# ISA infrastructure:
## collecting & managing functional genomics datasets with rich semantics

*Rencontres GenOuest, 2011, Rennes, 18th October 2011*

Philippe Rocca-Serra (Ph. D)
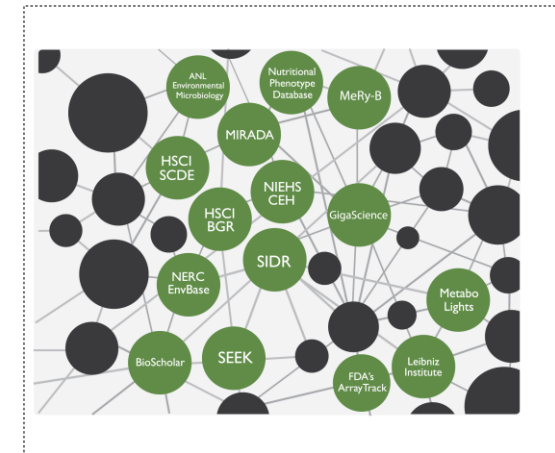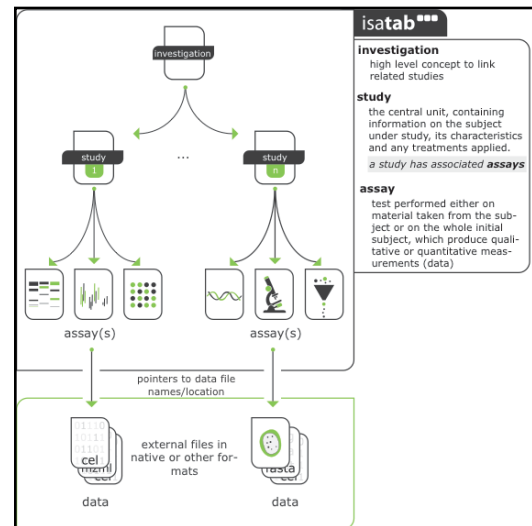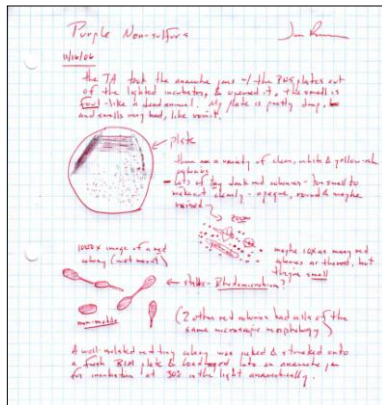
ISA Team

philippe.rocca-serra@oerc.ox.ac.uk

http://www.isa-tools.org

# Presentation Outline

- **<span style="color:green">Background information</span>**

- Rationale for developing ISA-tools

- ISA paradigm and interaction with ontologies

- Moving on: ISA future developments

# MAIN THEME:
## It is all about <u>structuring</u> experimental information to make it <u>available to computer</u> and software agents to enable <u>mining</u>.

## But let's proceed gradually…



Notes in Lab Books
(information for humans)

Spreadsheets and Tables
( the compromise)

Facts as RDF statements
(information for machines)

# What's wrong with free text in electronic records?

An example: {WT, wildtype, wild-type, sauvage, wildtypo}

- Hard to index

- Hard to search

- Poor query recall

+>Unhappy users & missed opportunities

Limit free text by means of controlled terminologies and ontologies

# Ontologies ?....it is about formalizing knowledge

- Organizing <u>types</u> into broad categories (e.g. `Objects,Subjects, Properties, Events` or `Material, Qualities, Processes`)

- Defining the properties of those types using sets of relations (e.g. `is_a, has_part/part_of, derives_from, located_in, participates_to`)

- Advanced software supporting the validation of those representations: Reasoners (Fact, Pellet, Hermit or ELK for OWL ontologies), .

# How can this be useful?

Just one simple example:

- It makes possibly things like <u>query expansion</u>:

  Searching for word 'cancer' should retrieve:

  {'carcinoma, adenocarcinoma, lipoma, sarcoma…}

  - How does it work?

  +> taking advantage of the 'is_a' relationship between those entities

# An example of query expansion

| | ID | Title | Assays | Species |
|---|---|---|---|---|
| ⊞ | E-TABM-1054 | miRNA expression profile between ER-beta- and ER-beta+ breast tumors. | 36 | Homo sapiens |
| ⊞ | E-MEXP-3025 | Spheres culture from lung adenocarcinoma pleural effusions | 19 | Homo sapiens |
| ⊞ | E-TABM-1055 | MicroRNA profiling by array of human MCF-7 breast cancer cells with ER-beta tagged at the C-terminal or N-ter... | 15 | Homo sapiens |
| ⊞ | E-TABM-1053 | Transcription profiling by array of human MCF-7 breast cell clones expressing ER-beta tagged with TAP-tag at th... | 12 | Homo sapiens |
| ⊞ | E-TABM-1052 | Transcription profiling by array of human MCF-7 cells with ER-beta tagged with TAP-tag at the C-term or N-term... | 12 | Homo sapiens |
| ⊞ | E-MEXP-3192 | Lapatinib and retinoic acid combination treatment of SKBR3 breast cancer cells | 32 | Homo sapiens |
| ⊞ | E-GEOD-27514 | Identification of a Potently Oncogenic CALM-AF10 Minimal-Fusion Mutant | 24 | Mus musculus |
| ⊞ | E-GEOD-27513 | Identification of a Potently Oncogenic CALM-AF10 Minimal-Fusion Mutant (mRNA) | 12 | Mus musculus |
| ⊞ | E-GEOD-27512 | Identification of a Potently Oncogenic CALM-AF10 Minimal-Fusion Mutant (miRNA) | 12 | Mus musculus |
| ⊞ | E-GEOD-25519 | Promoter methylation data: OHT/ICI-sensitive vs. -resistant cell lines | 2 | Homo sapiens |
| ⊞ | E-GEOD-15308 | Genomics of oral cancer cells | 11 | Homo sapiens |
| ⊞ | E-GEOD-24751 | Pulmonary gene and microRNA expression changes in mice exposed to benzo(a)pyrene by oral gavage | 45 | Mus musculus |
| ⊞ | E-GEOD-24520 | Epigenetic Based Enrichment of Cancer Stem Cells: Mechanistic and Clinical Implications for Liver Cancer | 38 | Homo sapiens |
| ⊞ | E-GEOD-32492 | Identification of Candidate Tumor Suppressor Genes Inactivated by Promoter Methylation in Melanoma | 24 | Homo sapiens |
| ⊞ | E-GEOD-23603 | Gene expression in ovarian cancer | 84 | Homo sapiens |

[Source:
http://www.ebi.ac.uk/arrayexpress/browse.html?keywords=cancer&expandefo=on.]

The system may suggest unseen association, could
help generate new hypothesis -> Happier users
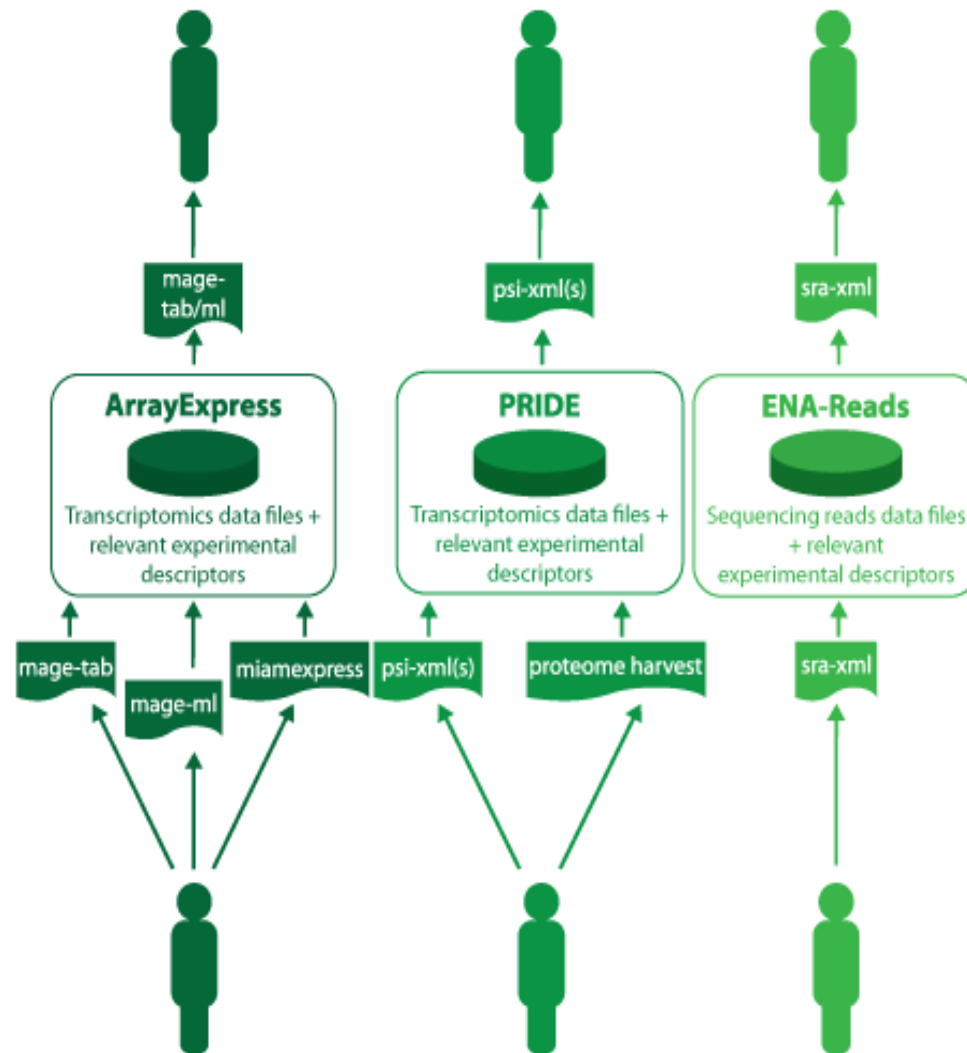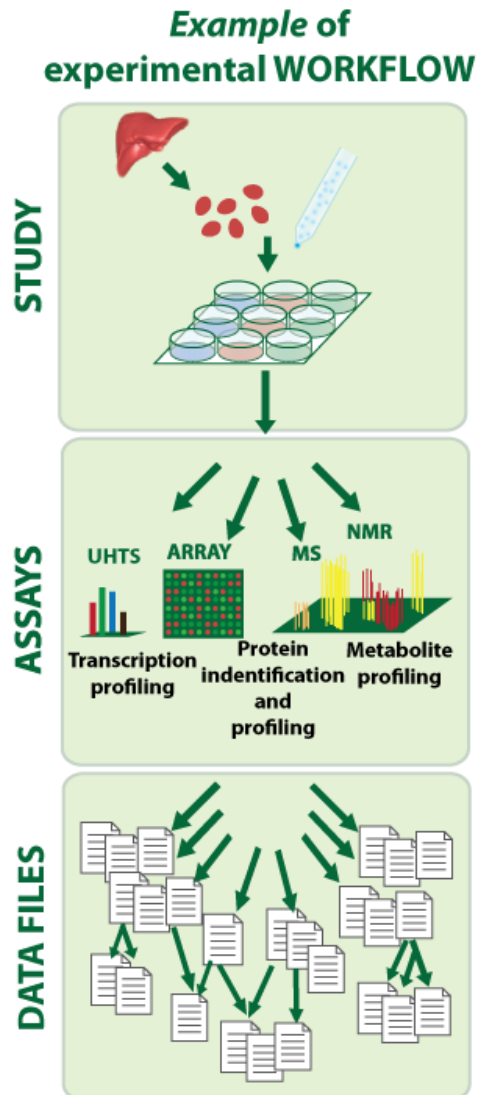
# Presentation Outline

- Background information

- Rationale for developing ISA-tools

- ISA paradigm and interaction with ontologies

- Moving on: ISA future developments

# Observations

- Experiments are expensive, often publicly funded, still many fail to see the light.

- Spreadsheets are the most common vehicle for so-called 'omics' (functional genomics) experimental metadata tracking

- technology centric repositories form de facto silos

- conversions are required to allow for deposition to public databases.

- submitting to common information across a series of repositories is inefficient

# Case Study

# Observations II

- A growing number to 'annotations requirements' (a.k.a MI checklists)

- Many different communities, many different needs.



- Creation of the MIBBI portal to harmonize and identify a core of common descriptors, create extensions where necessary.

# Many Requirements, Many Formats, Many ontologies …

- To support different fields of molecular biology:
  - Soil Metagenomics
  - Cancer genomics
  - Chromatin remodelling event and Stem Cell fate.
- To support various cases in data reporting & data management
  - Dealing with legacy data (spreadsheets hanging around)
    - Data Mapping and Import Function from files.
  - Dealing with de novo datasets:
    - Planning and Templating: reduce repetitive tasks by relying on patterns found in experimental designs

# Many ontologies, Many Formats, Many Requirements…

Credits: http://liverpoolsolfed.wordpress.com/resources/image-bank/demonstration/

# ISA framework overview

# A focus on standards...



OBO and OWL ontologies

# Presentation Outline

- Background information

- Rationale for developing ISA-tools

- ISA paradigm and interaction with ontologies

- Moving on: ISA future developments

# ISA syntax and Table definition

- Configuration files broadly define Material /Data workflows

**metaboliteprofiling_ms**

**fields**

**Sample Name**
**Protocol REF**
**Extract Name**
**Protocol REF**
**Labeled Extract Name**
**Label**
**Protocol REF**
**Parameter Value[instrument]**
**Parameter Value[ion source]**
**Parameter Value[detector]**
**Parameter Value[analyzer]**
**MS Assay Name**
**Raw Spectral Data File**
**Protocol REF**
**Normalization Name**
**Data Transformation Name**
**Derived Spectral Data File**
**Protocol REF**
**Metabolite Assignment File**
**Factors**

Input Material or Data Node

Output Material or Data Node

Characteristics[…]
Factor Value[…]

Characteristics[…]
Factor Value[…]

Protocol REF

Parameter Value […]

# List of supported assays in ISA default configuration

| ● measurement type | ● technology type |
|---|---|
| cell counting | flow cytometry |
| cell sorting | flow cytometry |
| clinical chemistry analysis | |
| copy number variation profiling | DNA microarray |
| DNA methylation profiling | DNA microarray |
| DNA methylation profiling | nucleotide sequen... |
| environmental gene survey | nucleotide sequen... |
| genome sequencing | nucleotide sequen... |
| hematology | |
| loss of heterozygosity profiling | DNA microarray |
| histology | |
| histone modification profiling | nucleotide sequen... |
| metabolite profiling | mass spectrometry |
| metabolite profiling | NMR spectroscopy |
| metagenome sequencing | nucleotide sequen... |
| protein-protein interaction detection | protein microarray |
| protein-DNA binding site identification | DNA microarray |
| metabolite profiling | NMR spectroscopy |
| metagenome sequencing | nucleotide sequen... |
| protein-protein interaction detection | protein microarray |
| protein-DNA binding site identification | DNA microarray |
| protein-DNA binding site identification | nucleotide sequen... |
| protein expression profiling | gel electrophoresis |
| protein expression profiling | protein microarray |
| protein expression profiling | mass spectrometry |
| protein identification | mass spectrometry |
| SNP analysis | DNA microarray |
| [Sample] | |
| transcription factor binding site identific... | DNA microarray |
| transcription factor binding site identific... | nucleotide sequen... |
| transcription profiling | DNA microarray |
| transcription profiling | real time PCR |
| transcription profiling | nucleotide sequen... |

Potential for Compliance with:
MIGS
MIMARKS
MIAME

Expanding the number of ISAconfigurations

Available from:
https://github.com/ISA-tools/Configuration-Files

# ISAconfigurator Tables

# ISAconfigurator Tables

# ISAconfigurator Tables



This is an example of a field definition created by the configurator. In this instance we are describing a label field, in particular, one used to describe the label used in a microarray experiment.

We have defined it to come from an ontology, and we recommend the ChEBI ontology. It is also required.

```
<field header="Label" data-type="Ontology term" is-file-field="false" is-multiple-value="false"
       is-required="true">
    <description>Indicates a chemical or biological marker, such as a radioactive isotope or a fluorescent dye
        which is bound to a material in order to make it detectable by some assay technology (e.g. P33, biotin,
        GFP).
    </description>
    <default-value/>
    <recommended-ontologies>
        <ontology id="1007" abbreviation="CHEBI" name="Chemical entities of biological interest"
                  version="46223"/>
    </recommended-ontologies>
</field>
```

# How do ISA tools access Ontology servers?

**Configuration**

Configuring fields to be defined by ontologies

ISAconfigurator is a tool for customizing annotation requirements within ISA-TAB syntax. ISAconfigurator relies on NCBO services when setting an ISA-TAB syntactic element to be of type Ontology Term. Superuser can select one or more ontology resources and within any given resource, select a node and children to restrict or define the annotation space.

**Ontology browsing & searching**

ISAcreator provides a unique spreadsheet embedded search and browse ontology functionality.

**Ontology tagging**

To complement this approach, ISAcreator makes use of NCBO Annotator service to allow end users to tag free text with ontology terms (in line with restrictions set in ISAconfigurations).

**Ontology Resource Manager**

**The resource manager** provides seamless searching of ontology resources, regardless of their origins, their underlying data schema or the mechanism (REST, SOAP or local file store) through which they are accessed.

**NCBO BioPortal**

Search, Hierarchy and Annotator services

**REST**

**Ontology Lookup Service (OLS)**

**SOAP**

**Plugin**

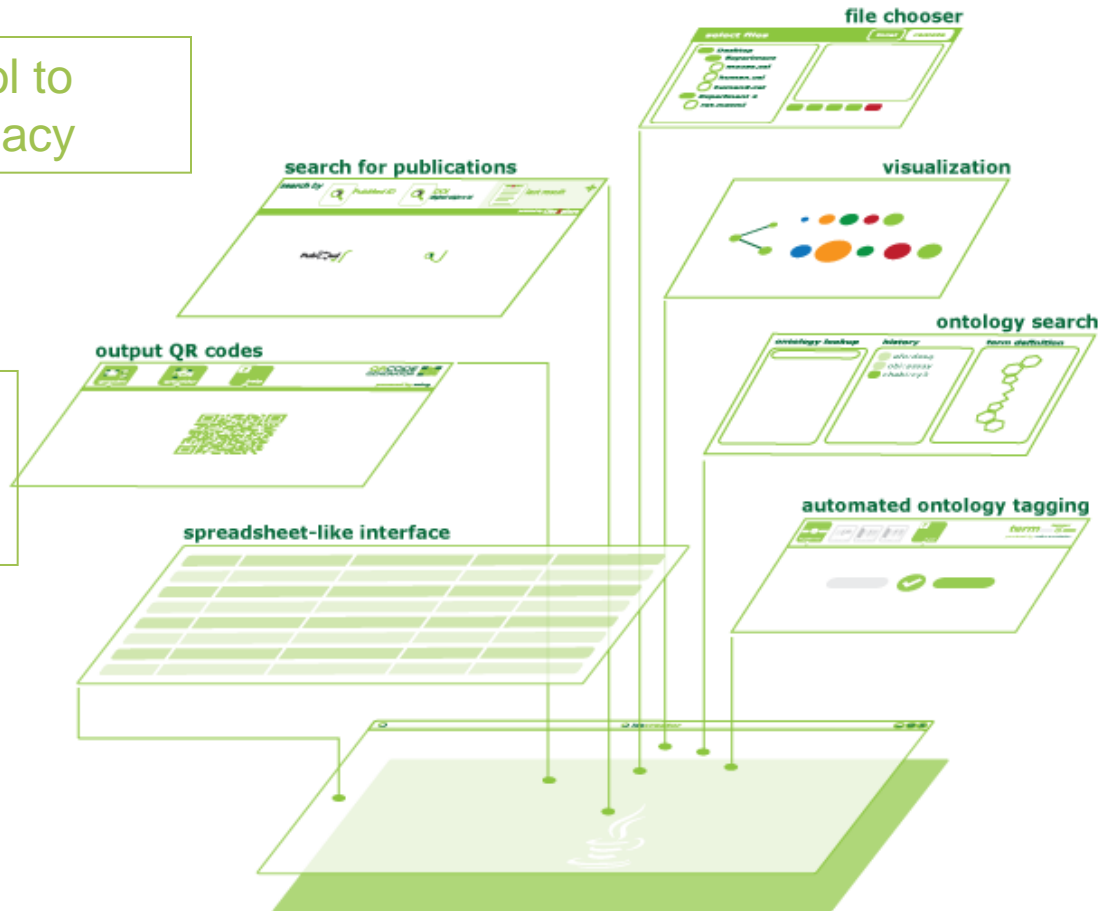**REST** **SOAP** **LOCAL**

# Anatomy of ISAcreator

create your first study | the anatomy of **isatab** | the anatomy of ○ isa**creator**

Mapping Tool to deal with legacy data

Experiment design wizard for templating

file chooser

search for publications

output QR codes

spreadsheet-like interface

visualization

ontology search

automated ontology tagging

powered by Java™
resulting a cross-platform tool. It is not limited to your operating environment so will work on

# Select and Annotate in ISAcreator

# Calling NCBO Annotator service from ISAcreator: Tagging free text

# Calling NCBO Annotator service from ISAcreator: Tagging free text

# Calling NCBO Annotator service from ISAcreator: Tagging free text

# Presentation Outline

- Background information

- Rationale for developing ISA-tools

- ISA paradigm and interaction with ontologies

- Moving on: ISA future developments

# Expand ISA community and welcome new members

- Metabolights, EBI's metabolomic data repository
- BBSRC funding for a Metagenomic Portal (EBI collaboration)
- Expand ISAconverter to support additional XML formats
  - GEO MiniML (ongoing development)
  - FuGE-ML (carried out at INIST [Magali Roux et al])
  - SRA XML maintenance  (regular schema updates)
- Scout for new domains of application

# Exposing Experimental Metadata on the semantic web

- Expansion of **ISAconverter** to provide RDF/OWL representation of experimental data

- Ontologies or Vocabularies matter again.

- Mapping of ISA elements to resources such as: OBI ontology classes, FOAF, Dublin Core.

- Expand ISAconfigurator to enable recording of Mapping to Vocabulary

# ISA2RDF work in progress

- Use case on W3C HCLS scientific discourse list
  - deciding on the granularity of representation
  - building on previous experience
  - Evaluating alternative representations.
- Participipation to the Biohackathon 2011
  - http://blogs.openaccesscentral.com/blogs/bmcblog/entry/biohackathon_2011_number_1
  - Discussing best practices
    - Use of URI supplied by **www.identifiers.org**
    - Avoid use of blank node as much as possible

# ISA2RDF: work in progress

```
<foaf:Organization rdf:about="http://dbpedia.org/NERC_CEH_Oxford">
    <foaf:member>
      <foaf:Person rdf:about="http://www.orcid.org/f85c348a-df04-4dc2-b35d-0cb89bbce664">
        <foaf:family_name>Tiwari</foaf:family_name>
        <foaf:mbox></foaf:mbox>
        <rdfs:label>Bela Tiwari</rdfs:label>
        <foaf:firstName>Bela</foaf:firstName>
      </foaf:Person>
    </foaf:member>
  </foaf:Organization>
  <dcmitype:Text rdf:about="http://www.protocol.org/9edbdca9-8db7-474b-808a-cfb51a37c049">
    <geo:lat>+45.3</geo:lat>
    <dcterms:references>pmid:14973331</dcterms:references>
    <dcterms:hasVersion></dcterms:hasVersion>
    <dcterms:description>DSN normalisation is described in Zhulidov et al, 2004 and enriches for full length transcript sequences and
equal transcript abundance.</dcterms:description>
    <dcterms:coverage>normalisation</dcterms:coverage>
    <rdfs:label>Duplex-Specific-Nuclease-normalisation</rdfs:label>
    <dc:language>en</dc:language>
    <dc:type>Protocol</dc:type>
  </dcmitype:Text>
```
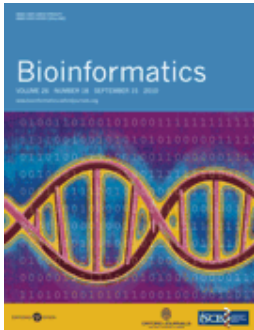
rdf:subject - rdf:predicate - rdf:object:

<lipoprotein>-<affects><inflammatory_cell>

<PRO:212342352>-<RO:543636><CL:84872762>

# Publication...

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

Philippe Rocca-Serra; Marco Brandizi; Eamonn Maguire; Nataliya Sklyar; Chris Taylor; Kimberly Begley; Dawn Field; Stephen Harris; Winston Hide; Oliver Hofmann; Steffen Neumann; Peter Sterk; Weida Tong; Susanna-Assunta Sansone

# MERCI!

Groups and individuals participating in:

MIBBI http://mibbi.org

ISA-Tab format http://isatab.sf.net

OBO Foundry http://obofoundry.org

OBI:  http://obi-ontology.org/page/Main_Page

ISA Infrastructure Team:

Eamonn Maguire (Oxford)

Philippe Rocca-Serra (Oxford)

Susanna-Assunta Sansone (Oxford)

Chris Taylor (EMBL-EBI)

Alumni:

Marco Brandizi (EMBL-EBI)

Nataliya Sklyar (EMBL-EBI)

collaborators at:

Cambridge University

EuNuGO

Harvard School for Public Health

FDAs NCTR

Leibniz Plant Institute

NERCs NEBC

SIDR, INIST

Metabolights, EMBL-EBI

Funders:

EU Carcinogenomics Project

UK BBSRC

# MERCI!

Groups and individuals participating in:
Dawn Field: NERC CEH Oxford
Winston Hide: HSPH
Oliver Hoffman: HSPH
Shannan Ho Sui : HSPH
Brad Chapman: HSPH
Christoph Steinbeck: Metabolights
Kenneth Haug: Metabolights
Paula de Matos: Metabolights
Magali Roux: INIST
Florian Mazur: INIST
Alain Zasadzinki: INIST
Marie Christine Jacquemot: INIST
And many more who have to forgive us!

# Questions: