# How data sharing leads to knowledge

M. Scott Marshall, Ph.D. W3C HCLS IG co-chair Leiden University Medical Center University of Amsterdam <u>http://staff.science.uva.nl/~marshall</u> <u>http://www.w3.org/blog/hcls</u>

### Motivation

→Science is based on **knowledge**: knowledge capture, knowledge sharing, i.e. communication of findings.

→Semantic Web provides a basis for knowledge sharing through machine-readable and *reason-able* annotation of resources.

### What is knowledge ?

"data", "information", "facts", "knowledge"

Knowledge is a statement that can be tested for truth.

(by a machine) Otherwise, computing can't add much

### **RDF** : a web format for knowledge

# RDF is a W3C language to express statements.



Graph of Knowledge:



# The Semantic Web is the New Global Web of Knowledge

It is about standards for publishing, sharing and querying knowledge drawn from diverse sources

> It makes possible the answering sophisticated questions using background knowledge



Source: Michel Dumontier

# Where is biomedical knowledge?

Can be extracted from:

- People
- Literature
- Diagrams
- Clinical reports
- Databases
- Excel sheets

Most of these sources of biomedical knowledge are *not* machine-readable

### Many tasks are still a challenge!

With existing Web and Health IT:

- Find and integrate information
  - "Although a plethora of resources (tools, databases, materials) for neuroscientists is now available on the web, finding these resources among the billions of possible web pages continues to be a challenge." [M. Martone, NCBO Seminar Series, 4 Nov 2009]
- Make multiple inferences based on background knowledge
  - to obtain more complete answers
  - to discover knowledge

### Examples

- in a medical record system
  - "find all patients whose radiology exhibits a fracture of femur"
- in genomic data

"find all genes annotated with a molecular function or any of its descendants and which is associated with any form of a given disease" (see genes associated with *muscular dystrophy* [Sahoo et al. 2007])

– find, share, annotate images

# Pistoia Alliance Vocabulary Services Initiative

"The life sciences industry currently operates in an environment where few of the basic components of its study (e.g. genes, proteins, cells, diseases, biomarkers, assays, drugs and technologies) are described using consistent, universally agreed-upon vocabularies."

## **Biological and medical ontologies**

😭 🍄 🌈 NCBO BioPortal				🔒 👌 🕶 🗟 👻 🖶	▼  Page ▼	0 <u>u</u> tils ▼ <sup>≫</sup>
O BioPortal Browse	Search	Projects	Annota	ate All Mappings	s All Resou	rces Alpha
NCI Thesaurus						E
Search all ontologies						
Fracture		Search	Categories (	Anatomy	▼]	
Include attributes in search		Help	Groups	All Groups	▼)	
		() Help	Filter	type filter text		
Contains      Exact Match			Ontologie	s: <u>Select All</u> <u>Select Non</u>	<u>ie</u>	
× Clear			Dictyoste	elium discoideum anatom		
			Drosoph	ila gross anatomy (FBbt)		
			Foundati	onal Model of Anatomy (	FMA)	
Selected Ontologies (139):			🔄 Fungal g	ross anatomy (FAO)		
All Ontologies			Human d	developmental anatomy,	abstract ver 🗸	
						<b>T</b>
http://bioportal.bioontology.org/search		•	Internet   Mod	de protégé : désactivé	Ð	ب 100% ح

## Some of the forces at work

- Pharmaceutical industry changing strategy
  - David Cox (Pfizer) Strategy: Academic / Industry partnership, wellness: rare variants that *protect* against disease
  - Pistoia Alliance, Vocabulary Services Initiative
- Personalized Medicine and EHRs
- US NIH NCBCs: NCBO and I2B2
- NCI Semantic Infrastructure
- European Innovative Medicine Initiatives (IMI)

# **Background of the HCLS IG**

- Originally chartered in 2005
  - Chairs: Eric Neumann and Tonya Hongsermeier
- Re-chartered in 2008
  - Chairs: Scott Marshall and Susie Stephens
  - Team contact: Eric Prud'hommeaux
- Broad industry participation
  - Over 100 members
  - Mailing list of over 600
- Background Information
  - <u>http://www.w3.org/blog/hcls</u>
  - <u>http://esw.w3.org/topic/HCLSIG</u>

# **Mission of HCLS IG**

•The mission of HCLS is to develop, advocate for, and support the use of Semantic Web technologies for

- Biological science
- Translational medicine
- Health care

•These domains stand to gain tremendous benefit by adoption of Semantic Web technologies, as they depend on the **interoperability of information** from many domains and processes for efficient decision support

### **Translating across domains**



### **Current Task Forces**

- **BioRDF** federating (neuroscience) knowledge bases
  - M. Scott Marshall (Leiden University Medical Center / University of Amsterdam)
- Clinical Observations Interoperability patient recruitment in trials
  - Vipul Kashyap (Cigna Healthcare)
- Linking Open Drug Data aggregation of Web-based drug data
  - Susie Stephens (Johnson & Johnson)
- Translational Medicine Ontology high level patient-centric ontology
  - Michel Dumontier (Carleton University)
- Scientific Discourse building communities through networking
  - Tim Clark (Harvard University)
- **Terminology** Semantic Web representation of existing resources
  - John Madden (Duke University)

### **BioRDF: Translating across domains**



### Provenance

- Data context (can be experimental context)
- Represent knowledge so that
  - others can discover where a fact (or triple) came from
  - and evaluate how to use it
- link facts to data as evidence

### Provenance types are perspectives on the data



Source: Helena Deus

### A Bottom-up Approach



Source: Helena Deus

### LODD: Translating across domains



### The Classic Web



- Single information space
- HTML describes presentation
- Built on URIs
  - globally unique IDs
  - retrieval mechanism
- Built on Hyperlinks
  - are the glue that holds everything together

### **Linked Data**

Use Semantic Web technologies to publish structured data on the Web and set links between data from one data source and data from another data sources



### **The Linked Data Cloud**



Source: Chris Bizer

### LODD



### Interlinking in LODD



#### http://esw.w3.org/HCLSIG/LODD/Interlinking

### **TripleMap**





### Homonyms

### PSA

- Prostate Specific Antigen
- PSoriatic Arthritis
- alpha-2,8-PolySialic Acid
- PolySubstance Abuse
- Picryl Sulfonic Acid
- Polymeric Silicic Acid
- Partial Sensory Agnosia
- Poultry Science Association



## **Shared Identifiers**

- Must use common URI's in order to link data
- Provenance related identifiers still needed:
  - Identifiers for people (researchers)
  - Identifiers for diseases
  - Identifiers for terms (Terminology servers)
  - Identifiers for programs, processes, workflows
  - Identifiers for chemical compounds
- Shared Names <u>http://sharednames.org</u>
- Bio2RDF

# Early semantic commitment: Map input data to concepts

Input table				Map to concepts	
Entrez Gene	Up or down	. Gene Name		Identifier column:	Entrez Gene
51182	down	HSPA14	▲		
8495	up	PPFIBP2	Ξ	Type of identifier:	Concept name 📃
64207	down	C14orf4			
3977	up	LIFR			Concept name 🔄
5899	down	RALB			Swiss-Prot
9967	up	THRAP3			Entrez-Gene 📐 📃
23534	down	TNP03			омім 📉 🗧
219404	down	MGC9850			FlyBase 📃
79029	up	SPATA5L1			Mouse Genome Dat:
9957	un	HS3ST1	<b>_</b>		Bet Conomo Detako
Add co	lumn Ac	ld row Pa	aste Load		Rai Genome Dalaba
					HUGO 🗾

Screenshot Anni: Martijn Schuemie

### **TMO: Translating across domains**



# Questions & Problems The Drug Development Pipeline



- The road is long, and costly.
- How do we contain costs and develop better drugs?

### Translational Medicine Ontology Mission

- Focuses on the development of a high level patient-centric ontology for the pharmaceutical industry. The ontology should enable data integration across discovery research, hypothesis management, experimental studies, compounds, formulation, drug development, market size, competitive data, population data, etc. This would enable scientists to answer new questions, and to answer existing scientific questions more quickly.
- This will help pharmaceutical companies to model patient-centric information, which is essential for the tailoring of drugs, and for early detection of compounds that may have sub-optimal safety profiles. The ontology should **link to existing publicly available domain ontologies**.

### Scope of the TMO



### **TMO Structure**



### **Translational Medicine KB**



### **TMO Query**



How many patients experienced side effects while taking Donepezil?

### **Discovery Questions and Answers**

What genes are associated with or implicated in AD?	Diseasome and PharmGKB indicate at least 97 genes have some association with AD.
Which SNPs may be potential AD biomarkers?	PharmGKB reveals 63 SNPs.
Which market drugs might potentially be repurposed for AD because they modulate AD implicated genes?	57 compounds or classes of compounds are used to treat 45 diseases, including AD, diabetes, obesity, and hyper/hypotension

### Clinical Trials Questions and Answers

Since my patient is suffering from drug- induced side effects for AD treatment, can an AD clinical trial with a different mechanism of action be identified?	Of the 438 drugs linked to AD trials, only 58 are in active trials and only 2 (Doxorubicin and IL-2) have a documented mechanism of action. 78 AD-associated drugs have an established MOA.
Find AD patients without the APOE4 allele as these would be good candidates for the clinical trial involving Bapineuzumab?	Of the 4 patients with AD, only one does not carry the APOE4 allele, and may be a good candidate for the clinical trial.
What active trials are ongoing that would be a good fit for Patient 2?	58 Alzheimer trials, 2 mild cognitive impairment trials, 1 hypercholesterolaemia trial, 66 myocardial infarction trials, 46 anxiety trials, and 126 depression trials.

### Physician Questions and Answers

What are the diagnostic criteria for AD?	There are 12 diagnostic inclusion criteria and 9 exclusion criteria
Does Medicare D cover Dopenezil?	Medicare D covers two brand name formulations of Donepezil: Aricept and Aricept ODT.
Have any AD patients been treated for other neurological conditions?	Patient 2 was found to suffer from AD and depression.

# Terminology: Translating across domains



# Terminology Ongoing Work

- RDF representation of clinical reports
- Mammogram: Represent both radiology and pathology report to discover discrepancies
- Use Translational Medicine Ontology, RadLex, SNOMED in the RDF
- Link to data about biomarkers and therapies

• There is a 1.2 cm x 1.3 cm round mass with an indistinct margin in the left breast at 9 o'clock. This round mass is isoechoic.

# these are pretty commonplace assertion types
 # hard to imagine much structural variation here
 # can pick a specific vocabulary later
 theMassAt9 size [dimension

[val "1.2"^^xsd:float; unit cm],

[val

```
"1.3"^^xsd:float; unit cm]].
```

# more sketchy, depends on chosen anatomy vocabulary
theMassAt9 location [a Location;

in thePatient,

[a Breast; laterality left],

[rdfs:label "9:00"]].

```
# very sketchy
# these are modeled loosey-goosey as value partitions
# need no idea how e.g. Radlex might do this
theMassAt9 shape round.
theMassAt9 margin indistinct.
```

Source: John Madden

# Barriers to data sharing: social, legal, and technical

- "Biologists would rather share their toothbrush than share a gene name"
  - Don't want to get "scooped" for a publication and potentially lose years of work and Ph.D. material.
  - Competition for grants
- Need clarity and transparency about threats to patient privacy
- Many data formats, example: CDISC and HL7 RIM
- Most researchers do not feel the need to look at data from neighboring domains (cross-disciplinary studies)

### Summary

- The data landscape for personalized medicine is highly fragmented
- Public vocabulary services can be used to connect data sets and make them accessible on the Web
- Data sharing can add value to data through linking
- Best practices for important data sources: microarray data, image data
- Data stewardship serve data back to community

### Acknowledgements

- W3C Health Care and Life Science Interest Group, http://www.w3.org/blog/hcls
- National Center for Biomedical Ontologies
- Concept Web Alliance
- Authors of all contributed slides

### The End

"Science is built up of facts, as a house is built of stones; but an accumulation of facts is no more a science than a heap of stones is a house."

*– Henri Poincaré, Science and Hypothesis, 1905*