



Health-e-Child Final Conference

Sestri Levante, Italy

April 24, 2010

Present and Future of Biomedical Ontologies



Olivier Bodenreider

Lister Hill National Center
for Biomedical Communications
Bethesda, Maryland - USA

Outline

- ◆ A few legitimate questions
 - What do you mean by *ontology*?
 - Why a talk on biomedical ontologies at the Health-e-Child conference?
- ◆ Biomedical ontologies
 - (A quick look at) The past
 - The present
 - The future

[Bodenreider, Brief Bioinf 2006]



What do you mean by *ontology*?

Ontology vs. other artifacts

- ◆ Ontology
 - Defining types of things and their relations
- ◆ Terminology
 - Naming things in a domain
- ◆ Thesaurus
 - Organizing things for a given purpose
- ◆ Classification
 - Placing things into (arbitrary) classes
- ◆ Knowledge bases
 - Assertional vs. definitional knowledge



Ontology vs. other artifacts (revisited)

- ◆ Lexical and terminological resources
 - Mostly collections of names for biomedical entities
 - Often have some kind of hierarchical organization (e.g., relations)
- ◆ Ontological resources
 - Mostly collections of relations among biomedical entities
 - Sometimes also collect names

“Ontological spectrum”



Why a talk on biomedical ontologies
at the Health-e-Child conference?

Ontology in HeC presentations

◆ Nomenclatures

- O. Milenasi
 - International Paediatric and Congenital Cardiac Codes
 - Normalization efforts (in Europe and the EU)
 - Harmonization of the two nomenclatures

◆ Ontologies

- A. Tsymbol
 - Gene Ontology, KEGG
 - Reasoning based on ontologies (e.g., semantic similarity)
- A. Everett
 - Abstract clinical information from patient records
 - Facilitate the recruitment of patients for clinical trials



HeC and ontology

- ◆ HeC
 - Outcomes, diagnoses, procedures
 - Personalized medicine
- ◆ Sharing information requires **normalization**
 - Among healthcare practitioners
 - Through clinical research databases
 - Evidence-based medicine
 - Comparative effectiveness
- ◆ Analyzing information requires **aggregation**
 - Compensate for differences in granularity



Genotype vs. phenotype

◆ Genotype information

- “Exact” measurement
(+ context)
- Can be easily analyzed through mathematical models
 - Micro-arrays
 - Sequence similarity
 - SNP patterns

◆ Phenotype information

- Results from human observation
(+ interpretation / context)
- Requires normalization
- May require aggregation for analysis



Biomedical ontologies

(A quick look at) The Past

To support a theory of diseases

◆ Hippocrates

- Dismisses superstition
- Four humors
 - Blood
 - Phlegm
 - Yellow bile
 - Black bile

◆ Thomas Sydenham (1624-1689)

- *Medical observations on the history and cure of acute diseases (1676)*

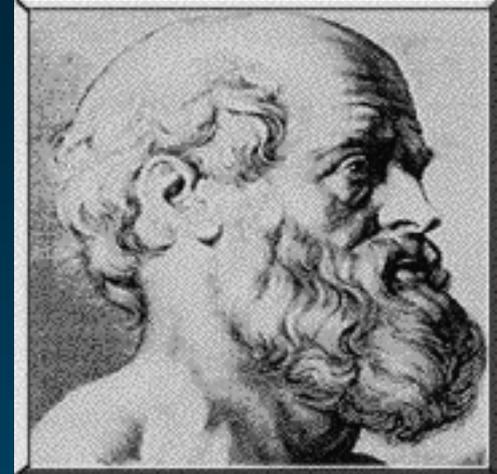
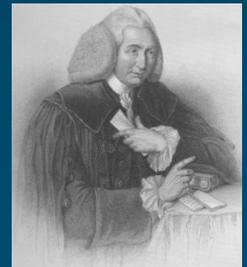
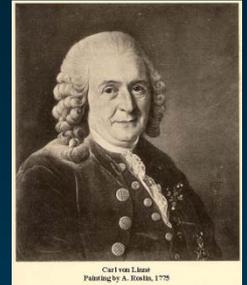


Figure 36 Thomas Sydenham (1624-1689)

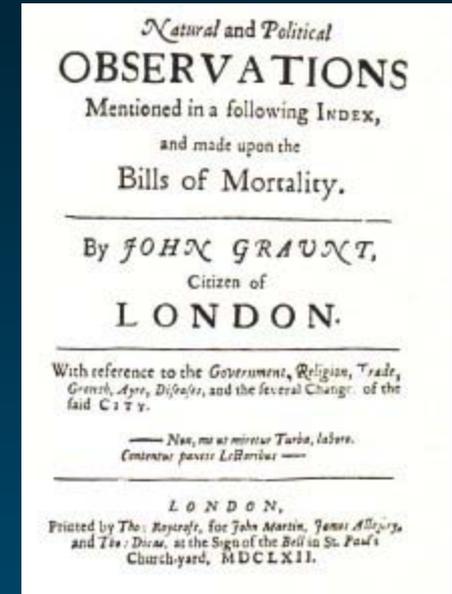
To classify diseases (and plants)

- ◆ Carolus Linnaeus (1707-1778)
 - *Genera Plantarum* (1737)
 - *Genera Morborum* (1763)
- ◆ François Boissier de La Croix
a.k.a. F. B. de Sauvages (1706-1767)
 - *Methodus Foliorum* (1751)
 - *Nosologia Methodica* (1763/68)
- ◆ William Cullen (1710-1790)
 - *Synopsis Nosologiae Methodicae* (1785)



To support epidemiology

- ◆ John Graunt (1620-1674)
 - Analyzes the vital statistics of the citizens of London
- ◆ William Farr (1807-1883)
 - Medical statistician
 - Improves Cullen's classification
 - Contributes to creating ICD
- ◆ Jacques Berthillon (1851-1922)
 - Chief of the statistical services (Paris)
 - Classification of causes of death (161 rubrics)



London Bills of Mortality

LONDON'S Dreadful Visitation:
Or, A COLLECTION of All the
Bills of Mortality
 For this Present Year:
 Beginning the 27th of December 1664. and
 ending the 19th of December following:
 As also, The GENERAL or whole years BILL:
 According to the Report made to the
 KING'S Most Excellent Majesty,
 By the Company of Parish-Clerks of London. &c

LONDON:
 Printed and are to be sold by E. Cotes living in Aldersgate-street.
 Printer to the said Company 1665.

A general Bill for this present year, ending the 19 of December 1665. according to the Report made to the KING'S most Excellent Majesty. By the Company of Parish Clerks of London, &c.

The Diseases and Casualties this year.

A Bortive and Stillborne	517	Executed	21	Palfie	30
Aged	1545	Flux and Small Pox	655	Plague	68598
Aque and Peaver	5257	Found dead in Streets, fields, &c.	2	Plasmod	6
Apoplex and Suddenly	116	French Pox	86	Pluritic	19
Bedric	10	Frighted	23	Posioned	4
Blind	1	Gout and Sciatica	27	Quinse	35
Bleeding	16	Grief	46	Rickets	137
Bloody Flux, Scouring & Flux	185	Griping in the Guts	228	Killing of the Lights	397
Burnt and Scalded	8	Hang'd & made away themselves	7	Lapitate	14
Colicure	3	Headmole shot & Moxie fallen	14	Scurvy	109
Cancer, Gangrene and Fillula	56	jaundies	120	Shingles and Swine pox	2
Canker, and Thrush	121	Imposume	227	Sores, Ulcers, broken and healed	82
Childbed	623	Kill'd by severall accidentes	46	Lambs	82
Christomes and Infants	1258	Kings Evill	28	Spleen	14
Cold and Cough	62	Leproric	2	Spotted Fever and Purples	1929
Collick and Winde	124	Lechary	14	Scopping of the stomack	332
Consumption and Tiflick	4808	Liverg-town	21	Stone and Stranguy	28
Convulsion and Morice	1052	Meagrom and Headach	1	Sucket	1100
Distacted	3	Mealles	7	Teeth and Worms	1014
Droove and Terpany	1476	Mothered and Shot	9	Worming	51
Drwaed	3	Overjaed & Starved	45	Vunn	7
Colicures	5114				
Childbed	4853				
In all	9567				
		Males	48569		
		Females	48737		
		In all	97306		
		Of the Plague		68598	
Increased in the Burials in the 130 Parishes and at the Pest-houses this year				79009	
Decreased of the Plague in the 130 Parishes and at the Pest-houses this year				88590	

Limitations of existing classifications

“The advantages of a uniform statistical nomenclature, however imperfect, are so obvious, that it is surprising no attention has been paid to its enforcement in Bills of Mortality. Each disease has, in many instances, been denoted by three or four terms, and each term has been applied to as many different diseases: vague, inconvenient names have been employed, or complications have been registered instead of primary diseases. The nomenclature is of as much importance in this department of inquiry as weights and measures in the physical sciences, and should be settled without delay.”

– William Farr

First annual report.

London, Registrar General of England and Wales, 1839, p. 99.

Biomedical ontologies

The Present

Many biomedical ontologies

- ◆ About 200 biomedical ontologies available in various repositories
- ◆ Over 2M biomedical concepts
- ◆ Hundreds of millions of relations among them
- ◆ Limited interoperability
- ◆ Quality assurance issues

Many biomedical ontologies

◆ General vocabularies

- anatomy (FMA, Neuronames)
- drugs (RxNorm, First DataBank, Micromedex)
- medical devices (UMD, SPN)

◆ Several perspectives

- clinical terms (SNOMED CT)
- information sciences (MeSH, CRISP)
- administrative terminologies (ICD-9-CM, CPT-4)
- data exchange terminologies (HL7, LOINC)



Many biomedical ontologies (cont'd)

◆ Specialized vocabularies

- nursing (NIC, NOC, NANDA, Omaha, PCDS)
- dentistry (CDT)
- oncology (PDQ)
- psychiatry (DSM, APA)
- adverse reactions (MedDRA, WHO ART)
- primary care (ICPC)

◆ Terminology of knowledge bases (OMIM, QMR)



Too many biomedical ontologies?

- ◆ Examples of exotic or obsolete ontologies in biomedical ontology repositories
- ◆ Governance issues
 - e.g., Ontology developed by a doctoral student

Uses of biomedical ontologies

- ◆ Knowledge management
 - Annotating data and resources
 - Accessing biomedical information
 - Mapping across biomedical ontologies
- ◆ Data integration, exchange and semantic interoperability
- ◆ Decision support
 - Data selection and aggregation
 - Decision support
 - Natural language processing (NLP) applications
 - Knowledge discovery

[Bodenreider, YBMI 2008]



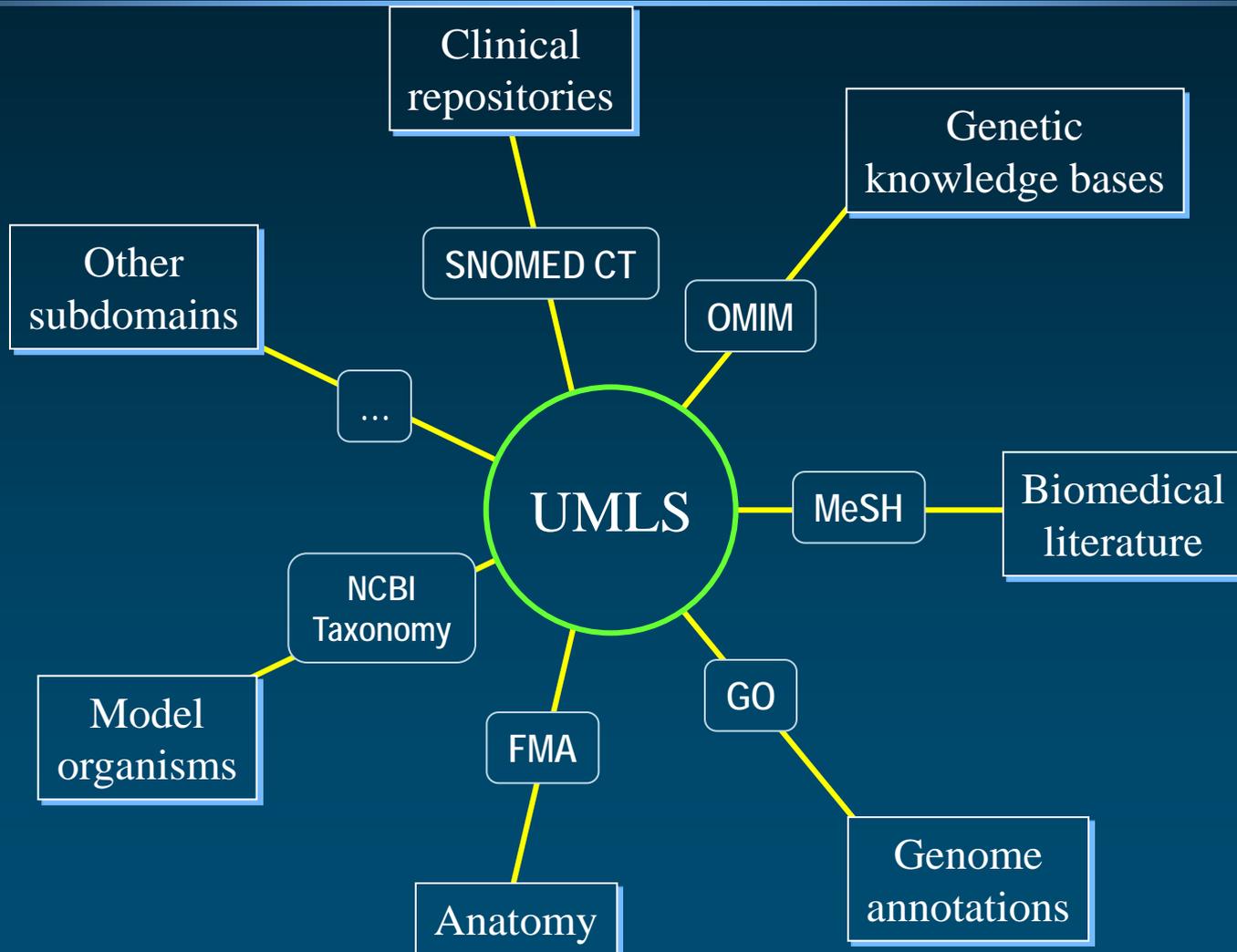
Development

- ◆ Still mostly uncoordinated
 - “Cottage industry”
 - Issues
 - Redundancy
 - Lack of consistence
 - Need for mapping
 - Exception: OBO Foundry
- ◆ Knowledge representation technology
 - Move towards description logics (e.g., OWL)
 - e.g., SNOMED CT [+ OBO ontologies]

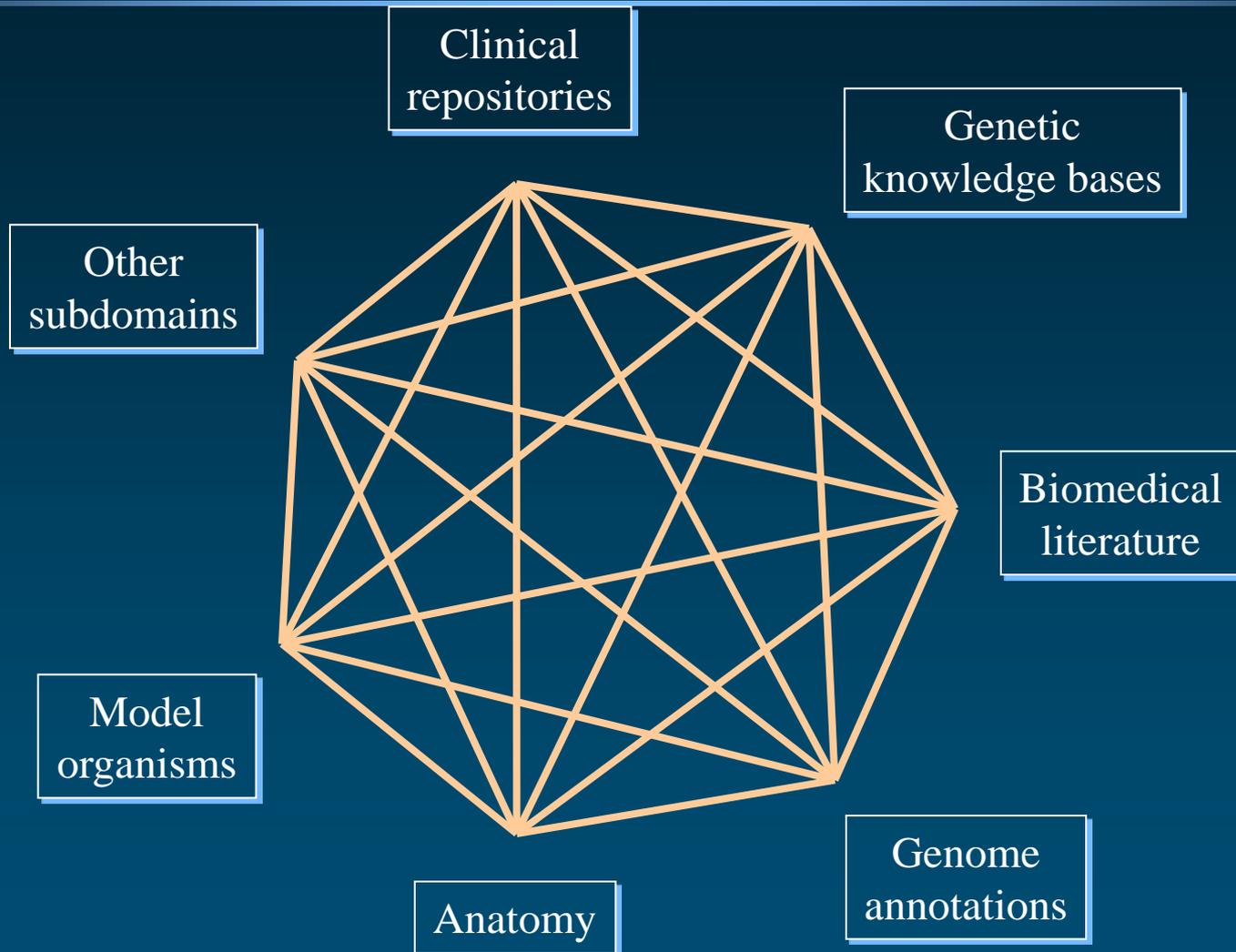
Loose integration

- ◆ Pairwise mappings
 - Unidirectional
 - Specific to a given purpose
 - Costly to create and maintain
- ◆ Integration through a reference
 - “Interlingua”
 - Identify which terms from different ontologies name the same entities and link them together
 - e.g., Unified Medical Language System (UMLS) Metathesaurus

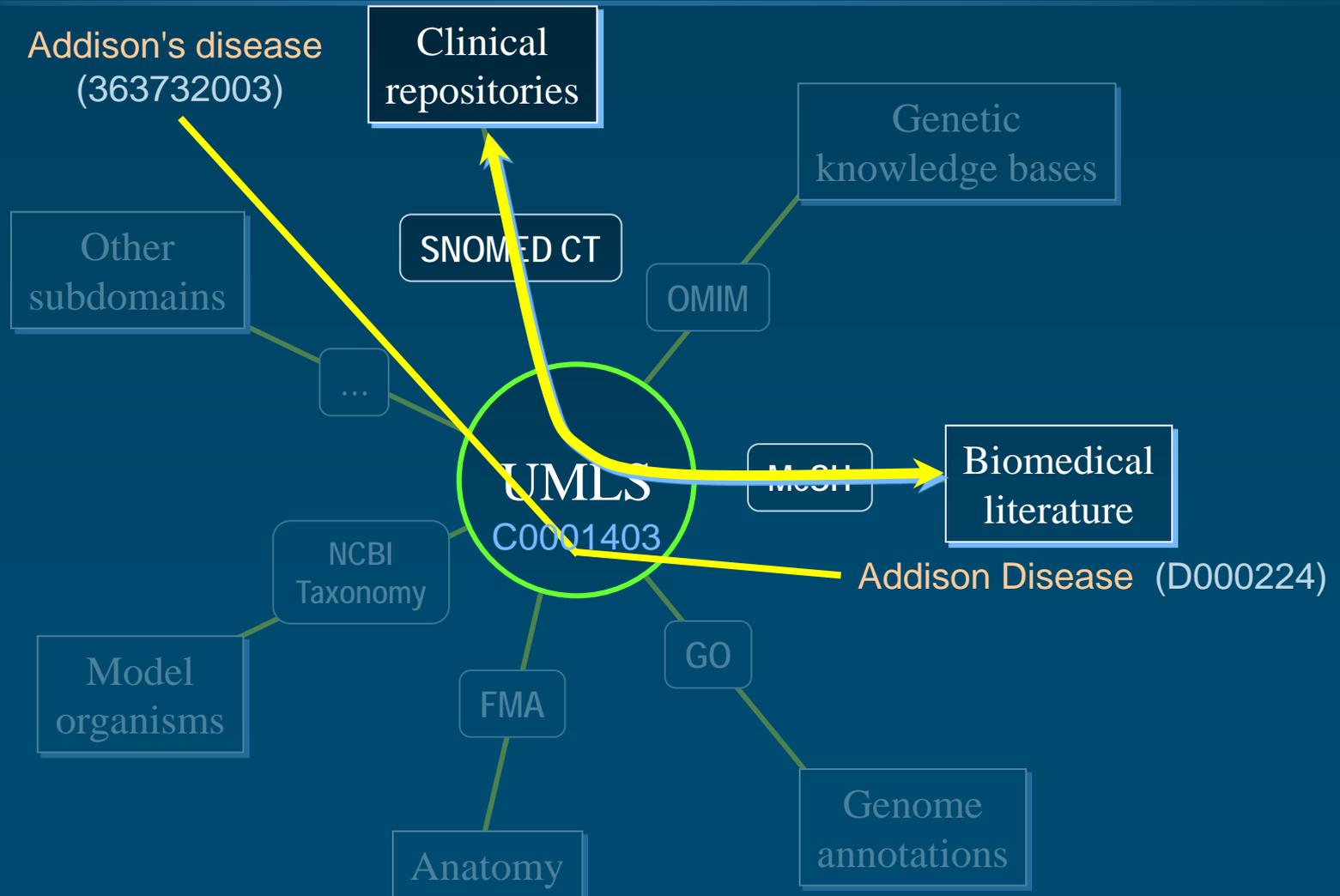
Integrating subdomains



Integrating subdomains



Ontology integration through the UMLS



(Integrated) concept repositories

- ◆ Unified Medical Language System

<http://umlsks.nlm.nih.gov>

- ◆ NCBO's BioPortal

<http://www.bioontology.org/tools/portal/bioportal.html>

- ◆ caDSR

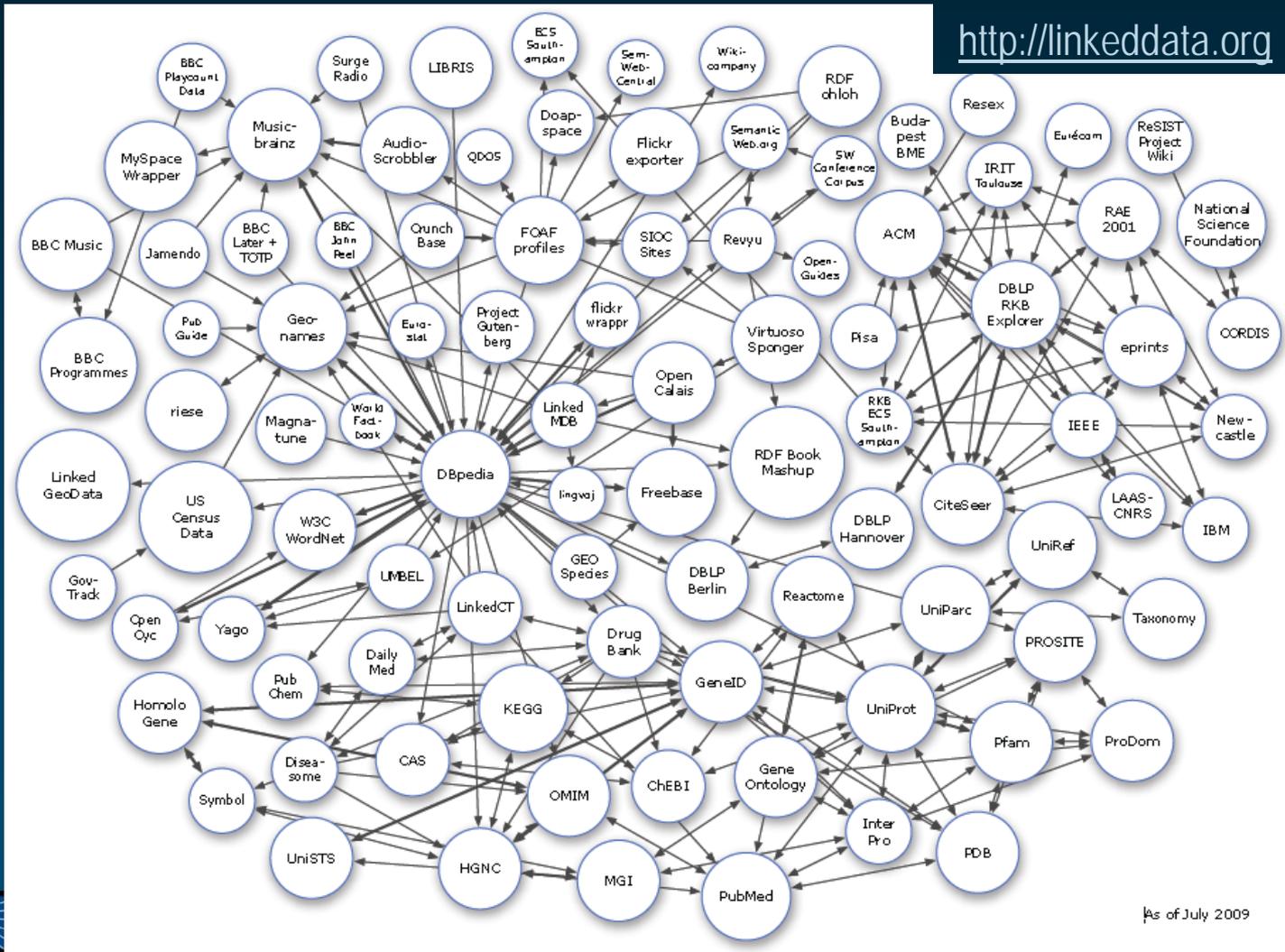
http://ncicb.nci.nih.gov/NCICB/infrastructure/cacore_overview/cadsr

- ◆ Open Biomedical Ontologies (OBO)

<http://obofoundry.org/>

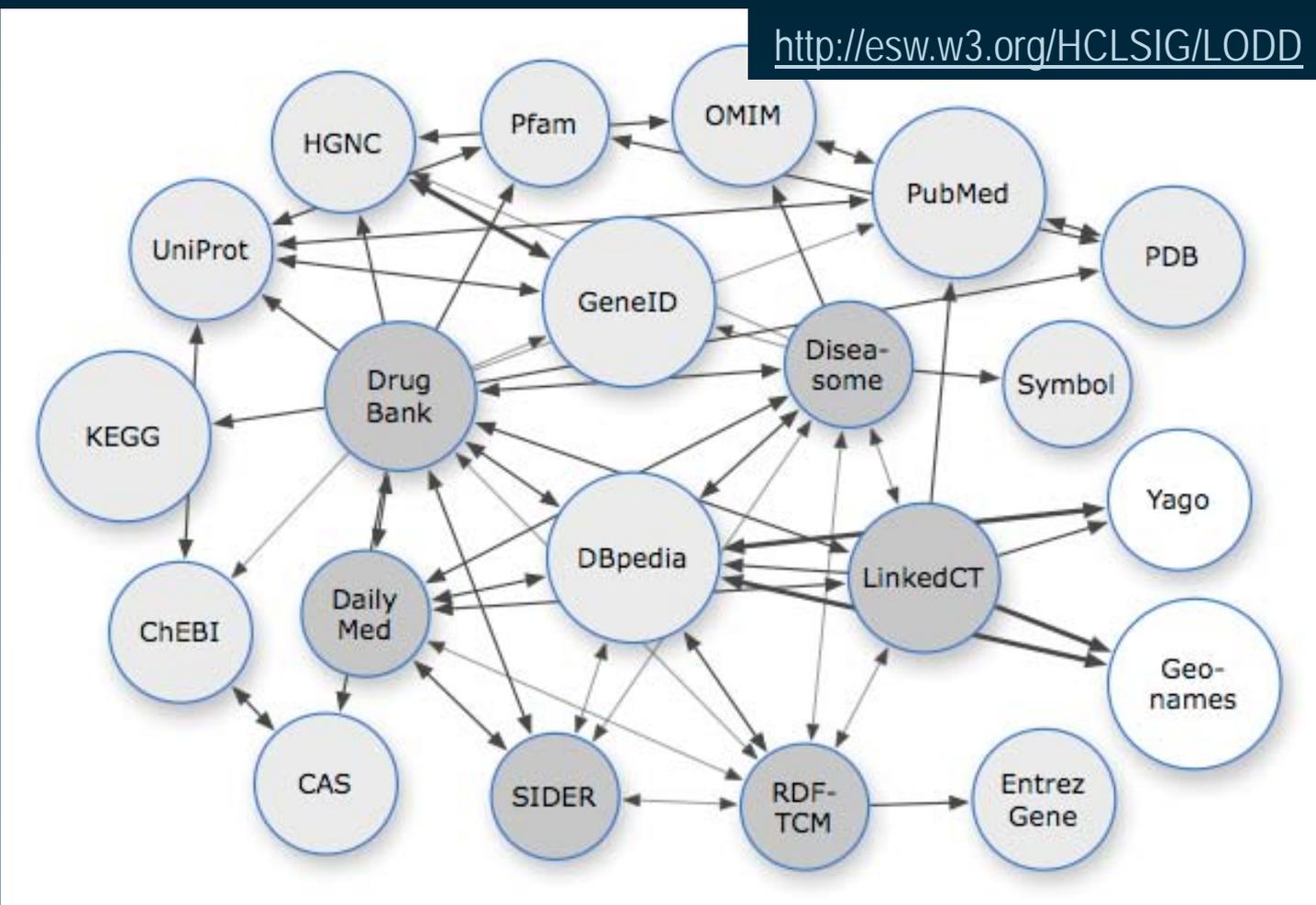


Ontology integration supports data integration



Linked Open Drug Data

<http://esw.w3.org/HCLSIG/LODD>



Linked data

- ◆ Semantic Web
- ◆ Resources available in RDF
 - Unique, unambiguous identifiers for entities
 - Explicit relations among entities
- ◆ Links across resources (federation)
 - Enabled by
 - Shared identifiers across resources
 - Global identifiers, resolvable on the web

Biomedical ontologies

The Future

Harmonization

- ◆ Collaboration among ontology developers
 - Prospectively
 - OBO Foundry model
 - Avoid redundancy
 - Foster collaboration
 - Retrospectively
 - SNOMED CT model
 - Seek agreement with other ontologies for specialized content (e.g., LOINC for observables)
 - Serve as an ontology backbone for classifications (e.g., ICD11)

Harmonization Benefits

- ◆ Fewer pairwise mappings
 - Not needed for concepts of the same level of granularity
 - Computable automatically for finer-grained concepts
- ◆ Increased interoperability
 - Among ontologies
 - Among datasets annotated to these ontologies
 - Among applications using these ontologies



Quality of biomedical ontologies

- ◆ Quality assurance in ontologies is still imperfectly defined
 - Difficult to define outside a use case or application
- ◆ Several approaches to evaluating quality
 - Collaboratively, by users (Web 2.0 approach)
 - Marginal notes enabled by BioPortal
 - Centrally, by experts
 - OBO Foundry approach
- ◆ Related issues
 - Quality of ontology integration (mappings)

Discoverability

- ◆ No universal repositories for biomedical datasets
 - Some datasets made available through portals (NCBI, EBI, NCBO)
- ◆ Ontology repositories
 - UMLS: 153 source vocabularies (biased towards healthcare applications)
 - NCBO BioPortal: 195 ontologies (biased towards biological applications)
 - Limited overlap between the two repositories
- ◆ Need for discovery services
 - Metadata for ontologies and biomedical datasets

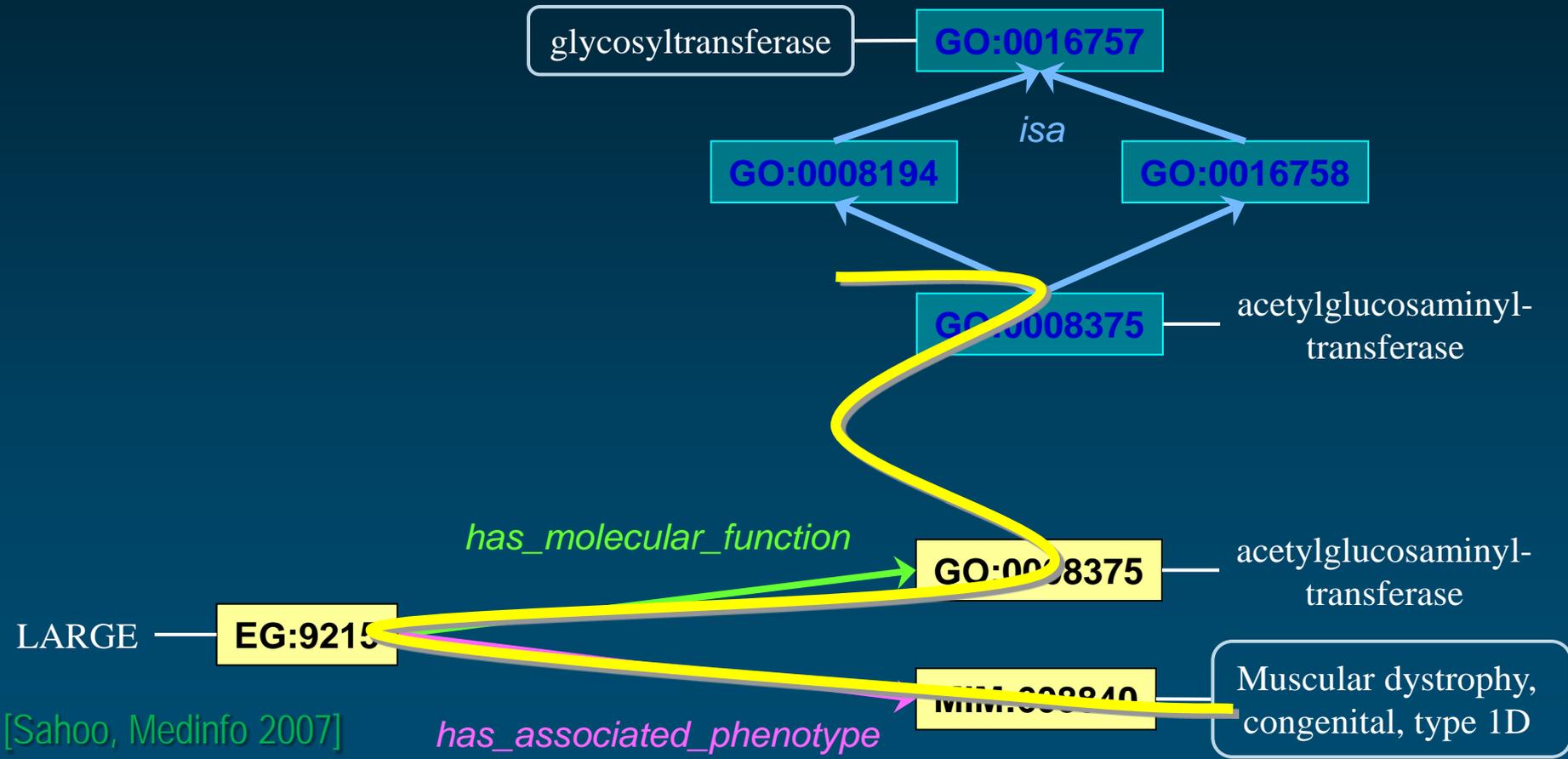
Common upper-level ontologies

- ◆ Formalize high-level ontological distinctions
 - Occurrents/continuants
 - Dependent/independent continuants
- ◆ Can be shared by multiple domain ontologies
- ◆ Make ontologies easier to integrate
 - Fewer essential differences in the organization

Reasoning with ontologies

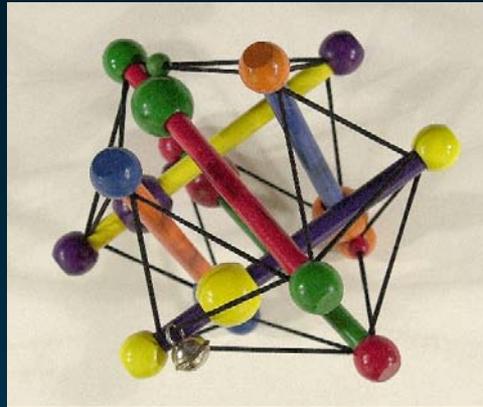
- ◆ Description logic (e.g., OWL) reasoners available, but few biomedical ontologies can fully take advantage of them
 - Limited expressiveness of the ontologies
 - Limited performance of the reasoners
- ◆ Subsumption reasoning
 - Useful for data aggregation
- ◆ Beyond subsumption reasoning
 - Rule-based systems (e.g., for clinical decision support)
 - Hypothesis generation and knowledge discovery

From *glycosyltransferase* to congenital muscular dystrophy



[Sahoo, Medinfo 2007]





Medical Ontology Research

Contact: olivier@nlm.nih.gov

Web: mor.nlm.nih.gov



Olivier Bodenreider

Lister Hill National Center
for Biomedical Communications
Bethesda, Maryland - USA