# Advancing Discovery Science
# Predictive, Evidential and Meta Analytical Methods

**Michel Dumontier**

Associate Professor of Medicine

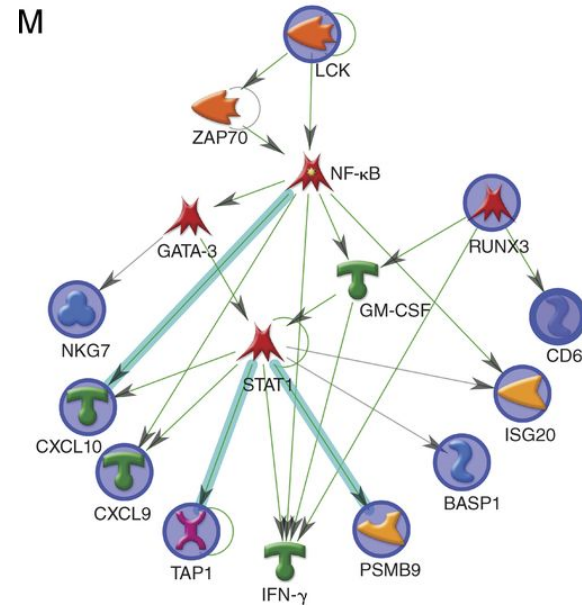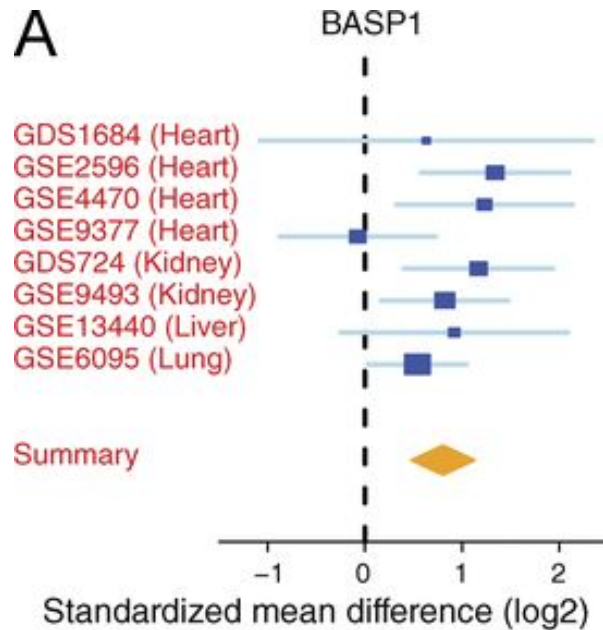Stanford Center for Biomedical Informatics Research

Stanford University

# New discoveries are being made by "research parasites" using other people's data



**A common rejection module (CRM) for acute rejection across multiple organs identifies novel therapeutics for organ transplantation**

**towards automated knowledge discovery**

**Development of an intelligent system for scientific inquiry using the <u>totality of web-accessible data and services</u>.**

# Challenge

Efficient and uniform **access** to *distributed, versioned* and *self describing* ***data*** and ***services***

for *reproducible* **analyses**

# A fundamental inability
# to easily query and mine structured knowledge

# Linked Open Data uses the web as a framework to share and link semantically annotated data



Linked Datasets as of August 2014

# Descriptions of the data (metadata) are often messy, incomplete or missing



NCBI > GEO > **Accession Display** ?

**GEO help:** Mouse over screen elements for information.

Scope: Self | Format: HTML | Amount: Quick |

**Series GSE35240**

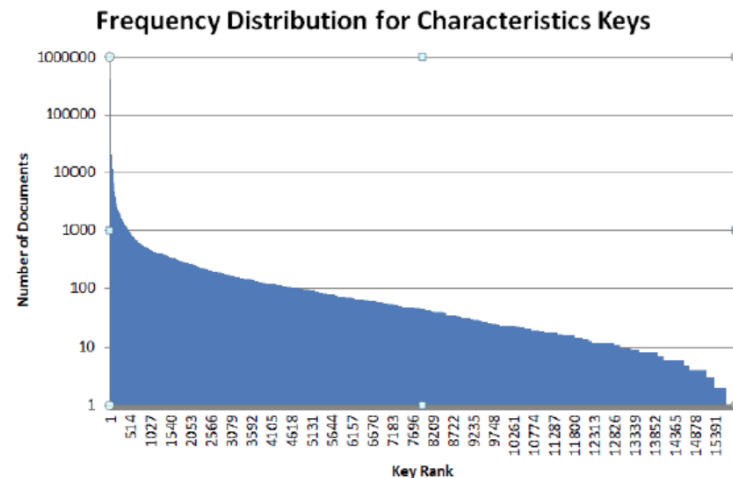| | |
|---|---|
| Status | Public on Aug 20, 2012 |
| Title | Gene expression in mitotic tissues of Drosop... too many centrosomes |
| Organism | Drosophila melanogaster |
| Experiment type | Expression profiling by array |
| Summary | Centrosome defects are a common feature can proceed through the majority of develop... amplified centrosomes in most of their cells. centrosome defects do not cause many prob... they can adapt to cope with any problems t and centrosome amplification fly b... assess how centrosome loss or centrosome a by profiling the global transcriptome of Droso... that either lack centrosomes or have too ma... |
| Overall design | Mitotic tissues (brains and imaginal discs Drosophila larvae of mutants lacking centros... with too many centrosomes (SakOE) and tv and OregonR). We extracted RNA from thre... used it for hybridisation to Affymetrix Dr... biological sample, material dissected fron... expression of the mutant strains was compar... |
| Contributor(s) | Baumbach J, Levesque MP, Raff JW |
| Citation(s) | Baumbach J, Levesque MP, Raff JW. Centrosc... dramatically perturb global gene expression in 15;1(10):983-93. PMID: 23213376 |

| Key | Count |
|---|---|
| age | 207147 |
| Age | 18089 |
| age (yrs) | 9891 |
| age (years) | 9272 |
| age (y) | 6226 |
| age in years | 1387 |
| age_years | 607 |
| AGE | 588 |
| age(years) | 558 |
| age (year) | 433 |
| age (yr) | 373 |
| Age (years) | 318 |
| age (in years) | 310 |
| Age(years) | 267 |
| age [year] | 97 |
| age [y] | 84 |
| age [years] | 83 |
| Age(yrs.) | 81 |
| age.year | 70 |
| age (yr-old) | 64 |
| age(yrs) | 59 |
| age of patient | 40 |
| Age, year | 39 |
| Age (yrs) | 36 |
| Age of patient | 33 |
| age, years | 24 |
| 'Age | 21 |
| Age (Years) | 20 |
| age (after birth) | 18 |
| age, yrs | 12 |
| age of subjects | 4 |



**Frequency Distribution for Characteristics Keys**

# Vast numbers of schemas and terminologies make for a *confusing* set of choices



**BioPortal**

**Statistics**

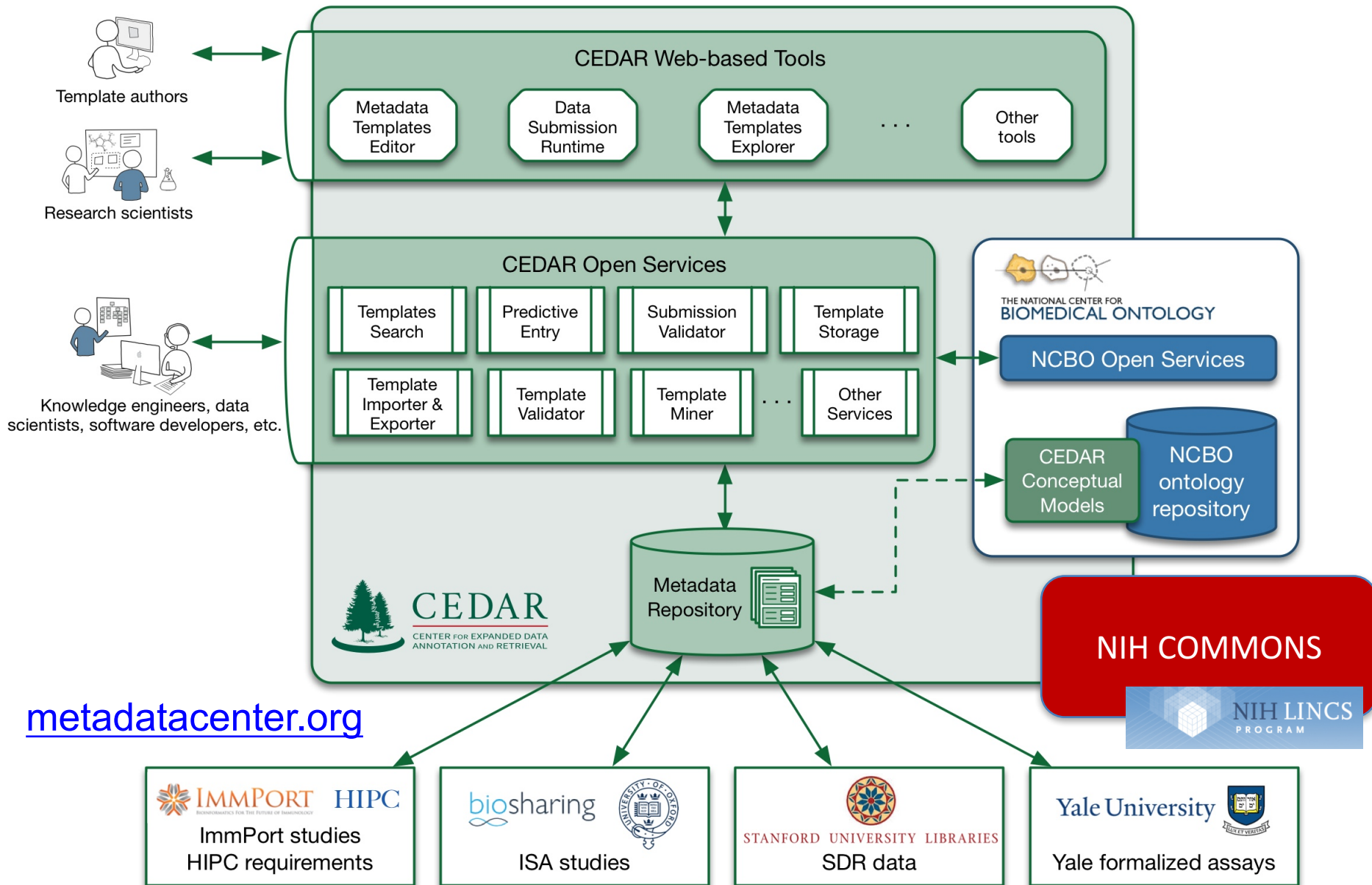| | |
|---|---|
| Ontologies | 513 |
| Classes | 7,996,220 |
| Resources Indexed | 48 |
| Indexed Records | 39,359,542 |
| Direct Annotations | 95,468,433,792 |
| Direct Plus Expanded Annotations | 144,789,582,932 |

Linked Open Vocabularies (LOV)

**542 Vocabularies in LOV**

# *Making it Easier, Possibly Even Pleasant, to Author Interoperable Experimental Metadata*



metadatacenter.org

# smartAPI: semantic and self describing REST APIs



SWAGGER

# Field auto-suggestion w/conformance

# Value auto-suggestion
## *calls the smartAPI field-specific suggestion service*

# Unify API data
# with Linked Open Data



a) Simplified MyGene.info object with JSON-LD context

```
1.    {
2.      "@type": "http:/identifiers.org/ncbigene/",
3.      "@context": {
4.        "_id": "@id",
5.        "name": "http://schema.org/name",
6.        "interpro": {
7.          "@id": "http:/identifiers.org/interpro/",
8.          "@type": "@id"
9.        },
10.       "description": "http://schema.org/description"
11.     },
12.     "_id": "1017",
13.     "symbol": "CDK2",
14.     "name": "cyclin-dependent kinase 2",
15.     "interpro": {
16.       "_id": "IPR000719",
17.       "description": "Protein kinase-like domain"
18.     }
19.   }
20.
```

b) Transformed JSON-LD object with semantic URIs included

```
1.    {
2.      "@id": "1017",
3.      "@type": "http:/identifiers.org/ncbigene/",
4.      "http://schema.org/name": "cyclin-dependent kinase 2",
5.      "http:/identifiers.org/interpro/": {
6.        "@id": "IPR000719",
7.        "http://schema.org/description": "Protein kinase-like domain"
8.      }
9.    }
```

# SCIENTIFIC DATA

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

Mark D. Wilkinson, Michel Dumontier  […]  Barend Mons

Affiliations | Contributions | Corresponding author

**About** *Scientific Data*

*Scientific Data* is an open-access, peer-reviewed journal for descriptions of scientifically valuable datasets. Our primary article-type, the **Data Descriptor**, is designed to make your data more discoverable, interpretable and reusable.

# FAIR: Findable, Accessible, Interoperable, Re-usable

# Applies to *all* digital *resources* and their *metadata*

# Challenge

Data will always be described using different schemas and vocabularies. Does that *still* matter? Can we *automate* the **integration of data**?

# Schema and vocabulary heterogeneity is a challenge for data retrieval and data mining



Three ways to model the relationship between a protein and the volume it occupies.

# New Mehtods for Data Integration

- Many elegant solutions for entity or concept mappings, but these only offer an incomplete solution when combined with schemas

- Need to *learn* robust **transformation patterns**
  - Subsumption, Similarity, Analogy, ML, Probability

- Evaluate these in the context of use cases
  - Query answering
  - Data mining
  - Prediction

# Challenge

## Can we *scale*

## the **validation** of **interesting findings**?

**Most published research findings are <span style="color:red">false</span>**
- John Ioannidis, Stanford University

Ioannidis JPA (2005) Why Most Published Research Findings Are False. PLoS Med 2(8): e124.

# The Problem of Reproducibility in Scientific Research

- Non-reproducibility of rates of **65–89%** in pharmacological studies and **64%** in psychological studies.

- Problem of <u>multiple testing</u> in high-dimensional experiments. For gene expression analyses, 26 of 36 (72%) genomic associations initially reported as significant were found to be **over-estimates of the true effect** when tested *in other datasets*

- Analytic focus has been on **significance** (P) values rather than **effect size** or **independent verification**.

# Support and Gap Analysis using Open Data and Open Services

- HyQue is a platform for knowledge discovery that uses data retrieval coupled with contradiction-based automated reasoning to validate scientific hypotheses

- Leverages semantic technologies to provide access to linked data, ontologies, and semantic web services

- Uses positive and negative findings, captures provenance

- Weighs evidence according to context

- Used to find **aging genes** in worm, assess **cardiotoxicity** of tyrosine kinase inhibitors

HyQue: evaluating hypotheses using Semantic Web technologies. J Biomed Semantics. 2011 May 17;2 Suppl 2:S3.

Evaluating scientific hypotheses using the SPARQL Inferencing Notation. Extended Semantic Web Conference (ESWC 2012). Heraklion, Crete. May 27-31, 2012.

Table 3 8 *C. elegans* genes that received the highest HyQue evaluations for their role in aging, the PubMed identifiers of papers describing their roles in regulating longevity, and the data evaluation functions that contributed to their scores

| WormBase identifier | Symbol | Score | PMID | Satisfied data evaluation function | | | | | | | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | | | | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 |
| WBGene00008205 | sams-1 | 0.89 | 16103914 | ✓ | ✓ | ✓ | ✓ | ✓ | ✓ | | ✓ | ✓ |
| WBGene00000371 | cco-1 | 0.78 | 21215371 | ✓ | ✓ | | | ✓ | ✓ | ✓ | ✓ | ✓ |
| WBGene00009741 | drr-1 | 0.78 | 16103914 | ✓ | ✓ | | ✓ | ✓ | ✓ | | ✓ | ✓ |
| WBGene00002178 | jnk-1 | 0.78 | 15767565 | ✓ | ✓ | | | ✓ | ✓ | ✓ | ✓ | ✓ |
| WBGene00004013 | pha-4 | 0.78 | 19239417 | | ✓ | | ✓ | ✓ | ✓ | ✓ | ✓ | ✓ |
| WBGene00004789 | sgk-1 | 0.78 | 15068796 | ✓ | ✓ | | | ✓ | ✓ | ✓ | ✓ | ✓ |
| WBGene00004800 | sir-2.1 | 0.78 | 21938067 | ✓ | | | ✓ | ✓ | ✓ | ✓ | ✓ | ✓ |
| WBGene00006796 | unc-62 | 0.78 | 17411345 | ✓ | ✓ | | | ✓ | ✓ | ✓ | ✓ | ✓ |

Text co-mentions
Gene ontology annotations
Differential gene expression
…

Access to open data
Linked to ontologies
Represented with a universal language
Queried using a portable language
Results stored with their provenance

Table 4 31 highest scoring *C. elegans* genes that received HyQue evaluation scores for their role in aging without existing aging-related annotations, and the data evaluation functions that contributed to their scores

| WormBase identifier | Symbol | Satisfied data evaluation function | | | | | | | | |
|---|---|---|---|---|---|---|---|---|---|---|
| | | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 |
| WBGene00000252 | bli-2 | | ✓ | | | | | ✓ | ✓ | ✓ |
| WBGene00000255 | bli-5 | | ✓ | | | | ✓ | | ✓ | ✓ |
| WBGene00000262 | bra-1 | | ✓ | | | | | ✓ | ✓ | ✓ |
| WBGene00000479 | cgh-1 | | | | | | | ✓ | ✓ | ✓ |
| WBGene00000915* | daf-21 | | | | | | | ✓ | ✓ | ✓ |
| WBGene00001165 | efn-4 | | ✓ | | | | | ✓ | ✓ | ✓ |
| WBGene00001428* | fkb-3 | | ✓ | | | | ✓ | | ✓ | ✓ |
| WBGene00001543* | gcy-18 | | ✓ | | | | ✓ | | ✓ | ✓ |
| WBGene00001578 | ges-1 | | ✓ | | | | ✓ | | ✓ | ✓ |
| WBGene00001746 | gsk-3 | | ✓ | | | | | ✓ | ✓ | ✓ |
| WBGene00001824 | hbl-1 | | ✓ | | | | ✓ | | ✓ | ✓ |
| WBGene00001974 | hmg-4 | | ✓ | | | | | ✓ | ✓ | ✓ |
| WBGene00001979 | hmp-2 | | ✓ | | | | | ✓ | ✓ | ✓ |
| WBGene00002005* | hsp-1 | | | | ✓ | | ✓ | | ✓ | ✓ |
| WBGene00002013* | hsp-12.6 | | ✓ | | ✓ | | | | ✓ | ✓ |
| WBGene00002069* | ikb-1 | | ✓ | | ✓ | | | | ✓ | ✓ |
| WBGene00002881 | let-756 | | ✓ | | | | | ✓ | ✓ | ✓ |
| WBGene00003029 | lin-44 | | ✓ | | | | | ✓ | ✓ | ✓ |
| WBGene00003058 | lov-1 | | ✓ | | | | | ✓ | ✓ | ✓ |
| WBGene00003210 | mel-28 | | | | | | ✓ | ✓ | ✓ | ✓ |
| WBGene00003473 | mtl-1 | | ✓ | | | | ✓ | | ✓ | ✓ |
| WBGene00003497 | mup-4 | | ✓ | | | | | ✓ | ✓ | ✓ |
| WBGene00003977* | pes-2.1 | | ✓ | | | | ✓ | | | ✓ |
| WBGene00004392 | rnr-2 | | | | | | ✓ | ✓ | ✓ | ✓ |
| WBGene00004765 | sel-8 | | ✓ | | | | | ✓ | ✓ | ✓ |
| WBGene00006789 | unc-54 | | ✓ | | | | | ✓ | ✓ | ✓ |
| WBGene00007036 | sod-5 | | ✓ | | | | ✓ | | ✓ | ✓ |

# Scaling Validation

- Automated experimentation (Adam & Eve)
- Crowdsourcing
  - As a simple task
  - As an open problem
- Automated discovery of viable methods
- Automated implementation of viable methods

# Key Research Challenges

- Scalable, shared, fault-tolerant, and readily re-deployable frameworks for **archiving** and **providing versioned and maximally FAIR biomedical (meta)data**

- Scalable methods for the *prospective* and *retrospective* **authoring, assessment, and repair of metadata**.

- Scalable methods to *learn* equivalent **representational patterns**

- Scalable frameworks for *open, transparent, reproducible* and *recurrent* **analysis** and **meta-analysis** of FAIR research data.

- Methods to identify **investigative** *biases* and **knowledge** *gaps*

- Scalable and reliable methods **for the prioritization scientific hypotheses using evidence gathered across scales and sources**

- Scalable methods for **validation** of research findings.