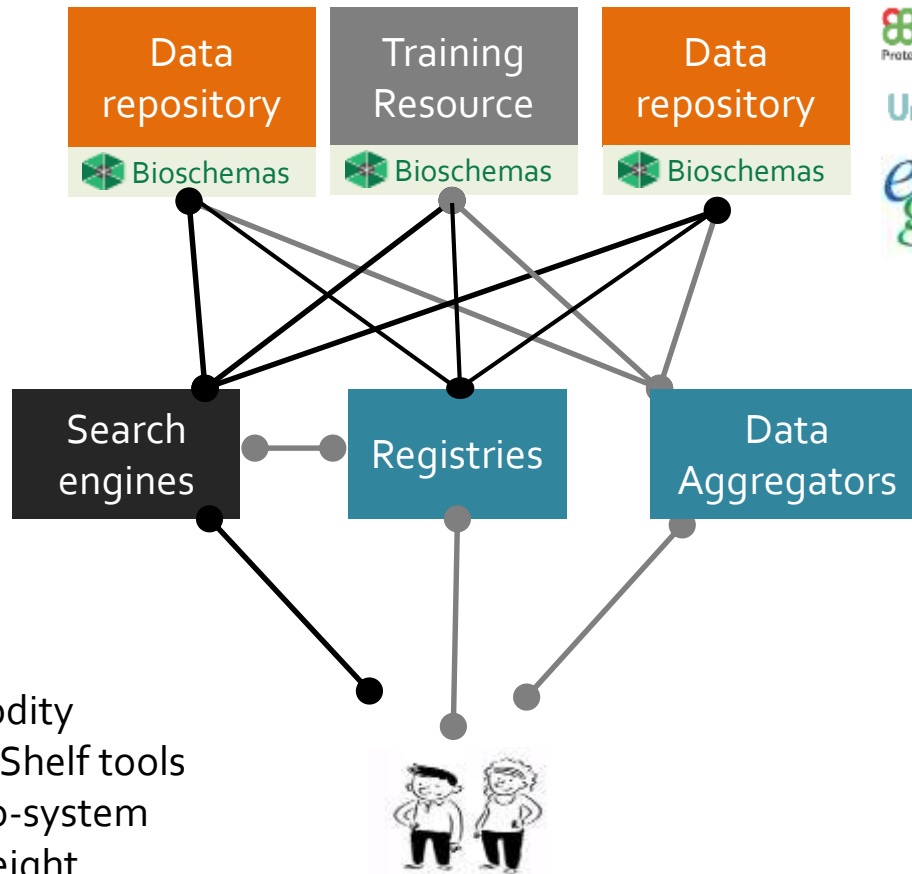
# Bioschemas.org

*simple **structured metadata markup** on web pages & sitemaps tailored to the Biosciences*

Standardised metadata mark-up

RDFa
JSON-LD
**Microdata**

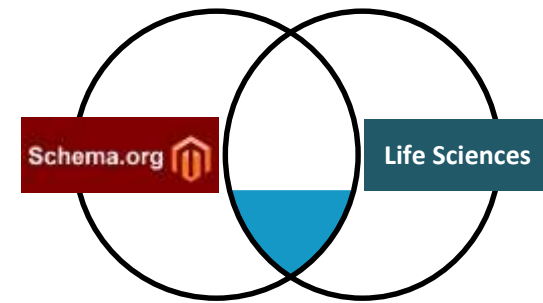Metadata published and harvested without APIs or special feeds

Commodity
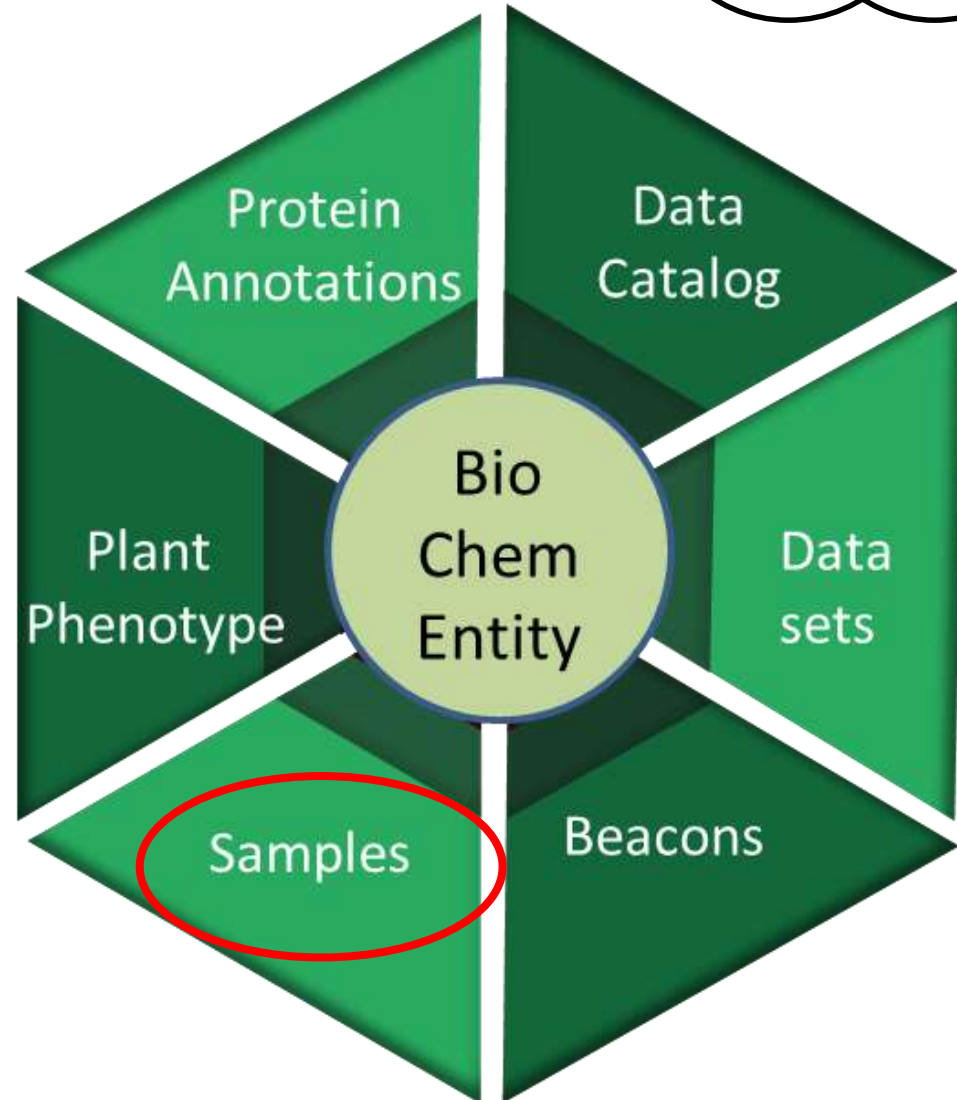Off the Shelf tools
App eco-system
Lightweight

# Bioschemas.org

*simple **structured metadata markup** on web pages & sitemaps tailored to the Biosciences*

First specifications:

- Bio data infrastructure
    - *DataCatalog*
    - *Datasets*

- Bio data types
    - *Human beacons*
    - ***Samples***
    - *Plant Phenotypes*
    - *Proteins*
    - *(Chemistry)*

- Bio stuff
    - *Training materials*
    - *Events*
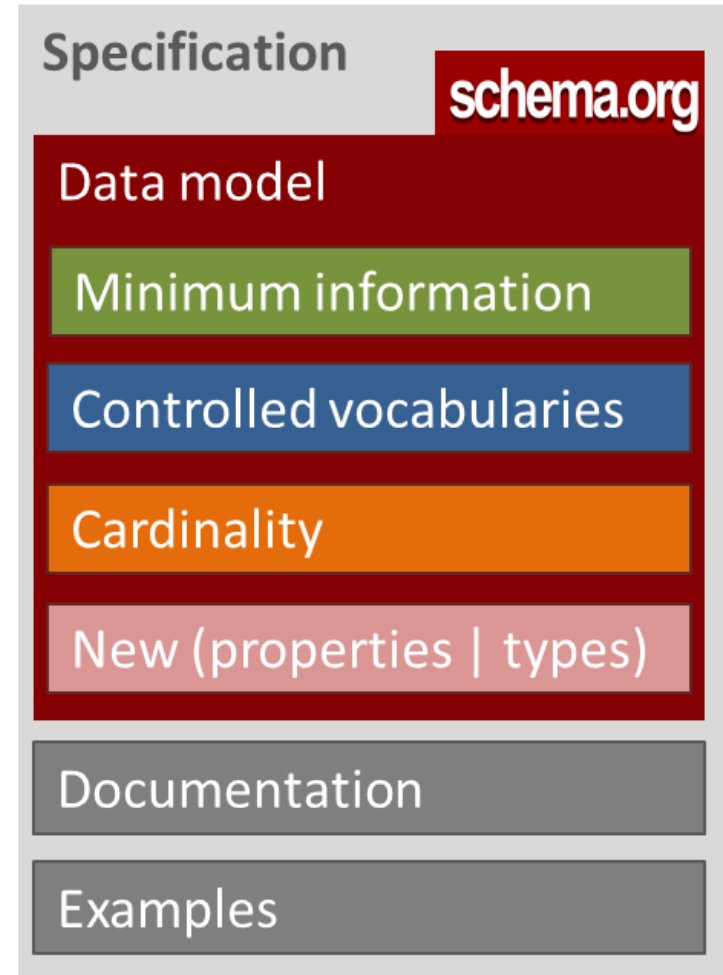    - *Laboratory protocols*
    - *(Workflows and Tools)*

# Bioschemas.org

*Tailored schema.org to improve*
***Findability*** *and* ***Accessibility*** *in Bioscience*

- # Specification on top of **schema.org**
- Introduce bioscience types
- Restricted use case
  - Finding data
  - Presenting search results
  - Metadata exchange
- Minimum properties – 6
- Link to domain ontologies not reinventing them

Specification

schema.org

Data model

Minimum information

Controlled vocabularies

Cardinality

New (properties | types)

Documentation

Examples

Layer of constraints +

documentation + extensions

- Name
- Description
- License
- Release
- Citation
- Metrics
- Tools
- …

# DataCatalog

# Bioschemas.org
*simple **structured metadata markup** on web pages & sitemaps tailored to the Biosciences*

Standardised metadata mark-up

RDFa
JSON-LD
**Microdata**

Metadata published and harvested without APIs or special feeds

Commodity
Off the Shelf tools
App eco-system

Data repository — Bioschemas

Training Resource — Bioschemas

Data repository — Bioschemas

Search engines

Registries

Data Aggregators

health RI research infrastructure

BBMRI.nl Biobanking and BioMolecular resources Research Infrastructure

HipSci

FAANG

BIOBANK

BIOBANK

EBiSC

MarRef

BioSamples

Tissue Directory and Coordination Centre

molgenis

RD Connect Sample Catalogue

BBMRI-ERIC Directory

# Samples

# BioSamples at the EBI



BioSamples stores and supplies descriptions and metadata about biological samples.

- 'reference' samples (e.g. from 1000 Genomes, HipSci, FAANG)
- used in an assay database such as the European Nucleotide Archive (ENA) or ArrayExpress.

BioSamples connects across resources

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

# Visibility & Credit – Find, Access, and Propagate my biobank metadata
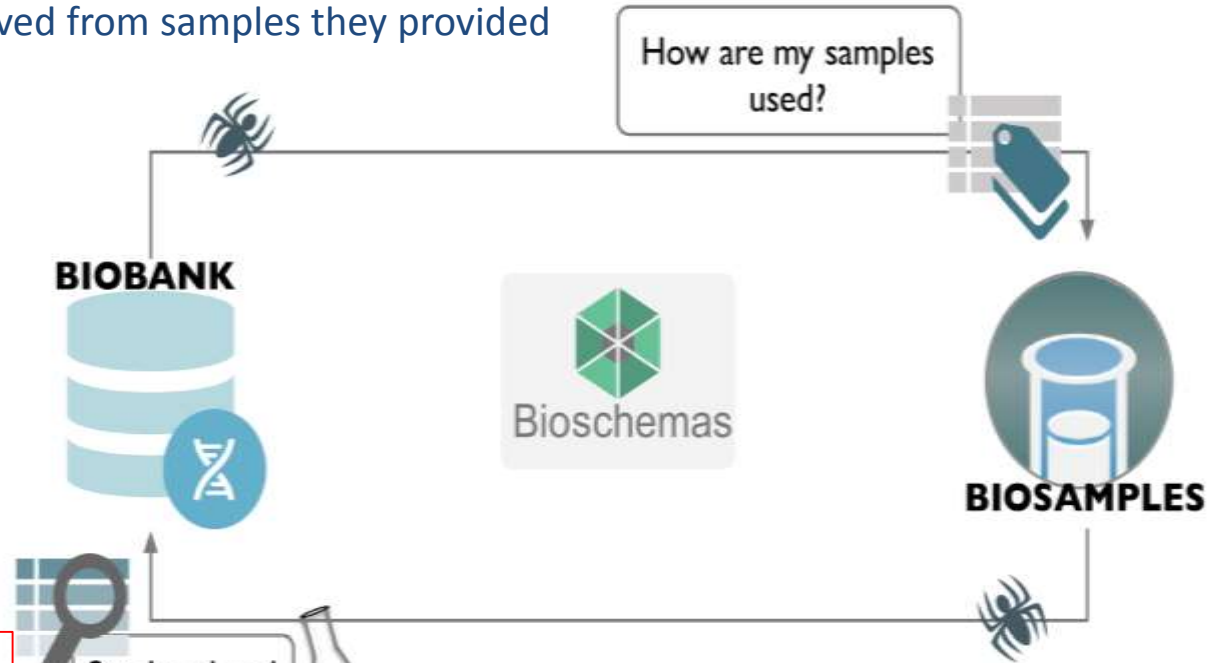
**Biobanks crawl BioSamples**
identify all the published (and searchable) datasets derived from samples they provided

Biobanks ensure only authorised metadata is visible & control access to restricted samples.

Only Biobanks know the specific samples connected to publicly available datasets

Help ingest sample metadata from data repositories (e.g. Biobank databases) into registries like BioSamples, BBMRI-ERIC Directory, the UKCRC Tissue Directory



How are my samples used?

BIOBANK

Bioschemas

BIOSAMPLES

Send updated metadata

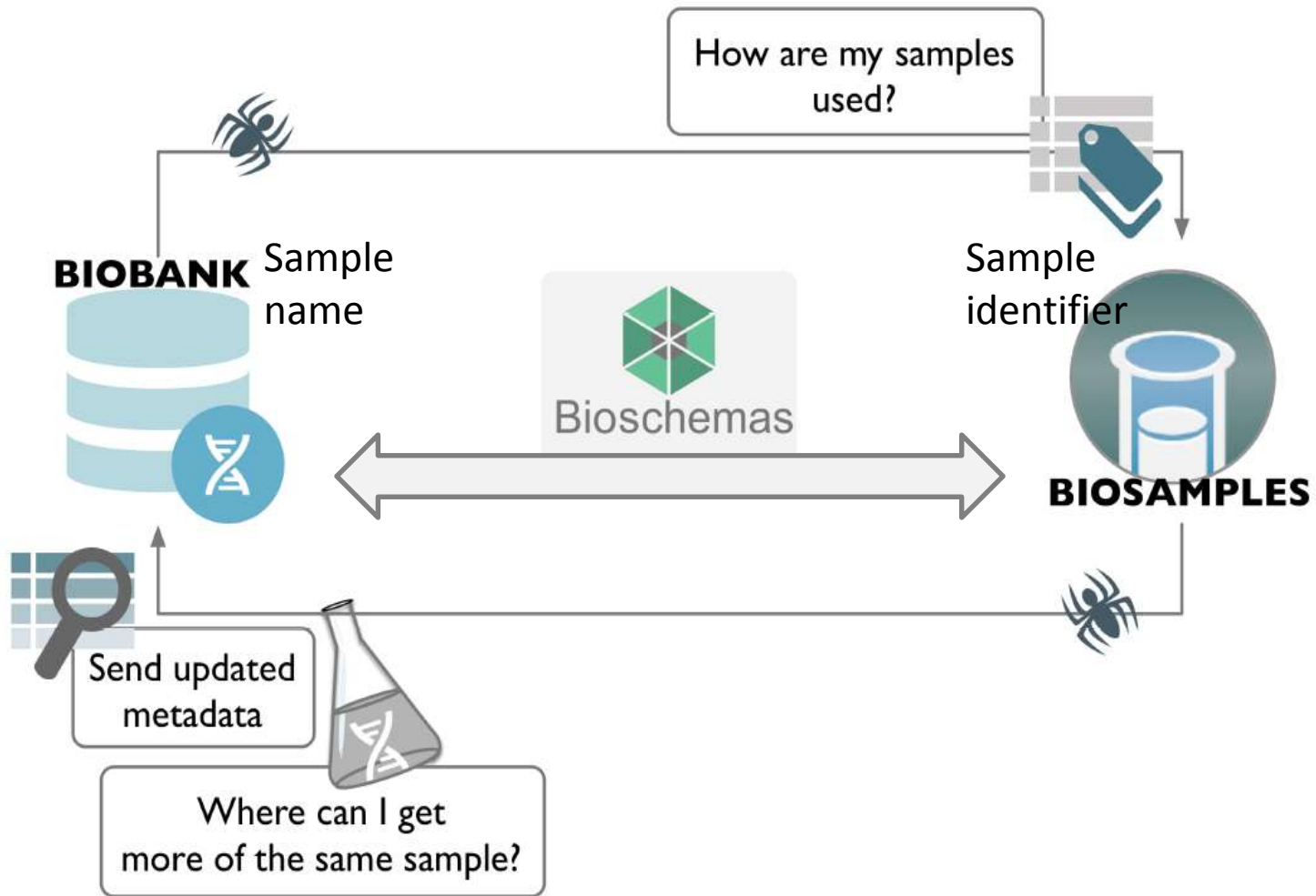Where can I get more of the same sample?

**BioSamples crawl Biobank websites**
- identify samples that are known to have public accessions in the BioSamples database *and* can be made publicly available
- link public samples to a provider ("where can I get more of this sample?").

# Visibility & Credit
## maintain a sanitised bi-directional link through identifiers

# Visibility & Credit – Find, Access, and Propagate Biobank metadata

- Customers
  - Direct: Technical providers of biobank and directory platforms
  - Indirect: The users, the Biobanks, the Directories
- Publishing and harvesting
  - Retrospective, prospective and lightweight mechanism easy to incorporate
- Scalable & sustainable
  - Exploits common web infrastructure
  - Standard Search engines can index

**BIOSAMPLES**

**BIOBANK**

Lightweight, hidden mark-up in the web page

# Open Public Process

MIABIS: Minimum Information About BIobank data Sharing (version 2.0)

schema.org

Find, Cite, Snippets, Metadata exchange

Ideally 6 concepts
Reuse ontologies

Real mark-up
Tools

Use cases          Mapping          Specification          Adoption



*2-3 Oct 2017, Hinxton, ~50 people*

Testing          Application

# Open Community driven Specifications

## Profiles

The Bioschemas profiles define a community agreed layer over the Schema.org model providing additional constraints. These constraints capture (i) the minimal information properties agreed by the community which are mandatory (M), recommended (R), or optional (O), (ii) the cardinality of the property, i.e. whether it is expected to occur once or many times, and (iii) associated controlled vocabulary terms drawn from existing ontologies.

| Name | Short description | Version | Group | Spec. | Folder | Mapping | Use Cases | Task & Issues | Examples |
|------|-------------------|---------|-------|-------|--------|---------|-----------|---------------|----------|
| Beacon | A convention for beacon to self-describe. | 0.2 | | | | | | | |
| DataCatalog | Bioschemas specification for describing data repositories, data registries, and data catalogues in the life-sciences. | 0. | | | | | | | |
| Dataset | Bioschemas specification for describing Dataset | 0. | | | | | | | |

| | | | | | | | | | |
|------|-------------------|---------|-------|-------|--------|---------|-----------|---------------|----------|
| Protein | Bioschemas specification describing a Protein (PhysicalEntity profile) in Life Sciences | 0.4 | | | | | | | |
| Sample | Bioschemas specification for describing Samples in the life-science. | 0.1 | | | | | | | |
| Standard | Bioschemas specification for describing Standards in the life-science. | 0.1 | | | | | | | |
| Tool | Bioschemas specification for describing SoftwareApplication in the life-science. | 0.1 | | | | | | | |
| TrainingMaterial | A specification for describing training | 0.2 | | | | | | | |

# Take Home

Data becomes findable and accessible to the biobank

F *Credit* indable  A *Track* ccessible  I *Understand* nteroperable  R *Reproduce* eusable

Create the framework for finding and accessing data through identifiers and metadata exchange using Bioschemas.org

Policies and incentives to get journals, data repositories, funders and authors to mark up and cite.

# Acknowledgements
## supportive projects and communities

http://elixir-europe.org

Bioschemas.org
http://bioschemas.org

http://fair-dom.org

FORCE11
The Future of Research Communications and e-Scholarship

bioCADDIE    http://biocaddie.org

http://force11.org