

Dataspaces: The Tutorial

Alon Halevy, David Maier
VLDB 2008
Auckland, New Zealand



Outline

- Dataspaces: why? What are they?
 - Examples and motivation
- Dataspace techniques:
 - Locating and understanding data sources
 - Creating mappings and mediated schemas
 - Pay-as-you-go: improving with time
 - Query processing for dataspace
- Research challenges on specific dataspace:
 - Science, the desktop, the Web

DBLife: a community dataspace

[Doan et al., U. Wisc; CIMPLE, UW + Y!]

- A data space for the database research community
- Started with 846 data sources in May 2005
 - researcher homepages, CS dept homepages, etc
- Immediately provided some basic service
 - crawl sources daily to obtain 11000+ pages
 - index & provide keyword search
- Incrementally extract & integrate data
 - provide more services & better services
 - leverage user feedback to further improve the system

Example Service: Create SuperHomepages



The Presidents of the USA - Enchanted Learning.com - Mozilla Firefox

http://www.enchantedlearning.com/history/us/pres/list.shtml

As a thank-you bonus, site members have access to a banner-ad-free version of the site, with print-friendly pages.

(Already a member? [Click here.](#))

[US Flags](#) [US History](#) [US Geography](#)

[A](#) [B](#) [C](#) [D](#) [E](#) [F](#) [G](#) [H](#) [I](#) [J](#) [K](#) [L](#) [M](#) [N](#) [O](#) [P](#) [Q](#) [R](#) [S](#) [T](#) [U](#) [V](#) [W](#) [X](#) [Y](#) [Z](#)

[African-American](#) [Artists](#) [Explorers of the US](#) [Inventors](#) [US Presidents](#) [US Symbols](#) [US States](#)

[President's Day Activities](#) [Enchanted Learning.com](#)

The Presidents of the United States of America

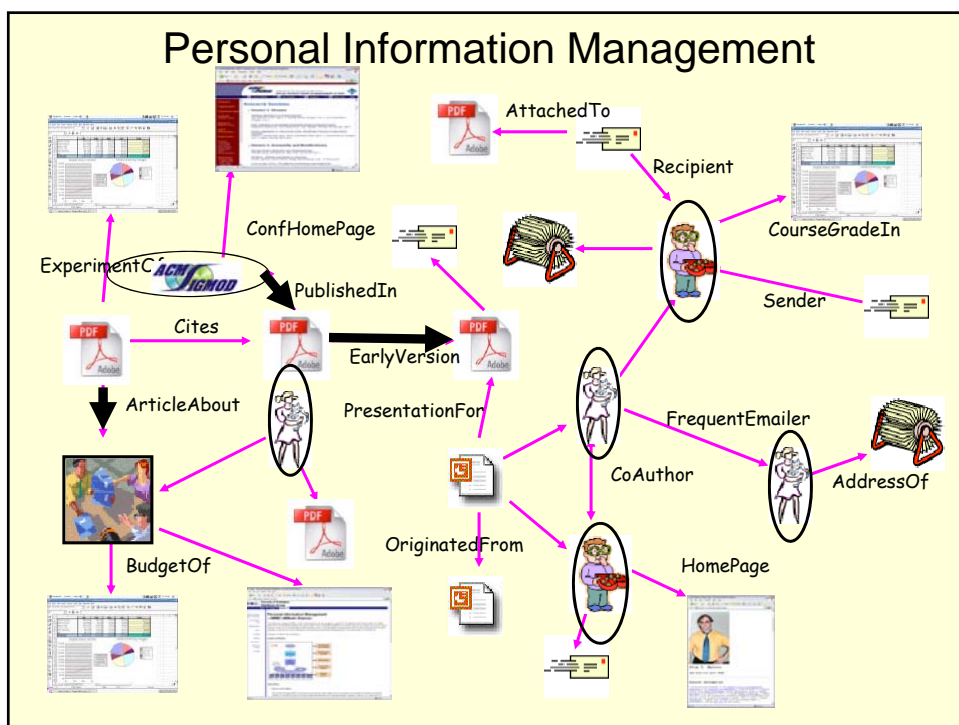
[In the order in which they served](#) [Alphabetical order](#) [Short table of Data](#)

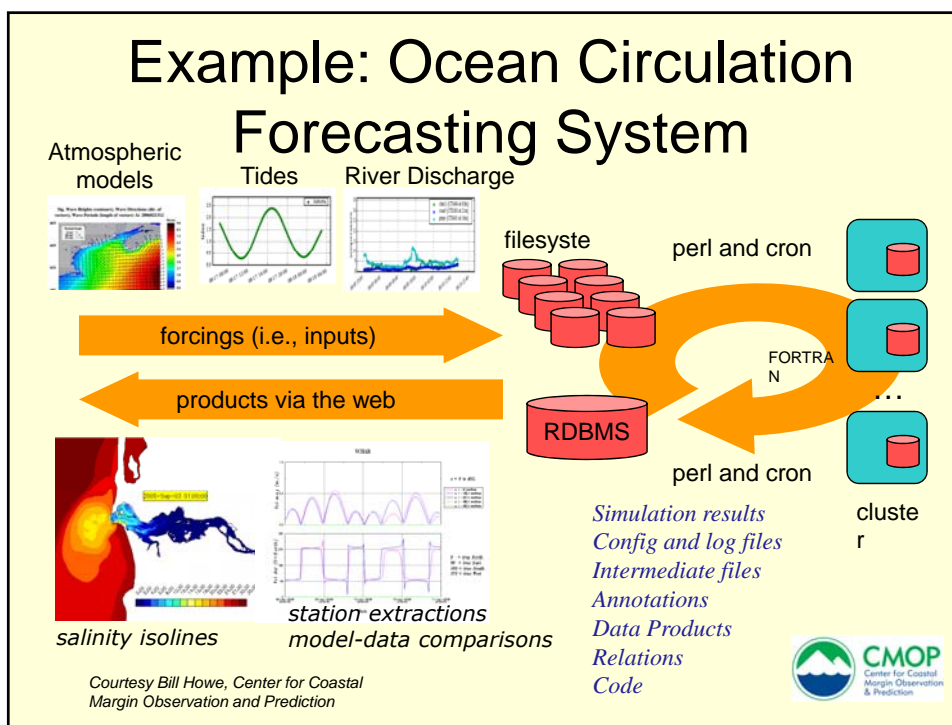
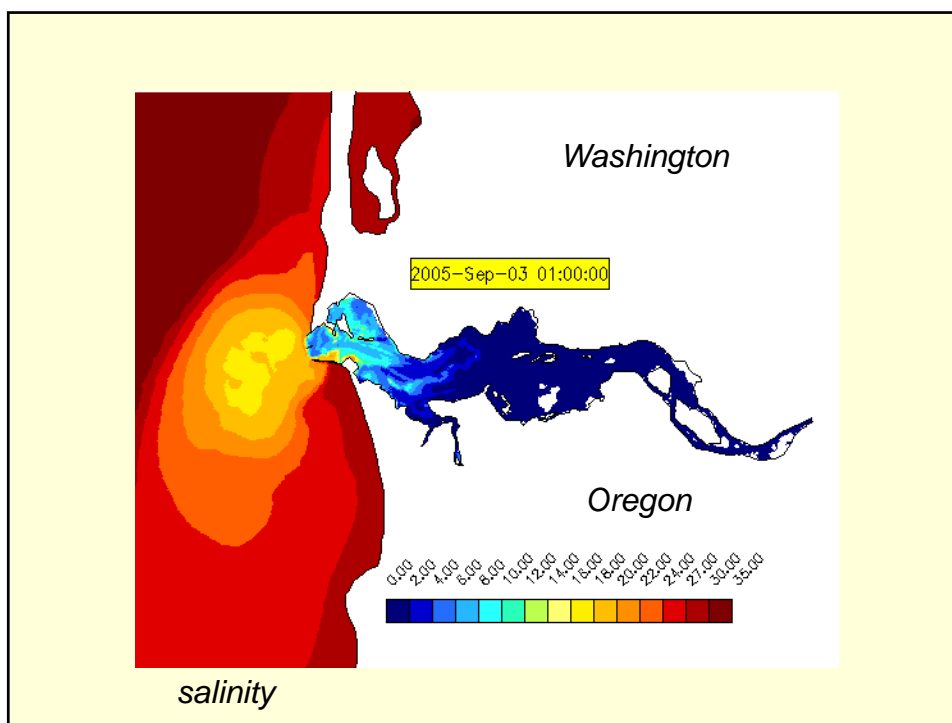
[Abraham Lincoln](#)

The President and Vice-President are elected every four years. They must be at least 35 years of age, they must be native-born citizens of the United States, and they must have been residents of the U.S. for at least 14 years. (Also, a person cannot be elected to a third term as President.)

President	Party	Term as President	Vice-President
1. George Washington (1732-1799)	None, Federalist	1789-1797	John Adams
2. John Adams (1735-1826)	Federalist	1797-1801	Thomas Jefferson
3. Thomas Jefferson (1743-1826)	Democratic-Republican	1801-1809	Aaron Burr, George Clinton
4. James Madison (1751-1836)	Democratic-Republican	1809-1817	George Clinton, Elbridge Gerry
5. James Monroe (1758-1831)	Democratic-Republican	1817-1825	Daniel Tompkins
6. John Quincy Adams (1767-1848)	Democratic-Republican	1825-1829	John Calhoun
7. Andrew Jackson (1767-1845)	Democrat	1829-1837	John Calhoun, Martin van Buren
8. Martin van Buren (1782-1862)	Democrat	1837-1841	Richard Johnson
9. William H. Harrison (1773-1841)	Whig	1841	John Tyler
10. John Tyler (1790-1862)	Whig	1841-1845	
11. James K. Polk (1795-1849)	Democrat	1845-1849	George Dallas
12. Zachary Taylor (1784-1850)	Whig	1849-1850	Millard Fillmore
13. Millard Fillmore (1800-1874)	Whig	1850-1853	
14. Franklin Pierce (1804-1869)	Democrat	1853-1857	
15. James Buchanan (1791-1868)	Democrat	1857-1861	

Cafarella et al, VLDB 2008

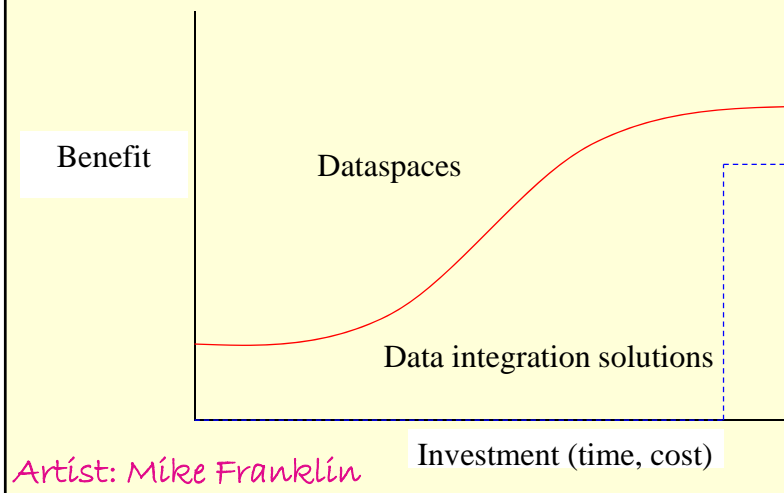




System Desiderata

- Manage all the data I have, not only what I explicitly put into it via a defined schema
 - Start working with data *where is, as is*.
- Pay-as-you-go:
 - Little or no setup time.
 - Provide *best-effort* services
 - Infers semantics to improve services. Gets help where it's most useful.

Key Principle



Dataspaces: A System and a Philosophy

- New kinds of systems:
 - A new kind of data integration system
 - Online data sharing and discussion systems
- Principles applicable anywhere:
 - Best-effort querying, consistency, ...
 - Pay-as-you-go data management
- This tutorial: introduces the principles in the context of data integration.

The Buzz on the Street

A significant long-term goal for our community is to transition from managing traditional databases consisting of well-defined schemata for structured business data, to the much more challenging task of managing a rich collection of structured, semi-structured and unstructured data, spread over many repositories in the enterprise and on the Web. This has sometimes been referred to as the challenge of managing dataspace.

Surajit Chaudhuri

The Buzz in the Press

The screenshot shows the redOrbit website interface. At the top left is the redOrbit logo. To its right are four red buttons: LOGIN, SIGN UP, EMAIL, and SUGGESTIONS. Further right is a large advertisement for 1-800-PetMeds, America's Largest Pet Pharmacy, featuring a 'FREE SHIPPING' offer on orders over \$39. Below the ad is a navigation bar with links: HOME, COMMUNITY, NEWS (highlighted), VIDEO, IMAGES, SPACE, SCIENCE, TECH, HEALTH, EDUCATION, FUN, SHOP, and SIT. Under the NEWS link is a sub-menu with categories: Space, Science, Technology (highlighted), Health, General, Sci-fi & Gaming, Oddities, International, Business, Politics, and Education. Below the navigation bar are social media and utility links: E-mail, Print, Comment, Font Size, Digg, del.icio.us, and a link to the redOrbit Knowledge Network. The main content area features the article title 'Dataspaces: Google's New Marching Order?' in bold. Below the title is the posting date 'Posted on: Sunday, 9 December 2007, 06:00 CST' and the author 'By Quint, Barbara'. The article text begins with 'Looks like that 800-pound gorilla is putting on weight again. I recently had a conversation with Stephen E. Arnold, one of the world's leading Google gurus...' and ends with 'Just a Few Predictions'.

redOrbit

LOGIN
SIGN UP
EMAIL
SUGGESTIONS

1-800-PetMeds
America's Largest Pet Pharmacy
on orders over \$39
Excludes refrigerated items

FREE SHIPPING

HOME COMMUNITY **NEWS** VIDEO IMAGES SPACE SCIENCE TECH HEALTH EDUCATION FUN SHOP SIT

Space Science **Technology** Health General Sci-fi & Gaming Oddities International Business Politics Education

E-mail Print Comment Font Size Digg del.icio.us Discuss in redOrbit Knowledge Network

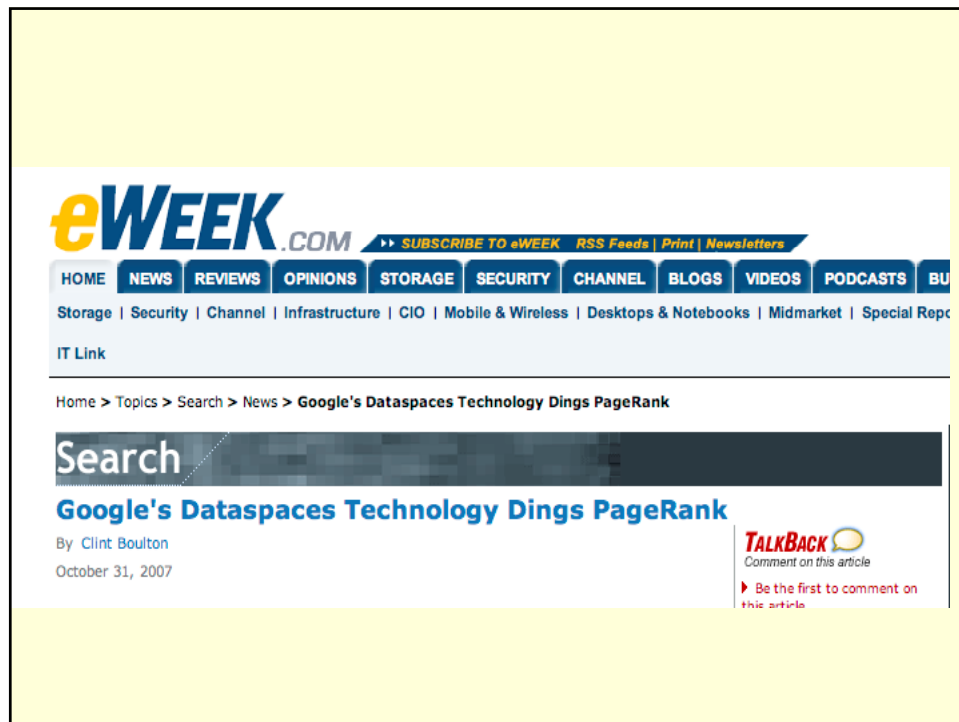
Dataspaces: Google's New Marching Order?

Posted on: Sunday, 9 December 2007, 06:00 CST

By Quint, Barbara

Looks like that 800-pound gorilla is putting on weight again. I recently had a conversation with Stephen E. Arnold, one of the world's leading Google gurus (his latest work on the subject is Google Version 2.0, www.infonortics.com/publications/google/google-pre-dator.html). Arnold draws on the background of a traditional information industry executive with a wide and varied experience. He discovered what he believes to be the emerging infrastructure that Google will use to design future information products, ones that will bring it more users and usage; bridge its products into other technologies, such as mobile phones; bring more ad revenue; and serve the causes of truth, justice, and world peace. In other words, Google has plans in place to do good and to do well, plans that may leave other information professionals gaping in awe or gasping for air.

Just a Few Predictions



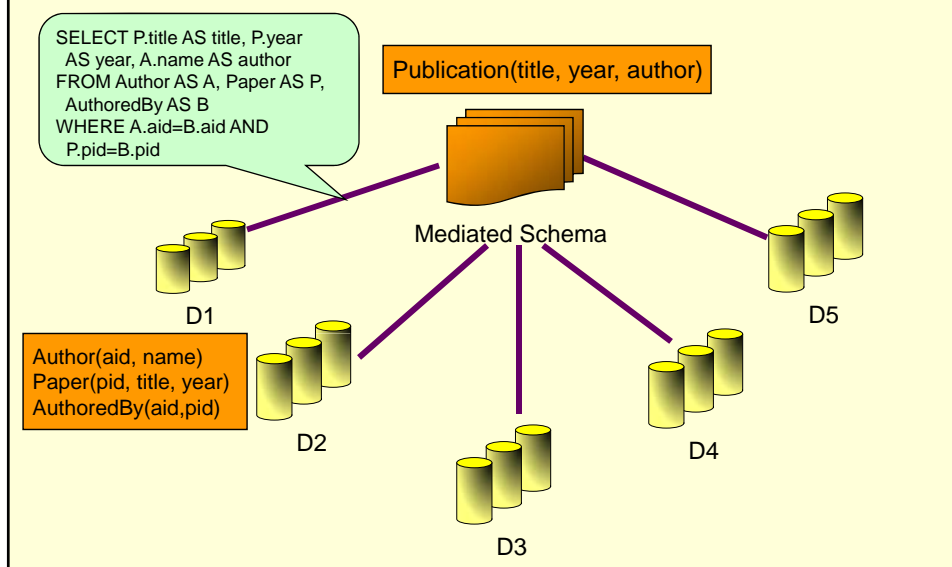
Outline

- ✓ Introduction
- Dataspace principles through data integration:
 - Recap of data integration
 - Locating and understanding data sources
 - Creating mappings and mediated schemas
 - Pay-as-you-go: improving with time
 - Query processing for dataspace
- Research challenges on specific dataspace:
 - The Web, Science, the desktop

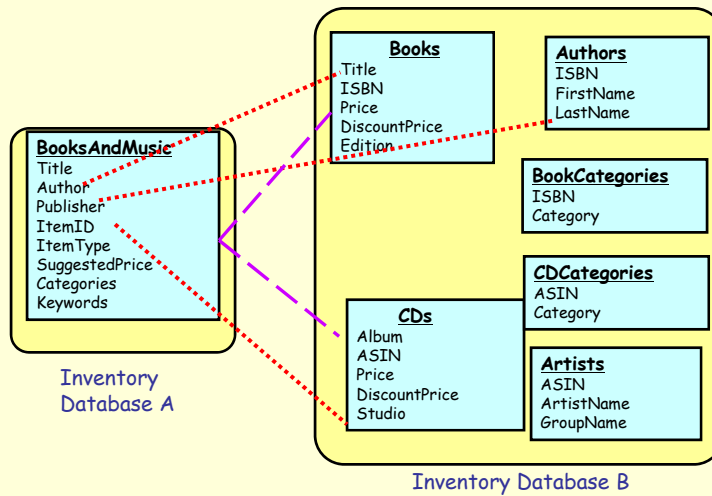
Recap of Data Integration

- Locating, understanding data sources
- Creating a mediated schema
- Creating schema mappings
- Querying and query processing

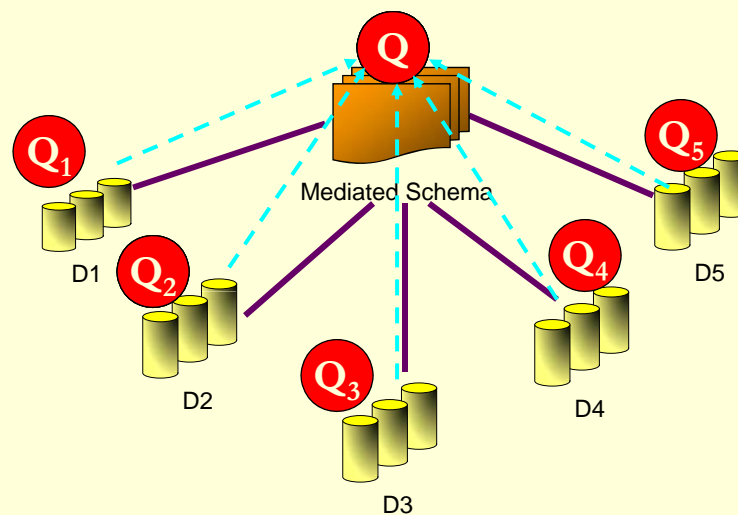
Traditional Data Integration Systems



Semantic Mappings



Querying on Traditional Data Integration Systems



Differences in Dataspaces

[to be covered next]

- Locating and understanding data sources is a challenge
- Semantic mappings will be created automatically:
 - May be approximate
- Create mediated schemas automatically
- Improve semantic mappings over time
- Different query mechanisms
 - Approximate and partial answers
 - User feedback at query stage

The Real Point

- Business data management: done.
- Now: data management = enabling collaboration.
- Collaboration = sharing + combining data, typically in ad-hoc manner.
- Ad-hoc sharing and combining = dataspace

Finding and Figuring Out Data Sources

Several activities here:

- Identification
- Familiarization
- Characterization
- Customization
- Selection

Steps are interwoven

Source Identification

- Deep-Web Search (e.g., CompletePlanet)
- Brute force?
 - Start with everything you can find (don't worry about relevance)
 - Postpone decisions on inclusion, e.g., until query time
- No generic, automated solution yet
 - How do you find relevant data that is “laying around”, e.g., in a spreadsheet on someone's laptop?

Other People's Data

Why is other people's data hard to understand?

- They aren't around to explain it
- What they wrote down about it wasn't quite true or complete in the first place
- The kinds of data in a source have expanded beyond the original intent:
schema drift

How do you understand the data?

25

Example Dataspace: RxSafe

Consolidated medication list for rural elders

Points in lifetime of a prescription

- Order (clinic, hospital)
- Dispensing (pharmacy)
- Approval (insurer)
- Administration (rehabilitation facility)

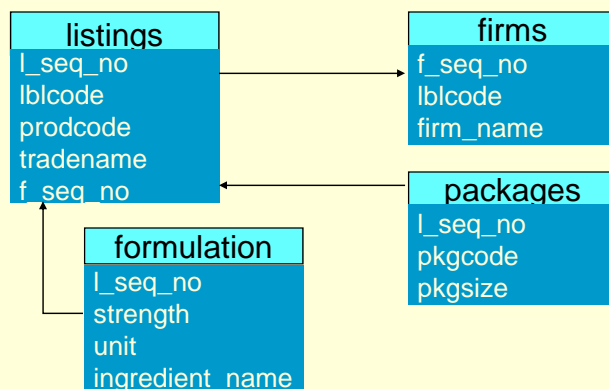
Relevant Standards

NDCCD, RxNorm, NDF-RT

Klaus Dittrich Symposium 2008

26

NDCD: National Drug Code Directory



Codes for drug packages

Sample NDC: 62584-023-00

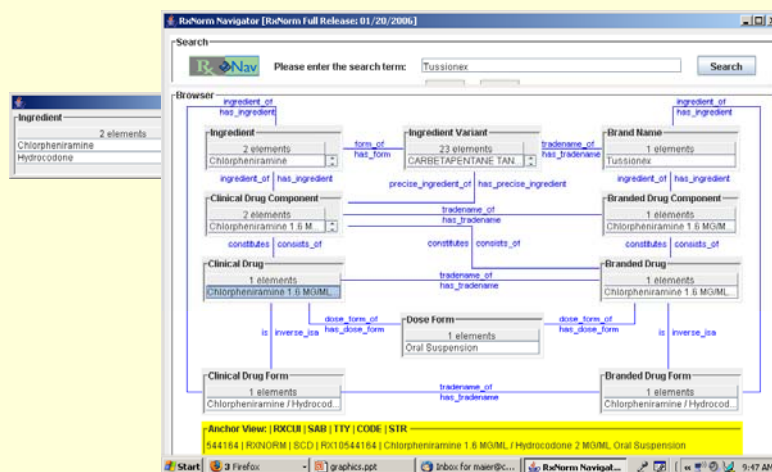
l_seq_no	lblcode	prodcod	tradename	f_seq_no
172062	62584	023	Vicodin tab	59064

l_seq_no	strength	unit	ingredient_name
172062	5	MG	Hydrocodone
172062	500	MG	Acetaminophen

l_seq_no	pkgcode	pkgsize
172062	00	100

f_seq_no	lblcode	firmname
59064	62584	Amerisource

RxNorm: Drug Nomenclature



RxNav from National Library of Medicine

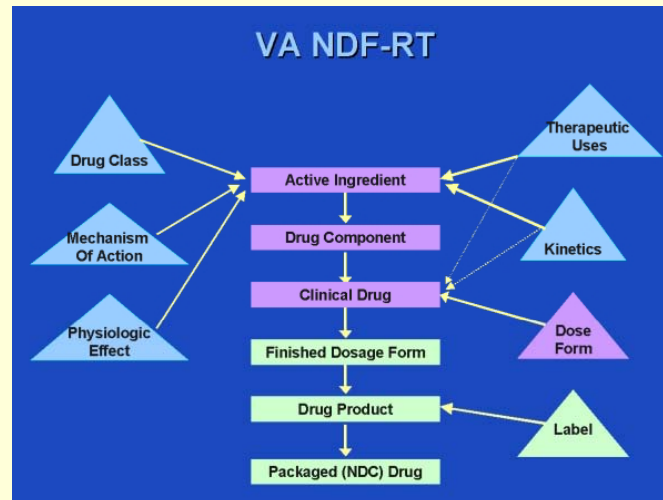
NDF-RT

National Drug File – Reference Terminology

From Veterans Affairs

- Drug class
- Chemical class
- Effects and actions

NDF-RT (Light Blue)



Goals of Understanding Based on Intended Purpose

- Grouping similar medications
- Connecting possible incarnations of same prescription
Generic – Brand Name
- Combining medication information for a given patient
Must be ***error preserving***

What You Know is Wrong

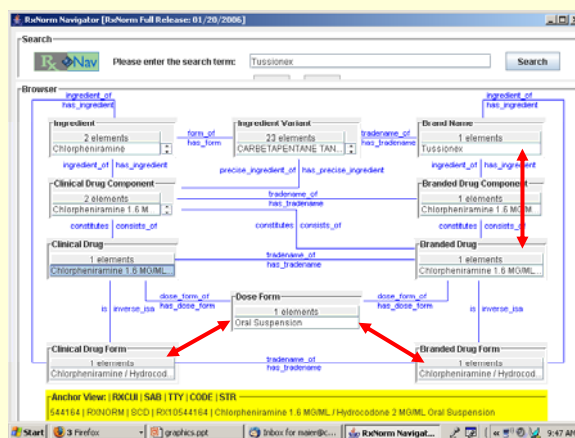
NDC and RxNorm talking about same things

- NDC tradenames: 18913
- RxNorm brand names: 7600
- Strings in common: 418

All RxNorm relationships have explicit inverses

33

What You Know is Incomplete



Doesn't mention atoms, attributes

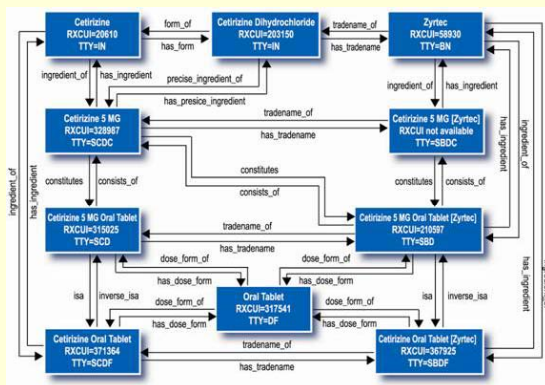
Doesn't include SY, ET, OCD, OBD

Source: National
Library of Medicine

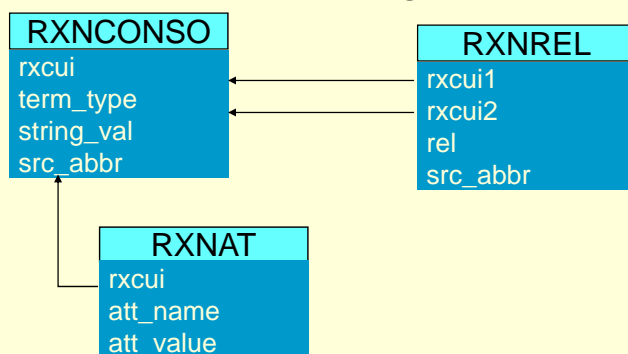
You Might Misunderstand What You're Told

RxNorm diagram is for instances

Multi-ingredient drug case not covered



The Structure Isn't What You Thought



RxNorm uses UMLS, not domain-specific

More complex than this – can have several *atoms* in each concept

How Do You Make Progress?

- Tools with minimal assumptions
that can help familiarize yourself with the data
- Ability to check hypothesis
to characterize what's true about the data (or almost true)
- Means to customize data to intended task
once you understand something about it

What I Want: Dataspace Charting Toolkit

Help with *Familiarization, Profiling, Enhancement*

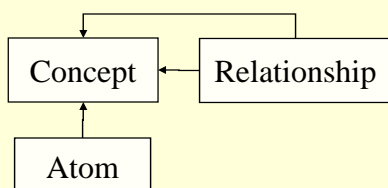
- Inspector for generic models
- Dataspace profiler
 - Assumption tracker and checker
 - Structure discovery techniques
- Customization to task based on discovered characteristics

Quarry: Scalable, Schema-less Browsing

Data Model [Howe, Rayner+ IIMAS 08]

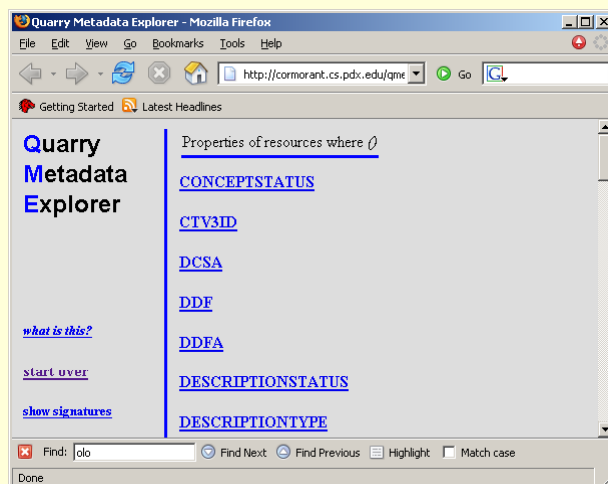
- resource, property, value
(subject, predicate, object) if you prefer
- no intrinsic distinction between literal values and resource values
can't necessarily tell the difference initially
- no explicit types or classes

Example: RxNorm

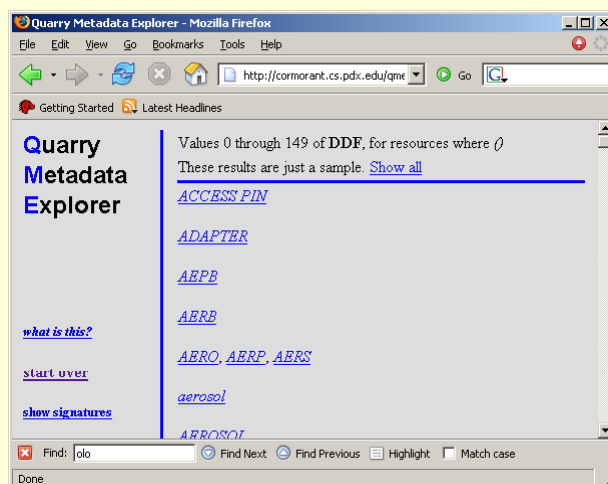


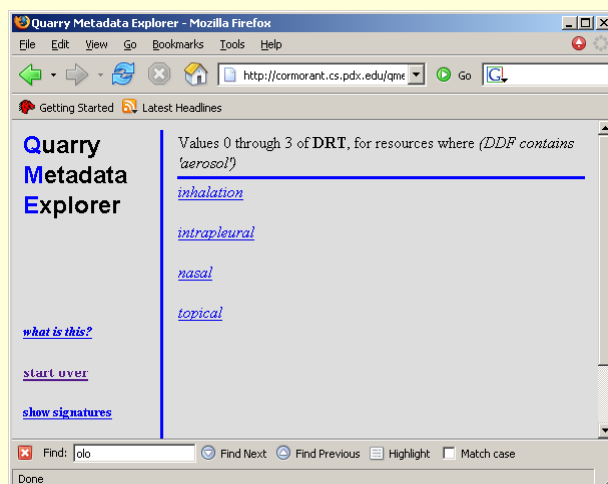
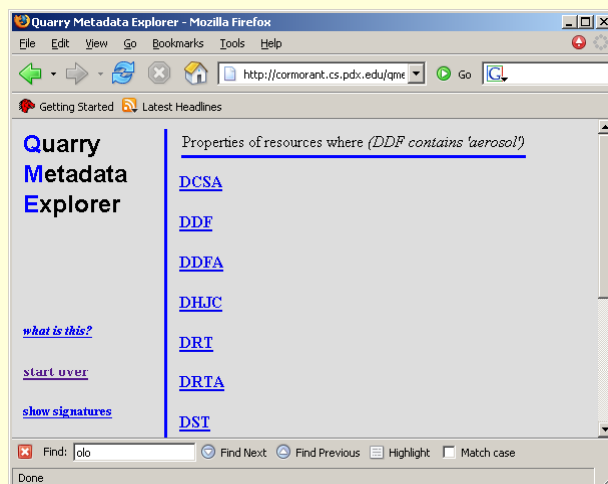
userkey	prop	value
10001	NDC	1
10001	ORIG_CODE	123
10001	ingredient_of	10004
10001	type	DC

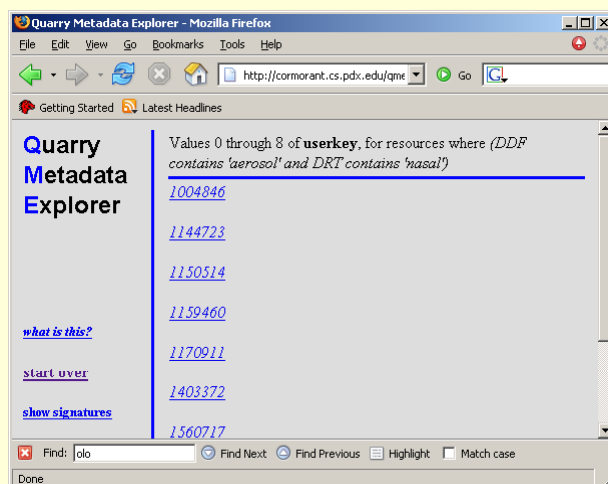
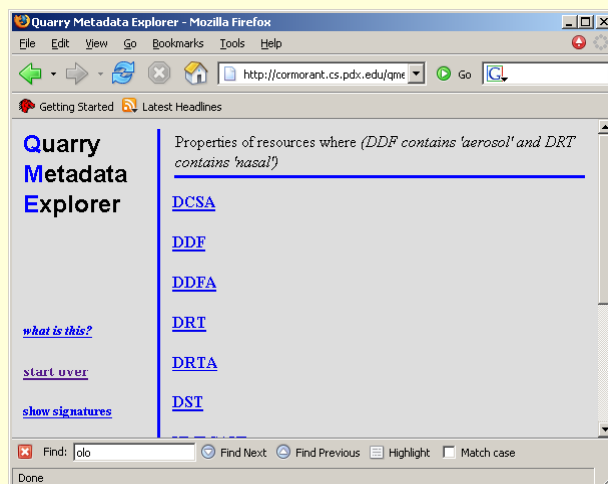
up to 23M triples describing 0.6M concepts and atoms

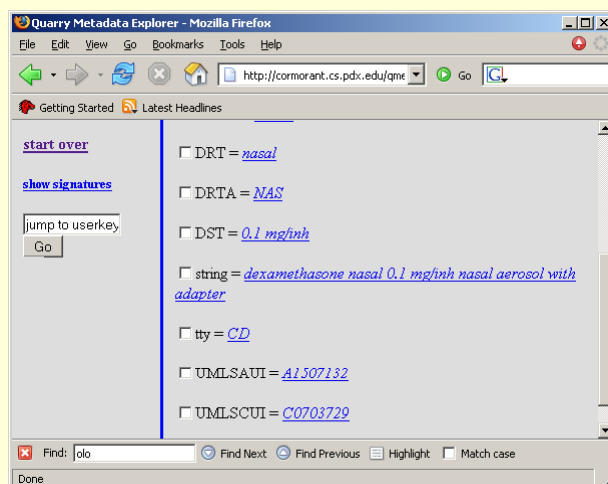
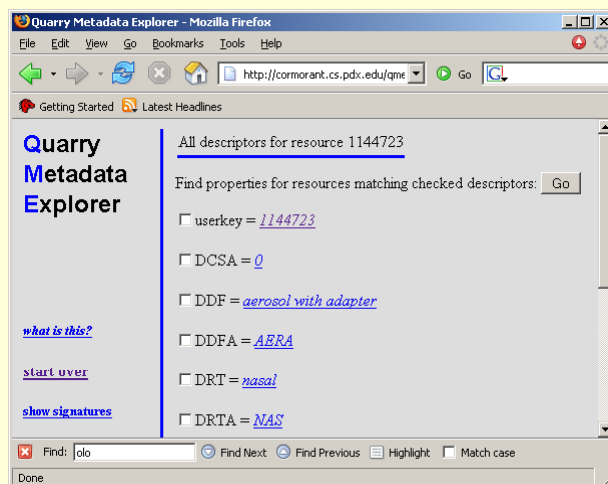


SKIP









Quarry API

.../2004/2004-001/.../anim-tem_estuary_bottom.gif

```

aggregate = bottom
animation = isotem
day = 001
directory = images
plottype = isotem
region = estuary
runid = 2004-001
year = 2004

```

Describe(key)

Properties(runid=2004-001)

:
.../2004/2004-001/.../amp_plume_2d.gif

```

day = 001
directory = images
plottype = 2d
region = plume
runid = 2004-001
year = 2004

```

**Values(runid=2004-001,
"plottype")**

Behind the Scenes

Signatures

- resources possessing the same properties clustered together
- Posit that $|\text{Signatures}| \ll |\text{Resources}|$
- Queries evaluated over Signature Extents

Dataspace Profiling

Commercial profilers: DataFlux, Infogix (ACR), KnowledgeDriver

- prep for cleaning, migration
- generally relational model, by table

Potter's Wheel [Raman, Hellerstein VLDB 2001]

- learning column transformation
- unfolding – data value to column labels

Cross-Source Profiles

Bellman [Dasu, Johnson+ SIGMOD 02]

- Find joinable columns (1-N, M-N)
- Is one field the composition of others?
- Part of T joins with T1, part with T2

Make the point that database schemas
“devolve” with time as business
processes change

Dataspace Profiling: NDC Examples

- `l_seq_no` is key of `listings` – yes
- `lblcode, prodcod` key of `listings` – no
45,953 45,972 (19)
- `firm_name` → `lblcode` – no
2931 2952 (21)
- each product `listing` should have >0
`packages` and >0 `ingredients`
44,972 (1180) 45,180 (792)

Courtesy of Nick Rayner

Checking Across Sources NDC vs. RxNorm

Ingredients

- 2794 ingredient names in NDCD
- 5145 ingredients in RxNorm
- 1570 equal strings

What to Do with Flawed Assumptions?

- Track exceptions
- Refine assumption
`firm_name, location` → `lblcode`
- Refine knowledge of world
RxNorm has *ingredient variants* (which have the same type as ingredients)
- Want to track assumptions as they evolve, results of checks

Customization: InfoSonde

Support customizations appropriate to discovered data characteristics

[Howe, Rayner+ IIMAS 2008]

Three-part modules

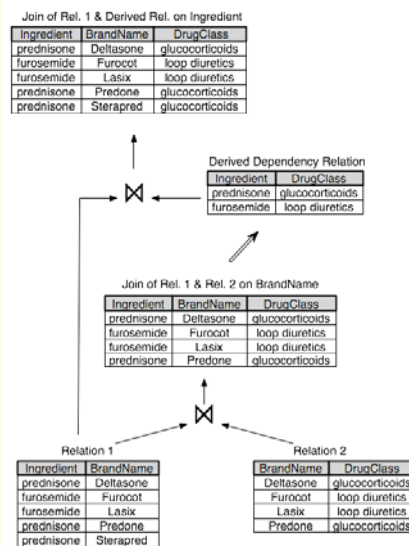
- Probe: Check or discover properties
- Switch: Present applicable customizations
- Check: Test that chosen switch is still valid

Functional Dependency Module

- Probe: Test for FD
- Switch:
 - FD holds: add constraint, decompose
 - FD fails: partition, repair
- Check:
 - Example – if using decompose, check that FD still holds

57

Linkage Extension Module



Want to extend a linkage based on a discovered functional relationship

Probe: Join satisfies FD here, ING → DC

Switch: Materialize functional relationship
can be used to extend original relation with DC

Check: Test that FD still holds

Related Work: Characterization and Enhancement

Structure Discovery [Andritsos, Miller+ SIGMOD 2005]

- Clustering of columns, rows
- Ranking FDs by corresponding redundancy

59

Other Structural Work

- Wide-table workbench [Chu, Baid+ VLDB 07]
 - Start with unstructured records & keyword search
 - **Extract**, **Integrate**, **Cluster** to get more structure
- LearnPADs [Fisher, Walker+ POPL 08]
 - Targets *ad hoc data*: semi-structured data with no toolset (e.g., log file)
 - Induces PADS grammar from examples, use to instantiate generic tools (parser, accumulator, formatter, translator)

Source Selection

- Obviously, familiarization and characterization can help in evaluating a single source
- But you may want to judge a combination of sources

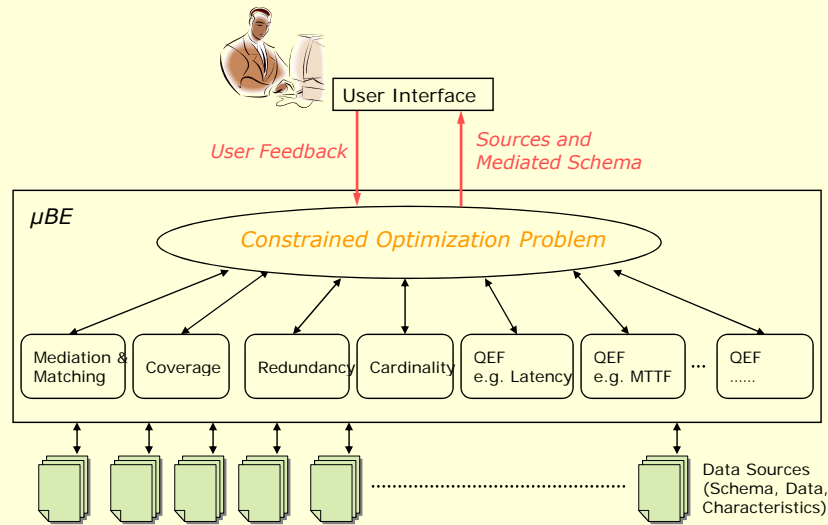
μ BE: User Guided Source Selection

Matching by Example

[Aboulnaga, El Gebaly ICDE 2007]

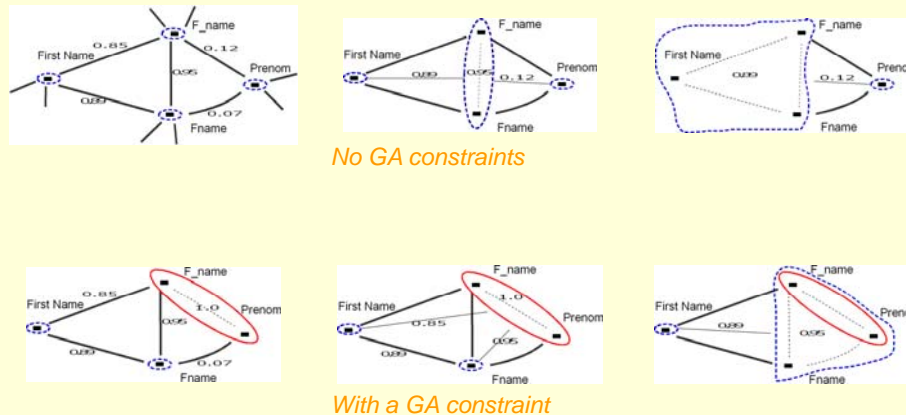
- Want to select set of sources based on matching quality, cardinality, coverage, redundancy and possibly other measures (latency)
- Mediated schema consists of global attributes = sets of source attributes
- User can constrain sources and global attributes

μ BE Architecture



Courtesy Ashraf Aboulnaga & Kareem El Gebaly

μ BE Matching Example



Courtesy Ashraf Aboulnaga & Kareem El Gebaly