# Data standards, sharing and publication

# in the life sciences

Susanna-Assunta Sansone, PhD

*Associate Director,*
*Principal Investigator*

*Board of Directors*

*Data Consultant,*
*Honorary Academic Editor*

OXFORD e-Research CENTRE    UNIVERSITY OF OXFORD

DRYAD

npg nature publishing group
SCIENTIFIC DATA

# ODIN mission

ODIN will build on the ORCID and DataCite initiatives to uniquely identify scientists and data sets and connect this information across multiple services and infrastructures for scholarly communication. It will address some of the critical open questions in the area:

- Referencing a data object

- Tracking of use and re-use

- Links between a data object, subsets, articles, rights statements and every person involved in its life-cycle.

# Outline of my talk

*Problem*:
Identification of datasets in pivotal. But meaningful sharing and (re)use also depend on how well described the datasets are.

*Status quo:*
In the life sciences there is a wealth of 'reporting standards' set to enhance and facilitate the experimental descriptions.

*Challenges:*
Identify 'reporting standards' and their organizations, track their use, usability and impact (e.g. linking them to datasets), credit their developers, users (e.g. curators)...

# My team's activities and groups we work with

**data management, biocuration and publication,
collaborative development of software, database, standards and ontology**

- environmental genomics
- metabolomics
- metagenomics
- nanotechnology
- proteomics

- stem cell discovery
- system biology
- transcriptomics
- toxicogenomics
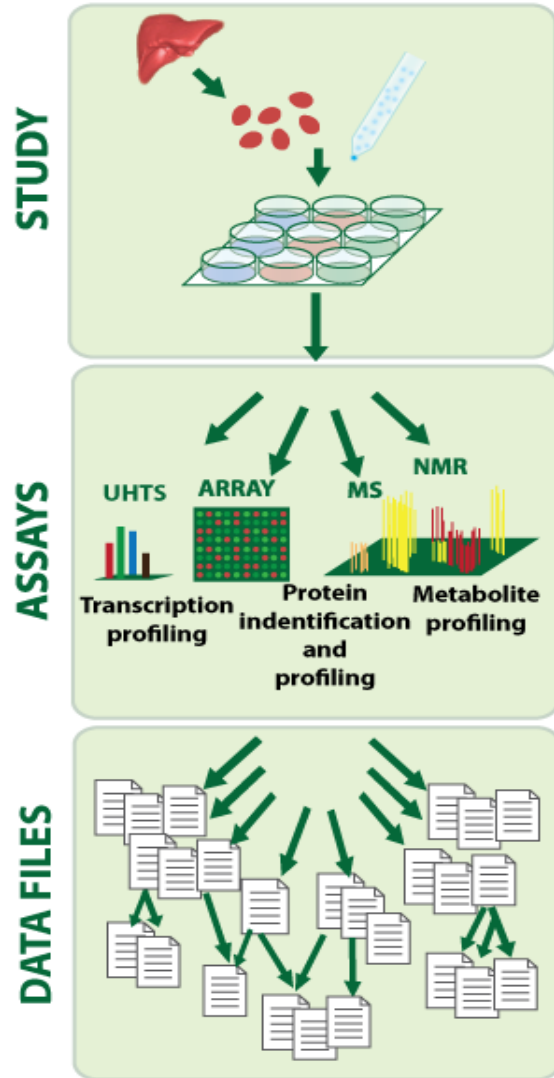- environmental health



env

agro

tox/pharma

health

COMPREHENSIVE
INTEROPERSI
REUSABLE

# Growing movement for reproducible research



*Example* of experimental WORKFLOW

STUDY

ASSAYS

UHTS  ARRAY  MS  NMR

Transcription profiling  Protein indentification and profiling  Metabolite profiling
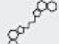
DATA FILES

- **Researchers** and **bioinformaticians** in both *academic* and *commercial* arenas, along with **funding agencies** and **publishers**, embrace the concept that to be *comprehensible*, *interoperable* and *reusable* shared datasets we should have **richly described**:

  - entities of interest

    *e.g.*, genes, metabolites, phenotypes, computational models, diseases ...

  - experimental steps

    *e.g.*, provenance of study materials, technology and measurement types, experimentalists and curators ...

**The *necessity* for well-annotated data and unambiguous experimental metadata was especially apparent**
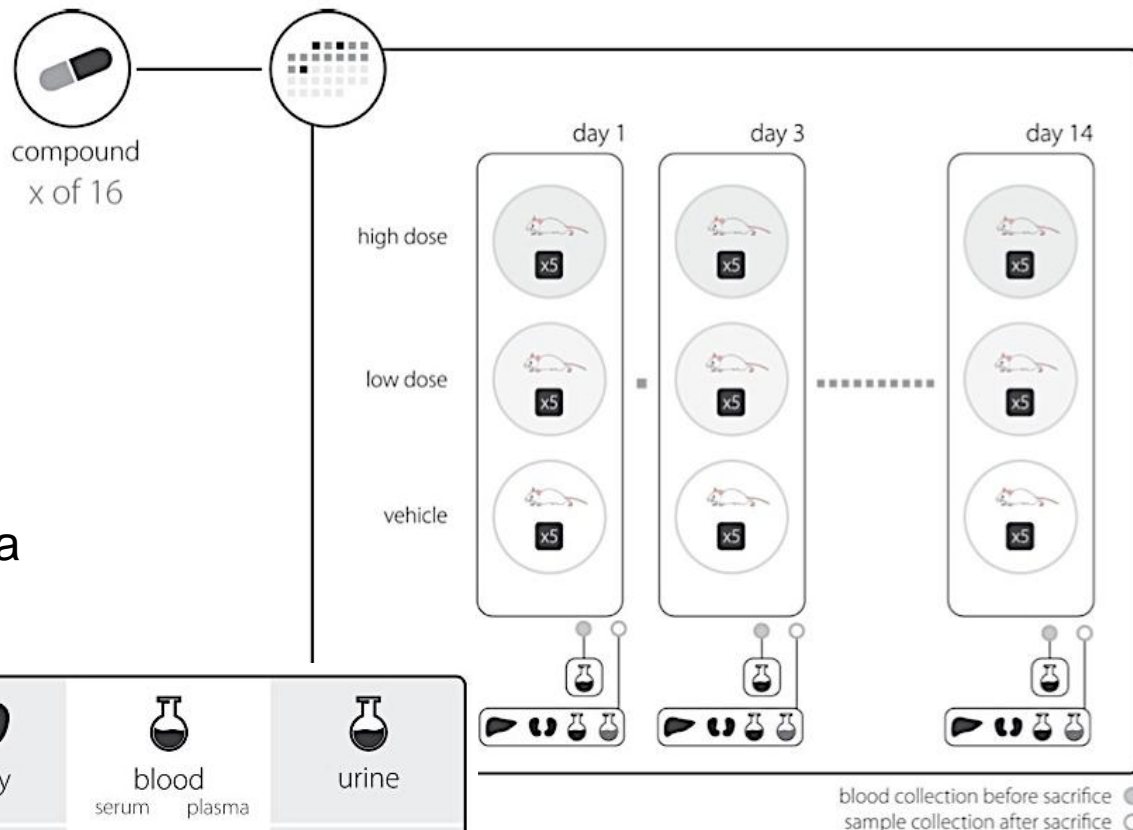
- during cross-study comparisons and data analysis
- in preparation for reformatting the datasets for submission to the different EBI repositories, requiring different level of information

compound
x of 16

day 1          day 3          day 14

high dose      x5             x5             x5

low dose       x5             x5             x5

vehicle        x5             x5             x5

blood collection before sacrifice ●
sample collection after sacrifice ○

|  | liver | kidney | blood serum plasma | urine |
|---|---|---|---|---|
| protein expression profiling by mass spectrometry | ✔ | ✔ | ✔ | ✔ |
| transcription profiling by dna microarray | ✔ | ✔ | ✔ ✔ | |
| metabolite profiling by mass spectrometry | ✔ | ✔ | ✔ | ✔ |
| metabolite profiling by nmr spectroscopy | ✔ | ✔ | ✔ | ✔ |
| histology | ✔ | ✔ | | |
| clinical chemistry | | | ✔ ✔ | ✔ |
| hematology | | | ✔ ✔ | |

experimental design

sample characteristic(s)

experimental variable(s)

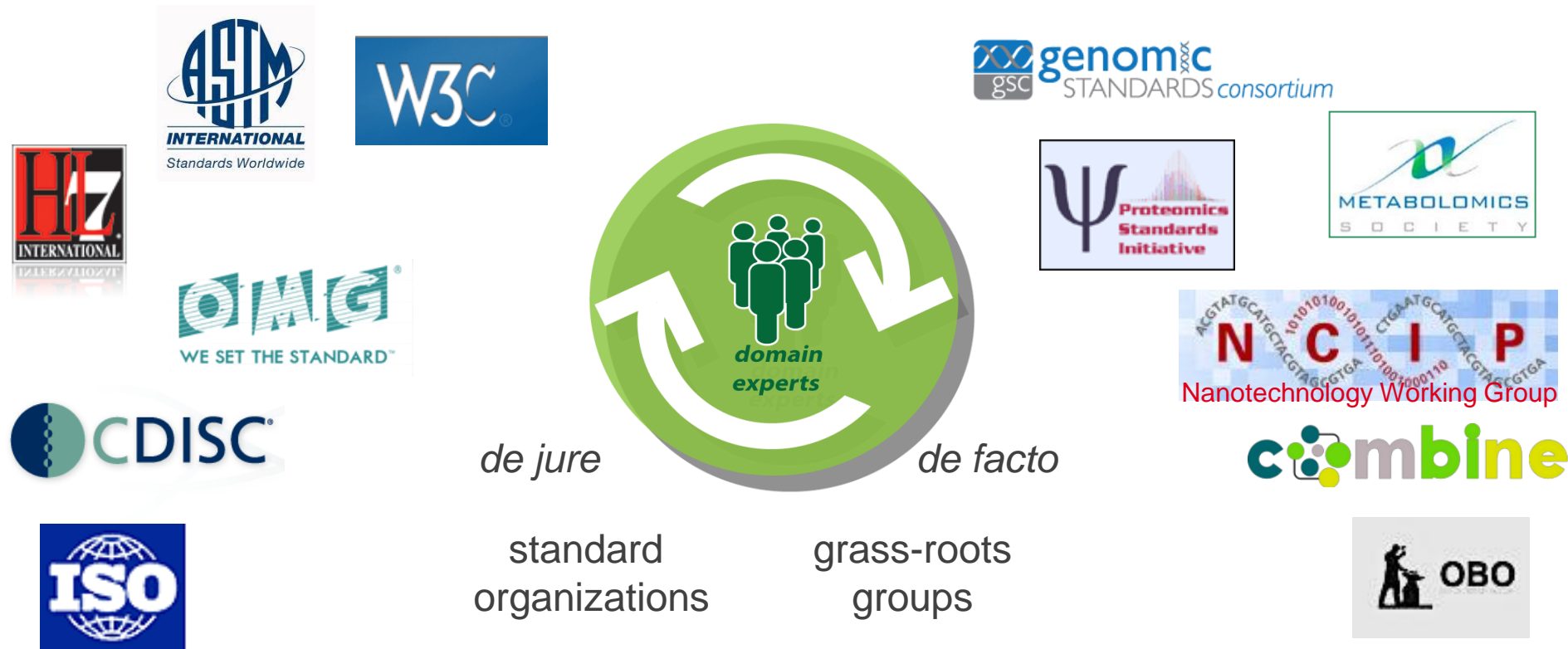technology(s)

measurement(s)

protocols(s)

data file(s)

- One must strike a balance between
    - depth and breadth of information; and
    - sufficient information required to reuse the data

- Capture all salient features of the experimental workflow

- Make annotation explicit and discoverable

- Structure the descriptions for consistency, tracking

# A community mobilization to develop standards, e.g.:



*de jure*

standard
organizations

*de facto*

grass-roots
groups

- Structural and operational differences
  - organization types (open, close to members, society, WG etc.)
  - standards development (how to formulate, conduct and maintain)
  - adoption, uptake, outreach (link to journals, funders and commercial sector)
  - funds (sponsors, memberships, grants, volunteering)

# Types of reporting standards



Including **conceptual model**, **conceptual schema** from which an exchange format is derived to allow data to flow from one system to another

Including **controlled vocabularies**, **taxonomies**, **thesauri**, **ontologies** etc. to use the same word and refer to the same 'thing'

Including **minimum information reporting requirements**, or **checklists** to report the same core, essential information

# Fragmentation, duplications and gaps



plant biology    epidemiology    microbiology

**Biologically**-delineated views of the world

**Generic** features ('common core')
- *description of source biomaterial*
- *experimental design components*

MS    MS

Arrays    Gels    NMR

Columns    FTIR

Scanning    Arrays & Scanning    Columns

**Technologically**-delineated views of the world

transcriptomics    proteomics    metabolomics

To compare and integrate data we need interoperable standards

# Growing number of reporting standards

*To track provenance of the information and ensure richness of data and experimental metadata descriptions, to maximize reusability*

**+ 303**

**+ 130**

**+ 150**

Estimated

Source: BioPortal

Source: MIBBI, EQUATOR

formats

terminologies

guidelines

?

Databases, annotation, curation tools

MAGE-Tab
GCDML
SRAxml
SOFT
CML          FASTA
        DICOM
GELML          SBRML
MITAB          MzML
ISA-Tab     SEDML…

AAO
CHEBI
OBI          VO
PATO          ENVO
MOD
TEDDY
XAO          BTO
DO     PRO      IDO…

miame
            MIAPA
MIRIAM
         MIX          MIQAS
                MIGEN     REMARK
         MIAPE          MIQE
            CIMR     CONSORT
         MIASE     MISFISHIE….

# But how much do we know about these standards?

- A *coherent, curated* and *searchable registry* of **standards** for describing and reporting experiments in life science, environmental, biomedical and biotechnological domains

- A *coherent, curated* and *searchable registry* of **standards** for describing and reporting experiments in life science, environmental, biomedical and biotechnological domains
- Progressively ***associate*** standards **to data policies** and **databases**
- ***Develop*** assessment **criteria** for usability and popularity of standards
- Help stakeholders to make informed decisions on e.g. what standards or databases to use or recommend; identify efforts they have funded

## ISA-Tab - Investigation Study Assay Tabular

General-purpose ISA-Tab file format - an extensible, hierarchical structure that focuses on the description of the experimental metadata (i.e. sample characteristics, technology and measurement types, sample-to-data relationships

### Database scope and data types

FUNCTIONAL GENOMICS | BIOLOGICAL MATERIAL | REPORT | EXPERIMENT | DEVICE | FILE | ASSAY | REAGENT | MATRIX

### Support

DOCUMENTATION  http://isatab.sourceforge.net/docs/ISA-TAB_release-candidate-1_v1.0_24nov08.pdf

Mailing List  isaforum[at]googlegroups.com

#### Contact Details

Philippe Rocca-Serra    Email

### Implementing Databases

**Stem Cell Discovery Engine**
Comparison system for cancer stem cell analysis

**Giga Science Database**
GigaDB primarily serves as a repository to host data and tools associated with articles in GigaScience; however, it also includes a subset of datasets that are not associated with GigaScience articles. GigaDB defines a dataset as a group of files (e.g., sequencing data, analyses, imaging files, software programs) that are related to and support an article or study.

**MetaboLights**
MetaboLights is a database for Metabolomics experiments and derived information. The database is

### Organisations

#### Maintainers

ISA community

#### Funders

BBSRC

### Publications

ISA software suite: supporting standards-compliant experimental annotation and enabling curation at the community level.

Bioinformatics 2010

[View Paper]

MetaboLights—an open-access general-purpose repository for metabolomics studies and associated meta-data

Nucleic Acid Research 2013

[View Paper]  [View Paper]

# bio**sharing**.org
## User profiles populated from ORCID...

### ORCID Profile

**Alejandra Gonzalez-Beltran**

Alejandra currently works in the ISA Team (http://www.isa-tools.org) at the Oxford e-Research Centre, University of Oxford, UK. Before that, Alejandra was at University College London, UK, working at the Computational and Systems Medicine project and the Department of Computer Science. Previously, she was awarded a PhD in Computer Science at Queen's University Belfast, UK and a Licentiateship (equivalent to MSc) from Universidad Nacional de Rosario, Argentina.

🌐 Websites

LinkedIn Profile

UCL Personal Website

OeRC Personal Website

View Alejandras profile on ORCID.

### Latest Publications

MetaboLights - An open-access general-purpose repository for metabolomics studies and associated meta-data

Read the paper   Get article metrics

Guidelines for information about therapy experiments: A proposal on best practice for recording experimental data on cancer therapy

Read the paper   Get article metrics

Establishing a knowledge trail from molecular experiments to clinical trials

Read the paper   Get article metrics

View the rest here

### My Standards

GIATE
**Guidelines for Information About Therapy Experiments**

👁 View Record

# biosharing.org
## ... credit for creating, contributing to, maintaining standards



**AGBELTRAN**

### ORCID Profile

Alejandra Gonzalez-Beltran

Alejandra currently works in the ISA Team (http://www.isa-tools.org) at the Oxford e-Research Centre, University of Oxford, UK. Before that, Alejandra was at University College London, UK, working at the Computational and Systems Medicine project and the Department of Computer Science. Previously, she was awarded a PhD in Computer Science at Queen's University Belfast, UK and a Licentiateship (equivalent to MSc) from Universidad Nacional de Rosario, Argentina.

### Websites

LinkedIn Profile

UCL Personal Website

OeRC Personal Website

View Alejandras profile on ORCID.

### Latest Publications

MetaboLights - An open-access general-purpose repository for metabolomics studies and associated meta-data

Read the paper    Get article metrics

Guidelines for information about therapy experiments: A proposal on best practice for recording experimental data on cancer therapy

Read the paper    Get article metrics

Establishing a knowledge trail from molecular experiments to clinical trials

Read the paper    Get article metrics

View the rest here

### My Standards

GIATE
Guidelines for Information About Therapy Experiments

👁 View Record

*Ownership* of open standards can be problematic in broad, grass-root collaborations

It requires improved models, to encourage *maintenance* of and *contributions* to these efforts, *rewards* and *incentives* need to be identified for all contributors to supporting the continued development of standards

# ... link to data records associated to publications

# ...and associated article-level metrics

# We need "standards impact metrics" to evaluate use/usability

▶ The standard itself
  – specification documentation
  – ease of implementation (eg, level of documentation, requirement for programmer support)
  – human and machine readability
  – formal structure
  – expressivity—the breadth of information that can be represented
  – ease of use, for example, minimal required fields, text-based interface familiarity to biologists.

▶ Adoption and user community
  – broad adoption and implementation, outside the initial group
  – support supplied by the user community
  – use by community databases
  – software development that supports the standard (eg, for curating, submitting to databases)
  – responsiveness to community requests
  – availability of examples of use
  – requirements of relevant authoritative bodies, for example, funders (NIH, National Science Foundation, Centers for Medicare & Medicaid Services), publishers, etc.

▶ Additional factors
  – integration/compatibility with other standards
  – extensibility and flexibility to cover new domains
  – conversion and mapping, when applicable
  – cost (eg, open vs licensing fee).

# biosharing.org working with data publication platforms:

# "Invisible" use of standards in data reporting tools

## isatools

isatab format

github SOCIAL CODING
GITHUB.COM/ISA-TOOLS

isacommons

Alejandra Gonzalez-Beltran, PhD

## ORCID CodeFest
Mash ups with the PUBLIC ORCID API
Thursday, May 23 2013    12:00 AM BST

One of the winners.
Project: integration of ORCID with the ISAcreator, the editor tool, helping curators and researchers to describe experiments following community standards.

### collect and curate, following standards
Describe the experimental steps using community-defined minimum reporting requirements and ontologies, where possible.

### store and browse, locally or publicly
create your own repository to search and browse the experimental description and associated data making it close or open.

### submit to public repositories
when required, reformat the experiments for submission to supported public repositories or directly export to those already using ISA-Tab.

### analyse with existing tools
upload experimental description and associated data to a growing number of well-known analysis systems, ISA connects with.

### release, reason and nanopublish
explore how to reason over your experiments, open them to the linked data universe, or publish nano-statements of your discoveries.

### publish data along your article
directly export your experiments to the new generation of data journals, accepting submission in ISA-Tab.

# ODIN mission

ODIN will build on the ORCID and DataCite initiatives to uniquely identify scientists and data sets and connect this information across multiple services and infrastructures for scholarly communication. It will address some of the critical open questions in the area:

- Referencing a data object

- Tracking of use and re-use

- Links between a data object, subsets, articles, rights statements and every person involved in its life-cycle.

# Summarizing my talk

*Problem*:
Identification of datasets in pivotal. But meaningful sharing and (re)use also depend on how well described the datasets are.

*Status quo:*
In the life sciences there is a wealth of 'reporting standards' set to enhance and facilitate the experimental descriptions.

*Challenges addressed by* biosharing
Identify 'reporting standards' and their organizations, track their use, usability and impact (e.g. linking them to datasets), credit their developers, users (e.g. curators)...

# Acknowledgements