# CEDAR
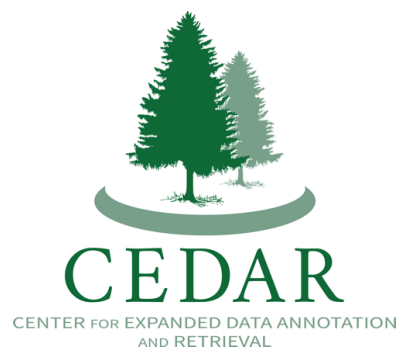
*Making it Easier to Use Ontologies to Author Experimental Metadata*

Mark A. Musen, M.D., Ph.D.

Stanford University



CEDAR
CENTER for EXPANDED DATA ANNOTATION AND RETRIEVAL

## Slide 1 (The Economist cover)

The Economist

Washington's lawyer surplus
How to do a nuclear deal with Iran
Investment tips from Nobel economists
Junk bonds are back
The meaning of Sachin Tendulkar

OCTOBER 19TH–25TH 2013    Economist.com

HOW SCIENCE GOES WRONG

99
Einsteinium

## Slide 2

Conducted

3600832 73931 96703291 194266 7555 938 6317409 81179881 1394773 307135 84439499 9931330 19729 807 3910094 0 27947 32252628914 3910 660718 9973066 57 58082 4375 36241878 572140 24375470 22926488 34853560 75499885 23391181 76156678 121030 36449393 2340635 65758283 67206007 96064598 52231472345 37613 25921356 77069117 75568973 831589502 68755674 38319872 565

3600832 73931 96
703291 194266 687674
38319 872 565

Captured

3600
832
73931

Collocated

**Biomedical Research Studies:
The Data Pipeline**

Consistent

6

## Minimum Information About a Microarray Experiment - MIAME

**MIAME** describes the **Minimum Information About a Microarray Experiment** that is needed to enable the interpretation of the results of the experiment unambiguously and potentially to reproduce the experiment. [Brazma et al., Nature Genetics]

The six most critical elements contributing towards MIAME are:

1. The raw data for each hybridisation (e.g., CEL or GPR files)
2. The final processed (normalised) data for the set of hybridisations in the experiment (study) (e.g., the gene expression data matrix used to draw the conclusions from the study)
3. The essential sample annotation including experimental factors and their values (e.g., compound and dose in a dose response experiment)
4. The experimental design including sample data relationships (e.g., which raw data file relates to which sample, which hybridisations are technical, which are biological replicates)
5. Sufficient annotation of the array (e.g., gene identifiers, genomic coordinates, probe oligonucleotide sequences or reference commercial array catalog number)
6. The essential laboratory and data processing protocols (e.g., what normalisation method has been used to obtain the final processed data)

For more details, see MIAME 2.0.

MIAME does not specify a particular format, however, obviously the data are more usable, if it is encoded in a way that the essential information specified by MIAME can be accessed easily. FGED recommends the use of MAGE-TAB format, which is based on spreadsheets, or MAGE-ML.

MIAME also does not specify any particular terminology, however for automated data exchange the use of standard controlled vocabularies and ontologies are desirable. FGED recommends the use of MGED Ontology for the description of the key experimental concepts, and where possible ontologies developed by the respective community for describing terms such as anatomy, disease, chemical compounds etc (see OBO page for more detail).

---

**The Good News:** Minimal information checklists, such as MIAME, are being advanced from all sectors of the biomedical community

**The Bad News:** Investigators view requests for even "minimal" information as burdensome



4

### MIBBI portal

article | discussion | view source | history

- Registration form ⬚ for the MIBBI Portal (please return to chrisftaylor[@]gmail.com)
- Summary spreadsheet ⬚ of all registered projects
- XML document ⬚ containing all registered projects (from this schema ⬚, same information as the Excel spreadsheet)

### Bioscience projects registered with MIBBI

| | |
|---|---|
| CIMR | Core Information for Metabolomics Reporting |
| GIATE | Guidelines for Information About Therapy Experiments |
| MIABE | Minimal Information About a Bioactive Entity |
| MIABiE | Minimum Information About a Biofilm Experiment |
| MIACA | Minimal Information About a Cellular Assay |
| MIAME | Minimum Information About a Microarray Experiment |
| MIAPA | Minimum Information About a Phylogenetic Analysis |
| MIAPAR | Minimum Information About a Protein Affinity Reagent |
| MIAPE | Minimum Information About a Proteomics Experiment |
| MIAPepAE | Minimum Information About a Peptide Array Experiment |
| MIARE | Minimum Information About a RNAi Experiment |
| MIASE | Minimum Information About a Simulation Experiment |
| MIASPPE | Minimum Information About Sample Preparation for a Phosphoproteomics Experiment |
| MIATA | Minimum Information About T Cell Assays |
| MICEE | Minimum Information about a Cardiac Electrophysiology Experiment |

---

UNIVERSITY OF OXFORD

npg nature publishing group · OXFORD UNIVERSITY PRESS · PLOS · BMJ · re3data.org REGISTRY OF RESEARCH DATA REPOSITORIES · elixir UNITED KINGDOM

F1000 FACULTY of 1000 · EMBOpress · BioMed Central The Open Access Publisher · DRYAD · FORCE11

### biosharing.org

**STANDARDS**

BioSharing standards have been partly compiled by linking to BioPortal, MIBBI and the Equator Network.
Or you can filter on MIBBI Foundry reporting guidelines or OBO Foundry terminology artifacts.

68 guidelines     168 formats

View as Grid | View as Table       40 records in view

**Standard Type**

| | |
|---|---|
| REPORTING GUIDELINE | 40 |
| EXCHANGE FORMAT | 0 |
| TERMINOLOGY ARTIFACT | 0 |

**Domains**

| | |
|---|---|
| ASSAY | 14 |
| DNA | 12 |
| RNA | 8 |
| PROTEIN | 7 |
| TRANSCRIPTOME | 5 |
| BIOCHEMISTRY | 4 |
| BRAIN | 4 |

**BioDBCore**
Core Attributes of Biological Databases
REPORTING GUIDELINE
Systems 1
Publications 1
1 Taxa types, including:
ALL
1 Data types, including:
DATABASE

**CIMR**
Core Information for Metabolomics Reporting
Systems 2
Publications 2
No taxa defined.
5 Data types, including:
METABOLITE

**GIATE**
Guidelines for Information About Therapy Experiments
Systems 0
Publications 5
No taxa defined.
2 Data types, including:
TREATMENT | ANTIBODY

**MIABE**
Minimum Information About a Bioactive Entity
Systems 1
Publications 1
No taxa defined.
2 Data types, including:
BIOACTIVITY | MOLECULAR ENTITY

## The Existing BioSharing Approach is Limited

- Emphasis traditionally has been on development of simple checklists of metadata elements
- Little practical consideration of
  - How to supply values for the metadata elements
  - Standard ontologies that might be used
- We need a more expressive—and *computable*—framework for describing metadata

## The ISA model

- Developed by BioSharing group and supported by a suite of tools
- Provides structure for metadata related to
  - Investigation
  - Study
  - Assay
- Is not easily extended within existing tool set
- Forms the foundation for the modeling of metadata in the CEDAR project

Example entity–relationship diagram
for describing metadata for annotating
multiplex bead array assays
(e.g., Luminex)

# A Metadata Ecosystem

- **HIPC investigators** perform experiments in human immunology
- **HIPC Standards Working Group** creates metadata templates to annotate experimental data in a uniform manner
- **ImmPort** stores HIPC data (and metadata) in its public repository
- CEDAR will ease
  - Template creation and management
  - The use of templates to author metadata for ImmPort
  - Analysis of existing metadata to inform the authoring of new metadata

# The CEDAR Approach to Metadata

# CEDAR technology will give us

- Mechanisms
  - To author metadata template elements
  - To assemble them into composite templates
  - To fill out templates to encode experimental metadata
- A repository of metadata from which we can
  - Learn metadata patterns
  - Guide predictive entry of new metadata
- Links to the National Center for Biomedical Ontology to ensure that metadata are encoded using appropriate ontology terms

# The National Center for Biomedical Ontology

- We **create and maintain a library** of biomedical ontologies and terminologies.
- We **build tools and Web services** to enable the use of ontologies and terminologies.
- We **collaborate with scientific communities** that develop and use ontologies and terminologies in biomedicine.

NATIONAL CENTER FOR
BIOMEDICAL ONTOLOGY

# The CEDAR Approach to Metadata

**Authoring of Metadata Templates**

Template authors (e.g., standards committees)

define

Metadata tempates

**Annotation of Data with Metadata**

Scientists

contribute | fill in

Metadata acquisition forms

**Exploration and Reuse of Datasets through Metadata**

search, reuse

Metadata repository

GEO — Gene Expression Omnibus
ImmPort
HMP HUMAN MICROBIOME PROJECT
The Cancer Genome Atlas — Understanding genomics to improve cancer care

---

CEDAR Repository Model
for ImmPort and ISA Studies

Contact
Publication
Investigation
Organization
Study
Study Factor
Study Protocol
Study Assay
Protocol Parameter
Study Group Population
Parameter Value
Process
Study Subject
Characteristic Value
output
input
Reagent
Study Time
Data File
Result Value
Sample
Factor Value

22

# The CEDAR Approach to Metadata

The CEDAR Approach to Metadata

Learning for Predictive Metadata Entry

Linked publications

Related data

Related templates

Experiment

Organism — Homo sapiens

Platform — HG-U133_Plus_2

Type — RNA

Cell — brain tumor tissue

?

Samples

Age — 20, 45, 57

Gender — M, F, F

?

Gene expression template

# How can we make
# metadata authoring better?

- Create an ecosystem based on searchable, "smart" metadata templates
- Predefine standard value sets to fill in the blanks
- Use machine learning to enable predictive metadata entry
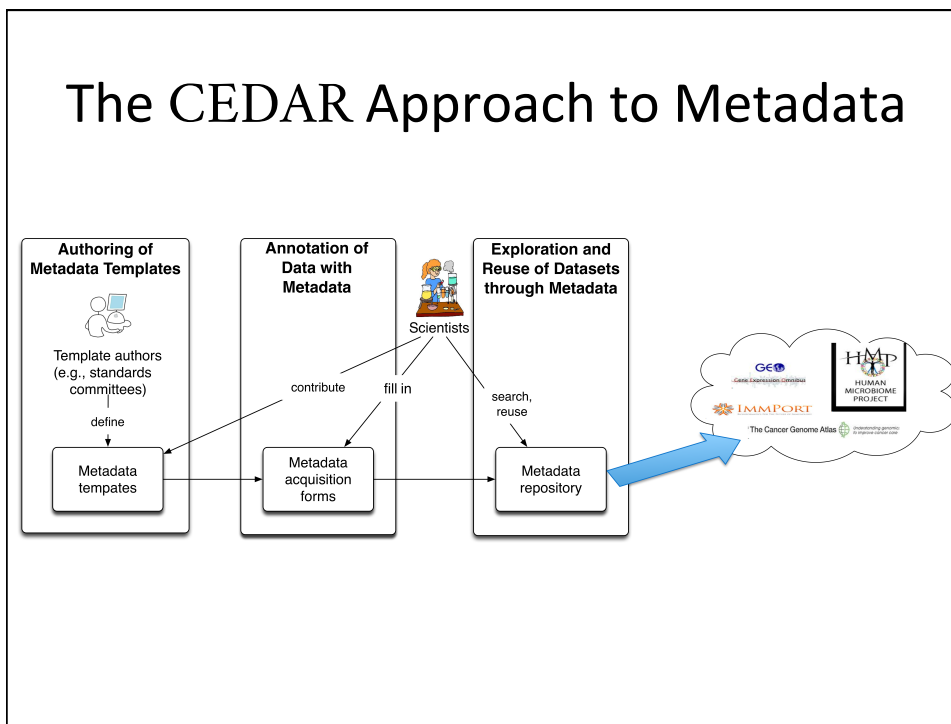- Use text processing to acquire metadata from written descriptions of the experiment (e.g., from PubMed and PubMed Central)

# How can we make
# metadata themselves better?

- Mirror metadata authored with CEDAR tools in our own metadata repository
- Augment those metadata with links to the published literature (including secondary analyses and retractions!)
- Augment those metadata with links to follow-up experiments (in online databases and in the literature)
- Allow the scientific community to comment on the experiment through structured metadata
- Learn from the metadata repository to ease the authoring of new metadata

## The CEDAR Approach to Metadata





http://metadatacenter.org