

IR & DB:

Toward Controversy and
Philosophy

Joe Hellerstein

UC Berkeley &
Intel Research

Controversial Statements

Controversial Statements

- Surajit was right

Controversial Statements

- Surajit was right
- I am an expert on this stuff

Controversial Statements

- Surajit was right
- I am an expert on this stuff
 - Led integration of Altavista and Cohera federated DB
 - Keyword search, Q-grams, auto-taxonomization, visual data cleaning & integration, etc.
 - With consulting from Marti Hearst (IR meets DB!)
 - IH & S, SIGMOD '01 industrial track I

Controversial Statements II

- This is mostly busyness, some re-search

Controversial Statements II

- This is mostly busyness, some re-search
 - Identify BKMs
 - Perform SI
 - Micro\$oft will work on TCO

Controversial Statements II

- This is mostly busyness, some re-search
 - Identify BKMs
 - Perform SI
 - Micro\$oft will work on TCO
- N\$F?

- SO: Let's "Uplevel the Dialogue"

Let's Get Philosophical

- We have a narrow definition of structure
- Structure is not about Relational Databases vs. Text Databases vs. "semi-structured"
 - In fact, even text is very structured!

Let's Get Philosophical

- We have a narrow definition of structure
- Structure is not about Relational Databases vs. Text Databases vs. "semi-structured"
 - In fact, even text is very structured!
- Human discourse based on "deep structure"
 - Ferdinand De Saussure, the father of Structural Linguistics
 - Extended by Chomsky
 - Also speak of "relational" nature of language

So what's the philosophical diff?

- Twofold
 - The SOURCE of the structure
 - The USE of the structure

Relational Databases

- ENGINEERED STRUCTURE
 - The design of structure is at the heart of the discussion
 - Codd's lessons: simple design for robust evolution
- LANGUAGE AS PROGRAMMATIC INTERFACE
 - Semantically strict queries provide predictable results
 - Suitable for computer interpretation
 - Programmer can reason about invariants
 - Hence good for embedding into application code
 - Relatively few users interface directly to a Database!

Information Retrieval

- "FOUND" STRUCTURE
 - Take a pile of information, and extract structure therefrom
- HUMAN-CENTRIC EXPLOITATION
 - Rough understanding of query intent
 - Interface more important than query language
 - user can browse/filter/interpret some results
 - Requires a human in the loop
 - Relatively few programs embed IR techniques invisibly

The Synergy

- This is not about semi-structured data!!!

It's not that semi-structured is bad...

It's not that semi-structured is bad...

It's just that semi-
structured is not semi-
structured

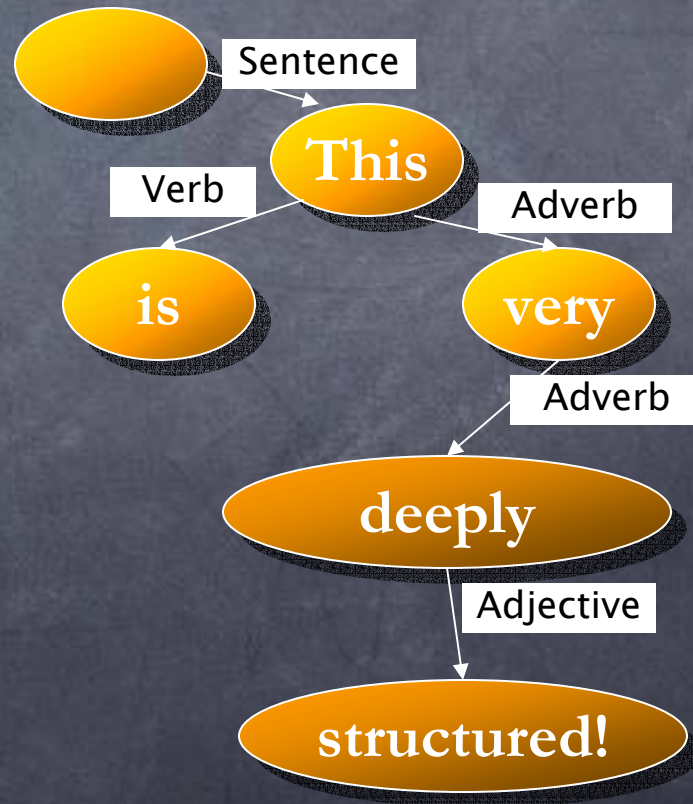
It's not that semi-structured is bad...

It's just that semi-structured is not semi-structured



It's not that semi-structured is bad...

It's just that semi-structured is not semi-structured



The Synergy

- DB folk working on FOUND STRUCTURE
 - E.g. text queries on "structured" data
 - Information extraction
 - User-centric data mining
 - Schema corpí
- DB folk working on HUMAN INTERFACES for handling uncertainty
 - E.g. keyword search of structure data
 - E.g. online aggregation, visual data cleaning cube navigation, other user-centric mining tasks



Logout Help About

Catalogs Users

- [-] Master Catalogs
 - [-] PC Hardware
 - [-] Categories
 - [-] Customers
 - [-] Prices
 - [-] Products
 - [-] Suppliers
 - [-] ACEtechnologies
 - [-] Loadable Price Tables
 - [-] Loadable Product Table
 - scb
 - [-] Supplier's Prices
 - [-] Supplier's Products
 - [-] HP
 - [-] IBM
 - [-] extra
 - [-] electrical
 - [-] Catalogs
 - [-] PC Hardware-ACEtechnologies-Su
 - [-] PC Hardware-Best Buy-Customer
 - [-] PC Hardware-Circuit City-Customer
 - [-] PC Hardware-HP-Supplier Catalog
 - [-] PC Hardware-IBM-Supplier Catalog
 - [-] PC Hardware-extra-Supplier Catalog

ADD TRANSFORM MAP LOAD CLASSIFY REPAIR

Transform Rules Column View Tools

Showing 1 - 116 of 116

		availability	manufacturer	category_code	category_name	
1	Add	Ships in 48 hours	Generic	awb120000002	Harddisk Controllers	1
2	Calc	Ships in 48 hours	Adaptec	awb120000002	Harddisk Controllers	2
3	Copy	Ships in 48 hours	Promise Technology	awb120000002	Harddisk Controllers	3
4	Convert	Ships in 48 hours	Microsoft	awb120000003	Keyboard	4
5	Divide	Ships in 48 hours	Logitech	awb120000003	Keyboard	5
6	Drop	Ships in 48 hours	CA-Datacom	awb120000004	Lanpack	6
7	Format	Ships in 48 hours	Powerplay	awb120000004	Lanpack	7
8	Merge	Ships in 48 hours	Mitsumi	awb120000005	Floppy Drives	8
9	Split	Ships in 48 hours	Sony	awb120000005	Floppy Drives	9
10	SQL	Ships in 24 hours.	ALI	awb120000006	Graphics	10
11	UDF	Ships in 24 hours.	ATI	awb120000006	Graphics	11
12		Ships in 24 hours.	Corsair	awb120000007	Memory	12
13	137.57	Ships in 24 hours.	Micron	awb120000007	Memory	13
14	73.57	Ships in 24 hours.	Panasonic	awb120000008	Mobile Power	14
15	117.57	Ships in 24 hours.	Compaq	awb120000008	Mobile Power	15
16	128.57	Ships in 24 hours.	Abit	awb120000009	Motherboard	16
17	803.57	Ships in 24 hours.	AsusTek	awb120000009	Motherboard	17
18	26.57	Ships in 24 hours.	Turtle Beach	awb12000000a	Audio	18
19	133.57	Ships in 24 hours.	Ensoniq	awb12000000a	Audio	19
20	64.57	Ships in 7 days	Sony	awb12000000b	Battery	20
21	613.57	Ships in 7 days	Battery Biz	awb12000000b	Battery	21
22	103.57	Ships in 7 days	Intel	awb12000000c	Accessories	22
23	203.57	Ships in 7 days	Vaster	awb12000000c	Accessories	23
24	22.57	Ships in 7 days	Maxtor	awb12000000d	CDROM	24

Save Previous Page Next Page Help

The Synergy

- IR folk working on ENGINEERED STRUCTURE
 - E.g. document design
 - But the simpler, the better (Codd)!!
- IR folk working on embedded systems
 - Google SOAP interface clients (Googlisms)
 - ??
- Frankly, I'm out of my comfort zone

For More Rambling on this Topic

- See the WebDB keynote on my home page
 - <http://www.cs.berkeley.edu/~jmh>

The Inktomi Search Engine

Result Set = [DocId, Score, URL, Date, Size, Abstract]

