



Hasso
Plattner
Institut

IT Systems Engineering | Universität Potsdam

Dr. Crowdsorce

or: How I Learned to Stop Worrying and Love Web
Data

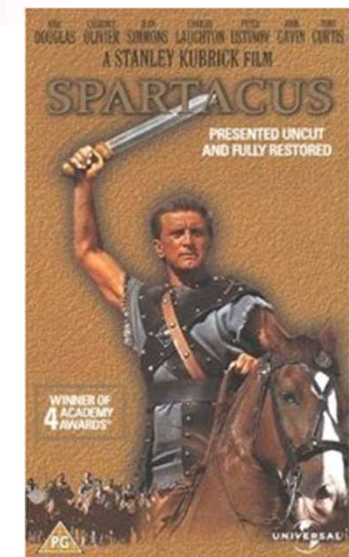
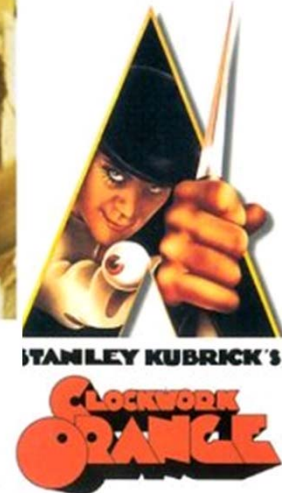
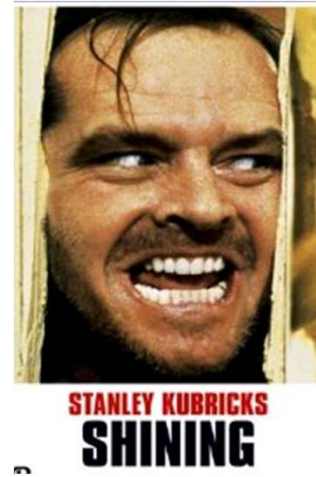
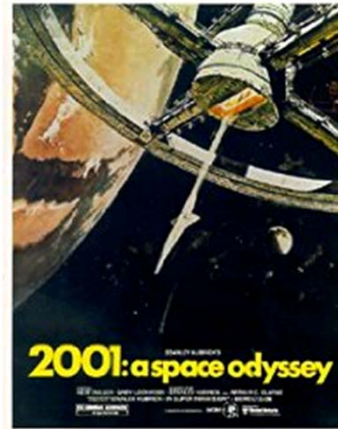
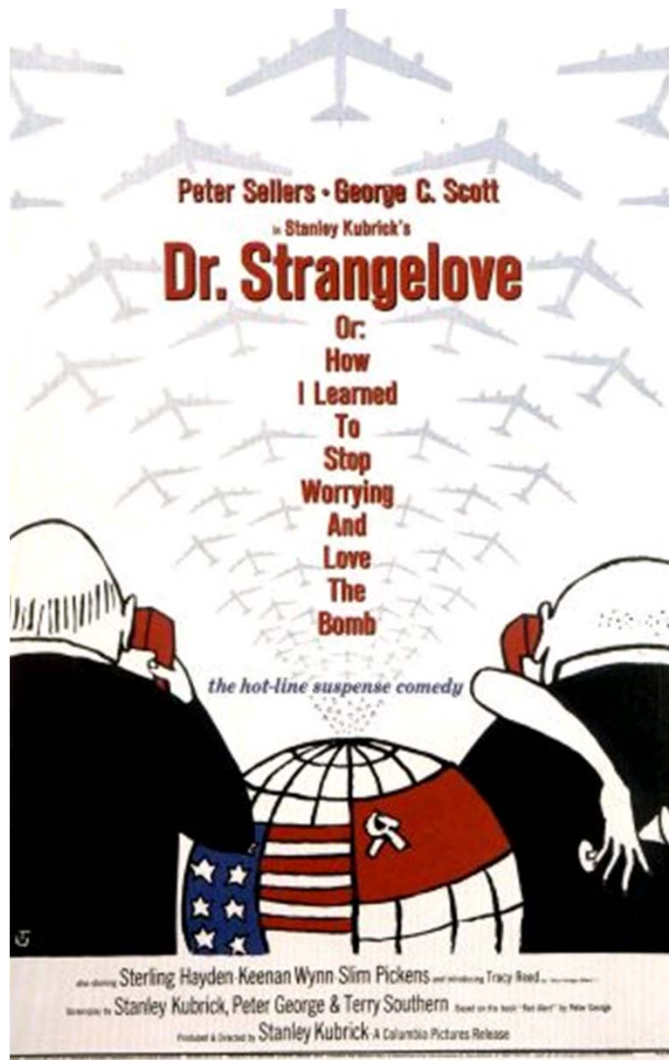
Invited Talk @ BEWEB 2011

25.3.2011

Felix Naumann

Dr. Strangelove – for the small fry...

2



Linked Data & Data Spaces – a database guy's PoV

5



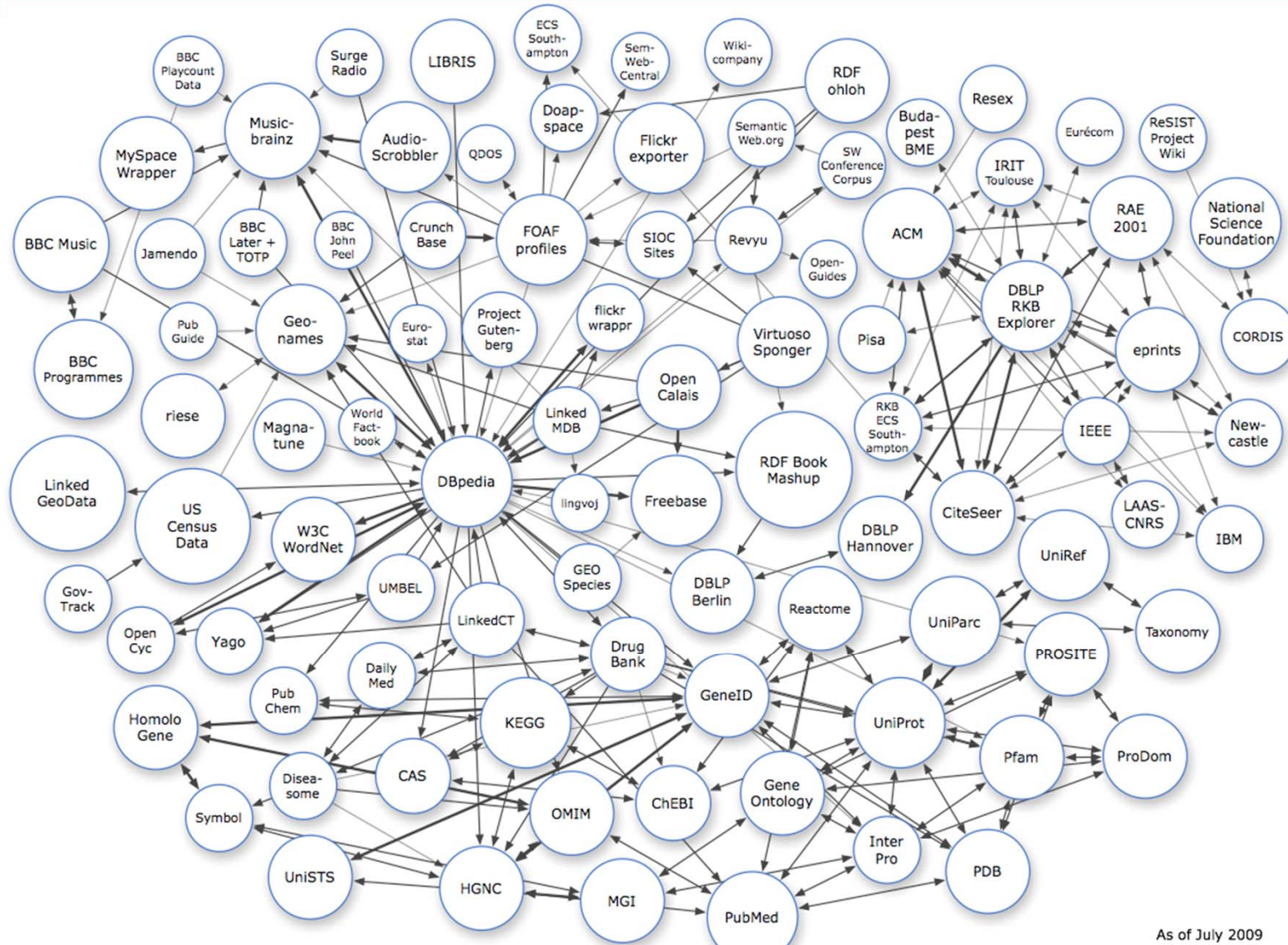
Linked data – 4 Principles, 7 Properties

6

1. Use **URIs as names** for things.
 2. Use **HTTP URIs** so that people can look up those names.
 3. When someone looks up a URI, **provide useful information**.
 4. Include **links to other URIs**, so that they can discover more things.
 - Many common things are represented in multiple data sets!
- The Good
 - Comes as triples
S: `http://.../Uppsala`
P: `location`
O: `http://.../Sweden`
 - Often user generated
 - Nice domains
 - Free
 - The Bad
 - Voluminous
 - Heterogeneous
 - The Ugly
 - Dirty, inconsistent, sparse

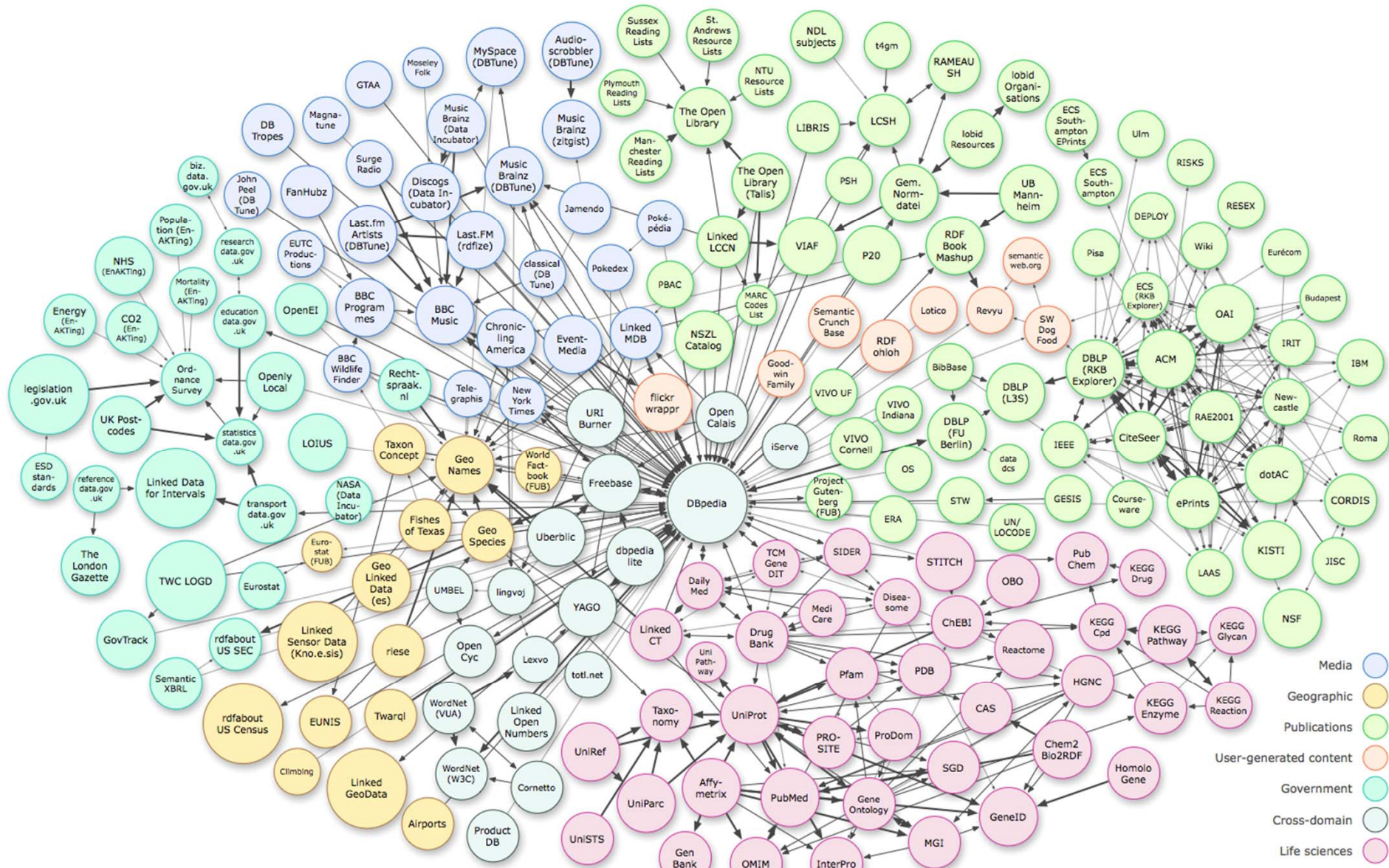
Linked Data Graph

7



As of July 2009

Linked Data Graph



DBpedia – Extraction

9

```

{{Infobox Non-profit
| Non-profit_name      = IEEE
| Non-profit_logo     = [[Image:IEEE logo.svg|200px]]
| Non-profit_type     = Professional Organization
| founded_date       = January 1, 1963
| founder            =
| location           =
| origins            = Merger of the American Institute of Electrical Engineers and
| key_people        = Mr. Pedro A. Ray, Current President
| area_served       = Worldwide
| focus             = Electrical, Electronics, and Information Technology [http://w
/visionmission.html]
| method            = Industry standards, Conferences, Publications
| revenue           = US$330 million
| endowment         =
| num_volunteers    =
| num_employees     =
| num_members       = 395,000+
| owner             =
| Non-profit_slogan  =
| homepage          = [http://www.ieee.org/ www.ieee.org]
| tax_exempt        =
| dissolved         =
| footnotes         =
}}

```

IEEE



Type	Professional Organization
Founded	January 1, 1963
Origins	Merger of the American Institute of Electrical Engineers and the Institute of Radio Engineers
Key people	Mr. Pedro A. Ray, Current President
Area served	Worldwide
Focus	Electrical, Electronics, and Information Technology [1] 
Method	Industry standards, Conferences, Publications
Revenue	US\$330 million
Members	395,000+
Website	www.ieee.org 

DBpedia statistics

10

1. Core Datasets

Dataset	en	de	fr	es	it	pl	nl	pt	sv	ja	ru	zh	fi	no
Titles (preview)	nt	nt	nt	nt	nt	nt	nt	nt	nt	nt	nt	nt	nt	nt
Short Abstracts (preview)	nt -	nt -	nt -	nt -	nt -	nt -	nt -	nt -	nt -	nt -	nt -	nt -	nt -	nt -
Extended Abstracts (preview)	nt -	nt -	nt -	nt -	nt -	nt -	nt -	nt -	nt -	nt -	nt -	nt -	nt -	nt -
Images (preview)	nt	nt	nt	nt	nt	nt	nt	nt	nt	nt	nt	nt	nt	nt
Links to Wikipedia Article (preview)	nt	nt	nt	nt	nt	nt	nt	nt	nt	nt	nt	nt	nt	nt
Articles Categories (preview)	nt	nt	nt	nt	nt	nt	nt	nt	nt	nt	nt	nt	nt	nt
External Links (preview)	nt	nt	nt	nt	nt	nt	nt	nt	nt	nt	nt	nt	nt	nt
Infoboxes (preview)	nt	nt	nt	nt	nt	nt	nt	nt	nt	nt	nt	nt	nt	nt
Properties (preview)	nt	nt	nt	nt	nt	nt	nt	nt	nt	nt	nt	nt	nt	nt
DBpedia Ontology (preview)	owl	owl	owl	owl	owl	owl	owl	owl	owl	owl	owl	owl	owl	owl
Ontology Infoboxes (preview)	nt	nt	nt	nt	nt	nt	nt	nt	nt	nt	nt	nt	nt	nt
Ontology Types (preview)	nt	nt	nt	nt	nt	nt	nt	nt	nt	nt	nt	nt	nt	nt
Homepages (preview)	nt	nt	nt	nt	nt	nt	nt	nt	nt	nt	nt	nt	nt	nt
Geographic Coordinates (preview)	nt	nt	nt	nt	nt	nt	nt	nt	nt	nt	nt	nt	nt	nt
Pagelinks (preview)	nt	nt	nt	nt	nt	nt	nt	nt	nt	nt	nt	nt	nt	nt
Persondata (preview)	nt	nt	nt	nt	nt	nt	nt	nt	nt	nt	nt	nt	nt	nt
Redirects (preview)	nt	nt	nt	nt	nt	nt	nt	nt	nt	nt	nt	nt	nt	nt
Disambiguation Links (preview)	nt	nt	nt	nt	nt	nt	nt	nt	nt	nt	nt	nt	nt	nt

672 million triples

286 million English

From 97 languages of Wikipedia

3.5 million things

364,000 persons

462,000 places

99,000 music albums

54,000 films

16,500 video games

<http://wiki.dbpedia.org/Datasets>

And more sources

11

- Government data
 - www.data.gov
 - data.gov.uk
 - ec.europa.eu/eurostat
- Finance / business data
- Scientific databases
 - www.uniprot.org
 - skyserver.sdss.org
- The Web
 - HTML tables and lists
 - General sources: DBpedia, freebase, ...
 - Domain-specific sources: IMDB, Gracenote, isbndb, ...
- ...



„Raw data now!“



Use cases

12

- General purpose integration: Create rich knowledge bases
 - Semantic Web
 - Improved search / question answering
 - Link creation and data enrichment
 - Cleansing: data correction and validation
- Domain specific integration
 - Creation of high quality data sets: Complete & accurate
 - Enhancement of organization-internal data
 - Create reference data sets
 - Mashups

KILLER APP?



Nineteen Eighty-Four

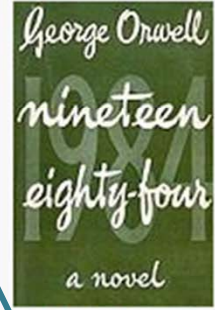
From Wikipedia, the free encyclopedia

This article is about the Orwell novel. For the year, see 1984. For other uses, see 1984 (disambiguation).

Nineteen Eighty-Four often abbreviated to **1984** is a classic **dystopian** novel by English author **George Orwell**. Published in **1949**, it is set in the **eponymous** year and focuses on a repressive, **totalitarian** regime. Orwell elaborates on how a massive **oligarchical** collectivist society such as the one described in *Nineteen Eighty-Four* would be able to repress any long-lived dissent. The story follows the life of one seemingly insignificant man, **Winston Smith**, a **civil servant** assigned the task of perpetuating the regime's **propaganda** by falsifying records and political literature so that it appears that the government is always correct in what it says. Smith grows disillusioned with his meager existence and so begins a **rebellion** against the system that leads to his arrest and torture.

The novel has become famous for its portrayal of pervasive government **surveillance** and control, and government's increasing encroachment on the rights of the individual. Since its publication, many of its terms and concepts, such as "thoughtcrime", and "Newspeak" have entered the popular lexicon. The term itself has come to refer to anything reminiscent of the novel. It is generally considered to be George Orwell's **magnum opus**.

Nineteen Eighty-Four (1984)



British first edition cover

Author	George Orwell
Country	United Kingdom
Language	English
Genre(s)	Dystopian , Political novel, Social science fiction
Publisher	Secker and Warburg (London)
Publication date	8 June 1949
Media type	Print (Hardcover & Paperback) & e-book, audio-CD
Pages	326 pp (Paperback edition)
ISBN	9780413024061



Contents [hide]

- 1 History
 - 1.1 Title
 - 1.2 Popular misconceptions
 - 1.3 Copyright status
- 2 Story
 - 2.1 Background
 - 2.2 Plot
- 3 Orwell's influences
- 4 Characters
 - 4.1 Major characters
 - 4.2 Minor characters
- 5 Fictional world
 - 5.1 Ingsoc (English Socialism)
 - 5.2 Ministries of Oceania
 - 5.3 Doublethink
 - 5.4 Political geography

navigation

- Main page
- Contents
- Featured content
- Current events
- Random article

search

Go Search

interaction

- About Wikipedia
- Community portal
- Recent changes
- Contact Wikipedia
- Donate to Wikipedia
- Help

toolbox

- What links here
- Related changes
- Upload file
- Special pages
- Printable version
- Permanent link
- Cite this page

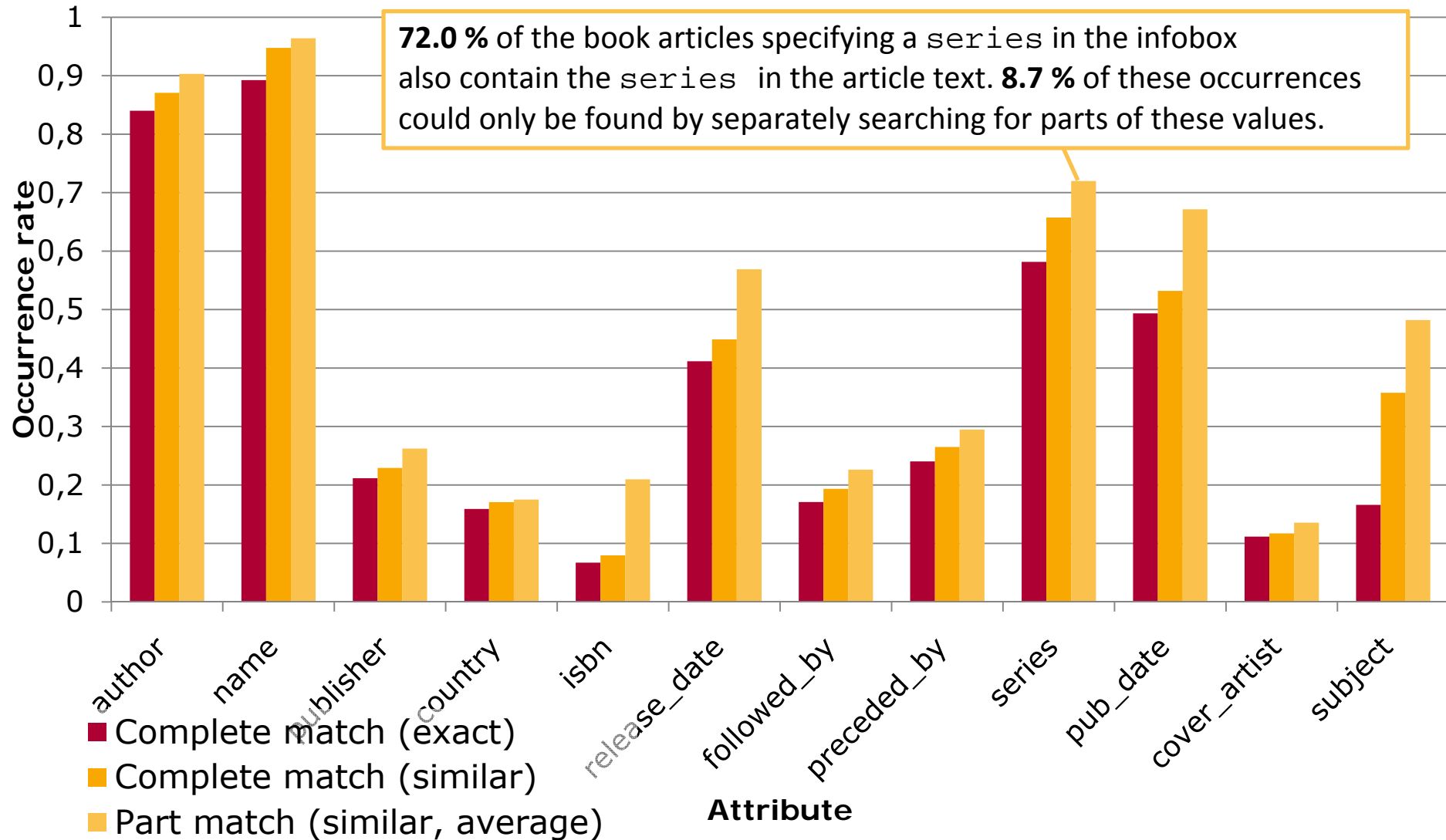
languages

- Anglo-Saxon
- العربية
- Беларуская
- Bosanski
- Български
- Català
- Česky

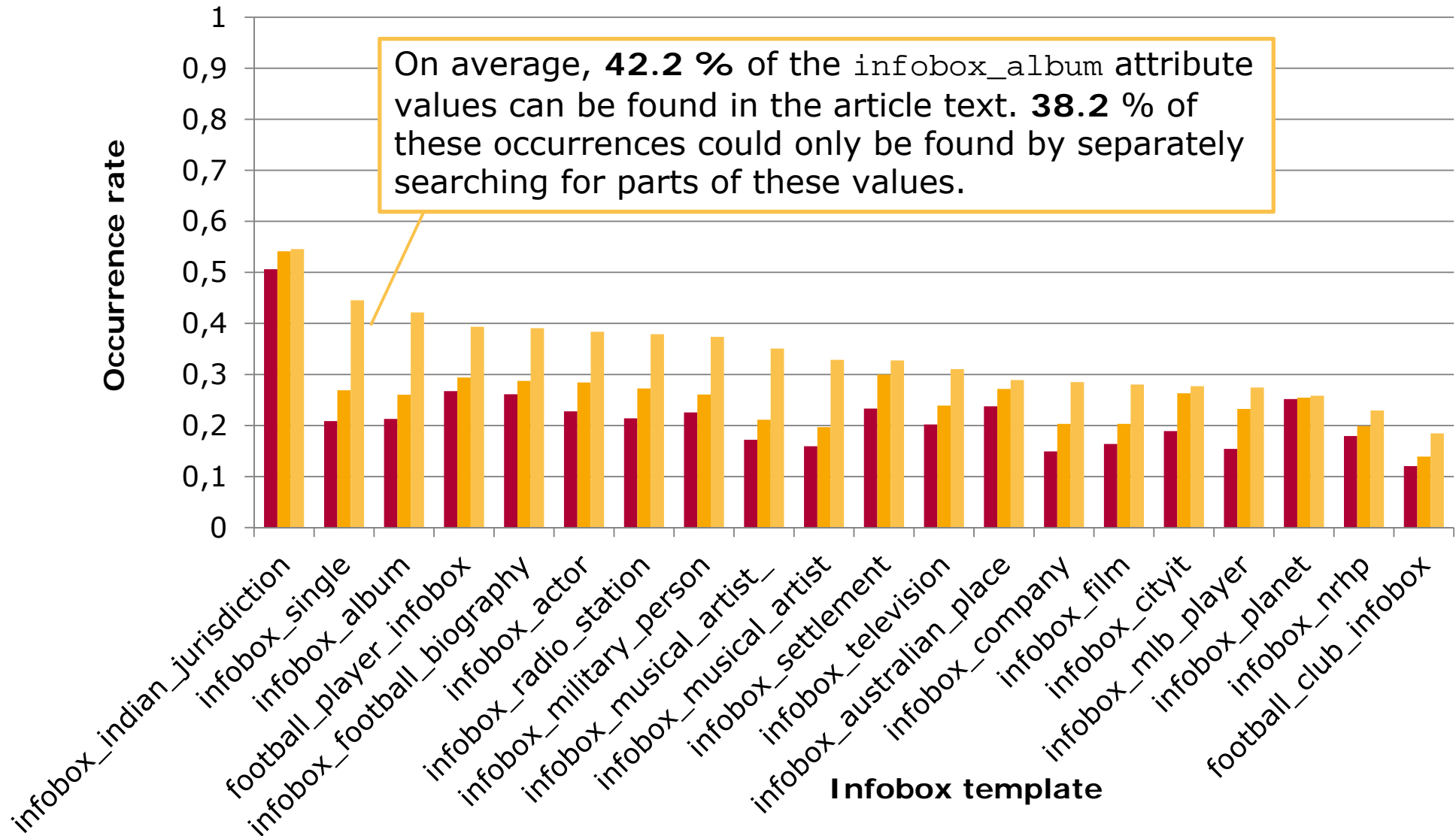
Master thesis by Dustin Lange
Now PhD student at HPI
Topic: similarity search



Occurrence of values in article text: 12 most frequent attributes in infobox_book

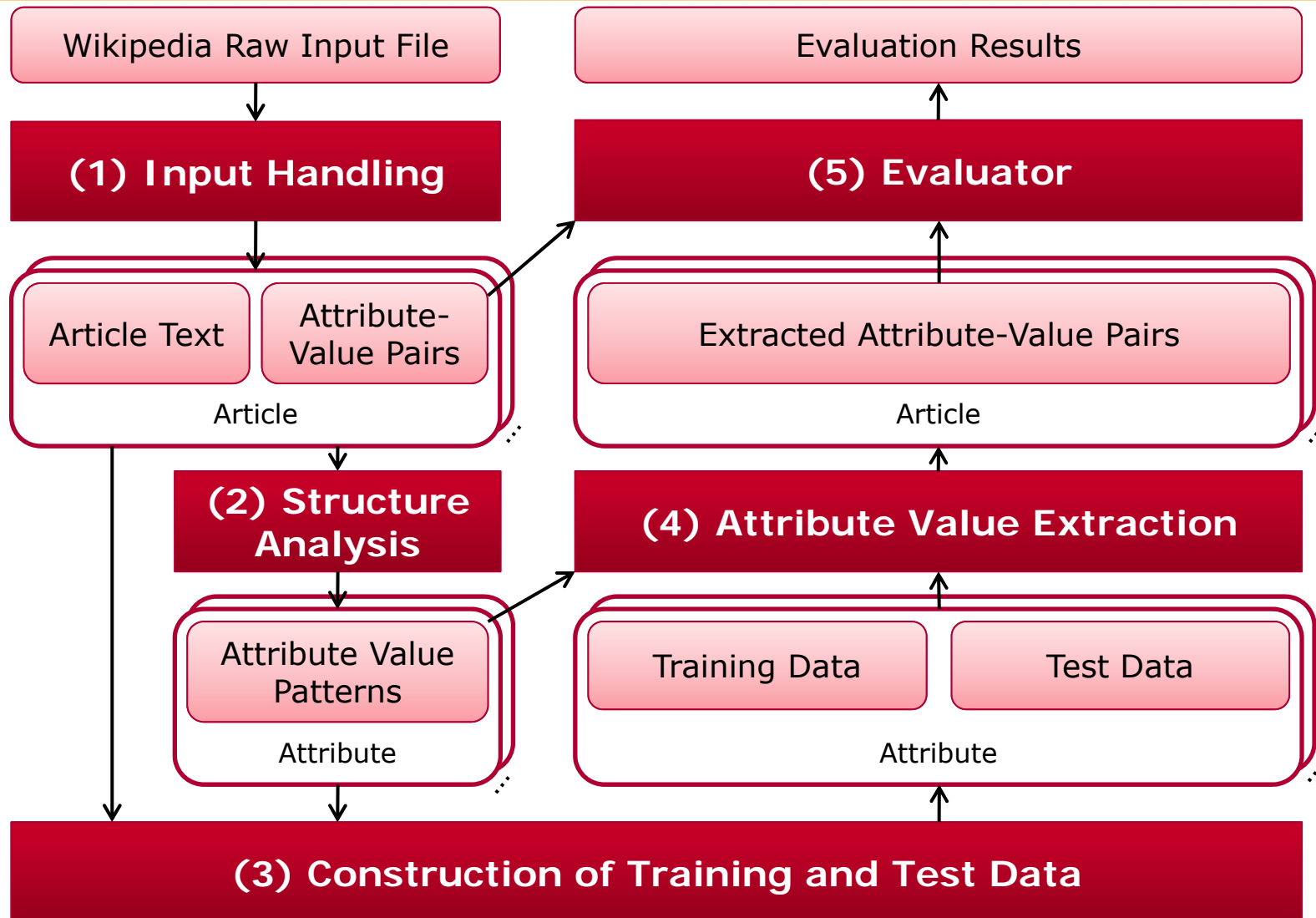


20 most frequent templates



Architecture of iPopulator

17



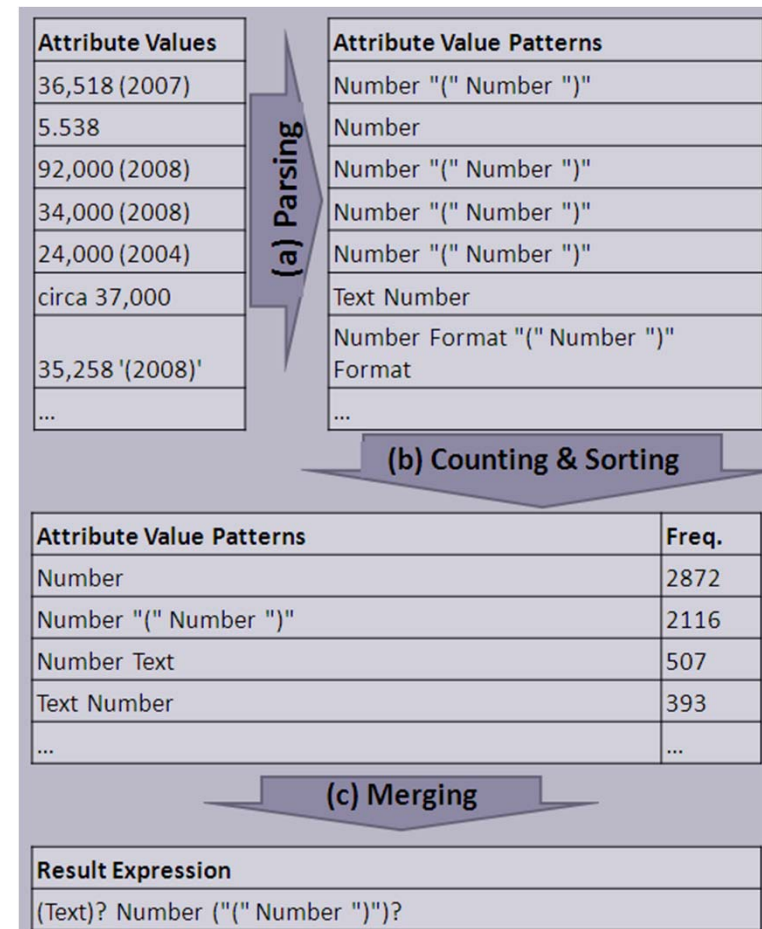
Structure Analysis

18

- Values of an attribute often share similar structure.
 - Extract value parts
 - Constructing homogeneous values from parts

- Determine common structure for each infobox template attribute

- Example: number_of_employees from infobox_company



Training Data and Extraction

19

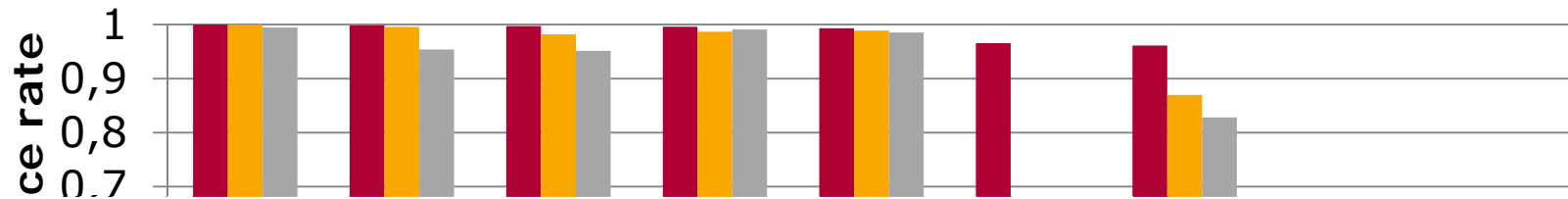
- Exploit existing infoboxes as training data
- Mark occurrences of infobox attribute values as training examples
 - Similarity measure to label fuzzy occurrences
- Automatic extraction method learns to recognize these occurrences by analyzing token (word-level) features
- Create extractors for thousands of infobox template attributes
- Extract parts of attribute values from different article text positions
- Arrange extracted value parts

Example: IBM employed 399,409 people in 2009.

Structure	Number	“(“	Number	”)”
Extracted value parts	399,409		2009	
Result value	399,409 (2009)			

Evaluation: infobox_planet

20

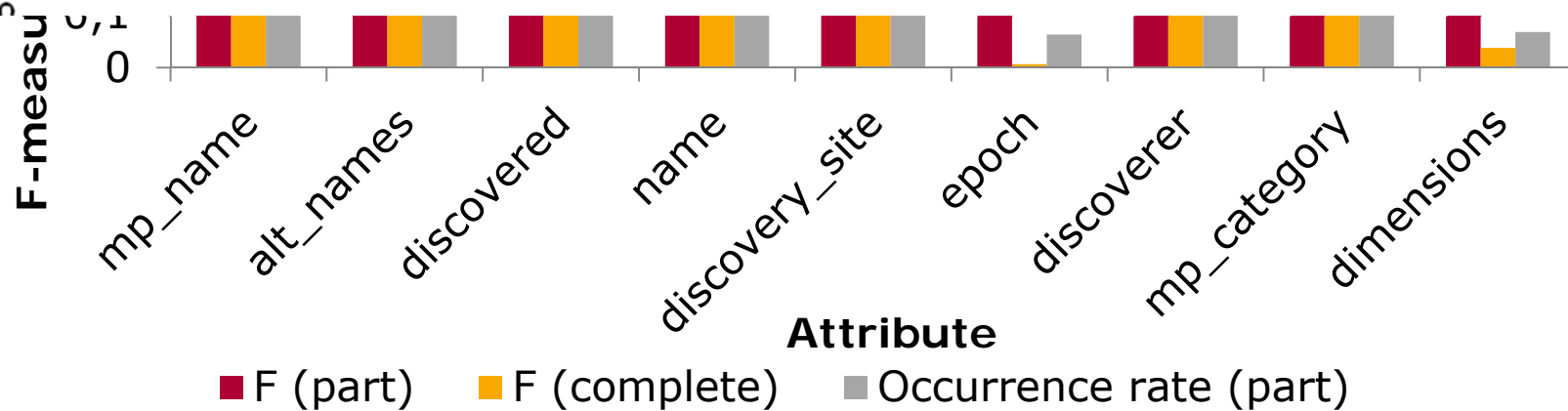


22032 Mikekoop

From Wikipedia, the free encyclopedia

22032 Mikekoop (provisional designation: **1999 XB₁₅₄**) is a [main-belt minor planet](#). It was discovered through the [Lowell Observatory Near-Earth-Object Search](#) at the [Anderson Mesa Station](#) in [Coconino County, Arizona](#), on December 9, 1999. It is named after Michael Walter Koop, an American electric engineer and amateur astronomer.

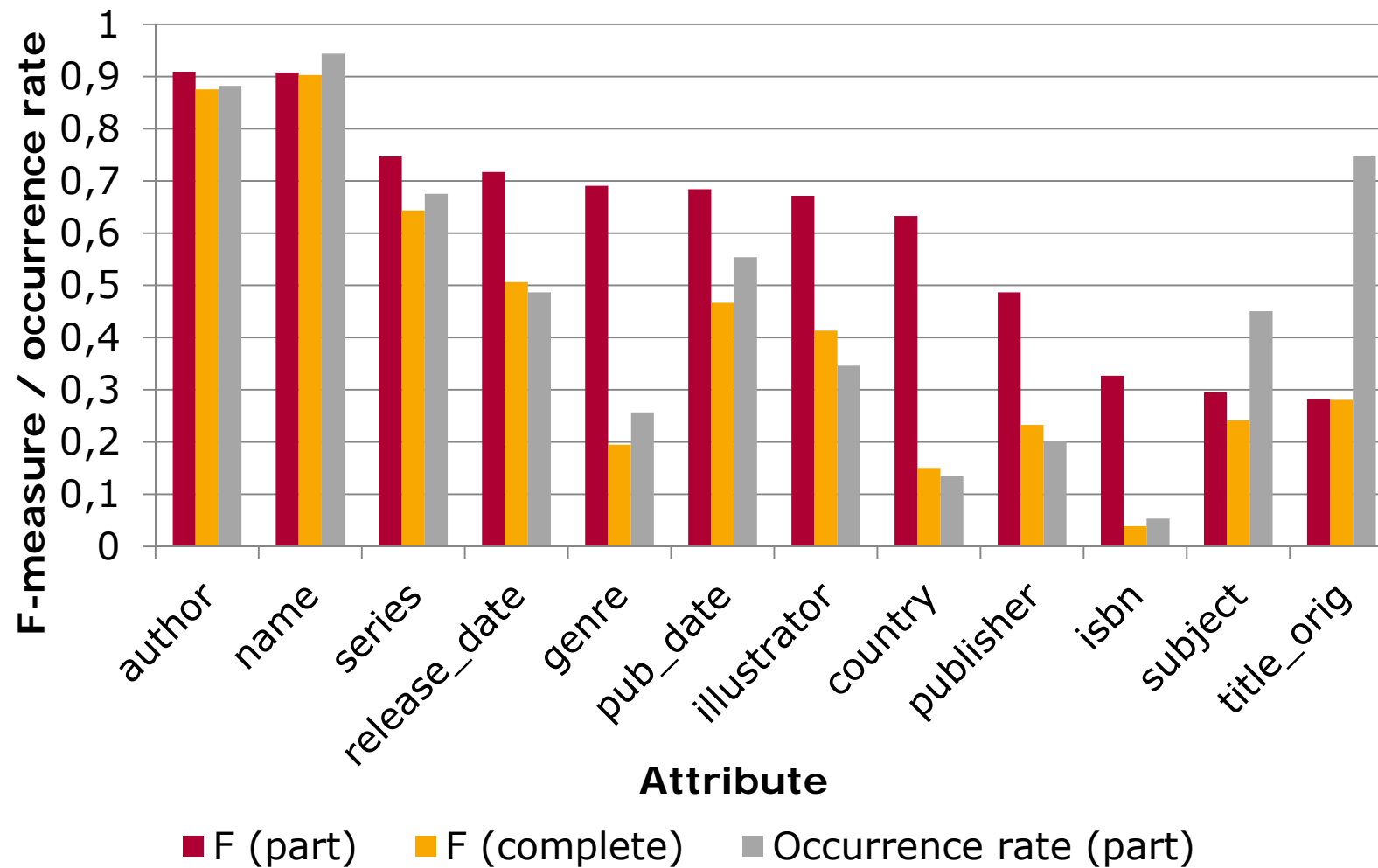
See also



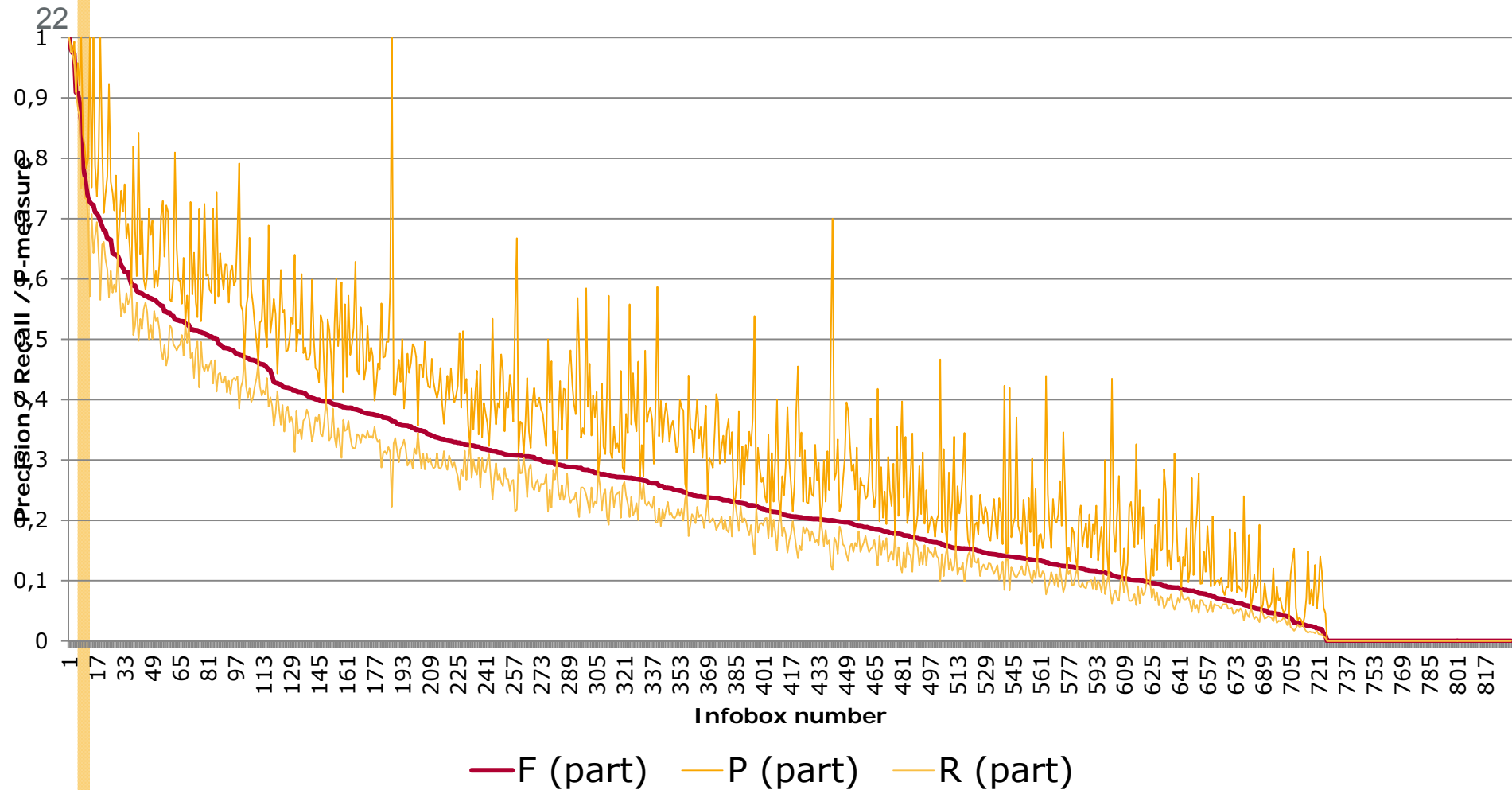
[edit]

Evaluation: infobox_book

21

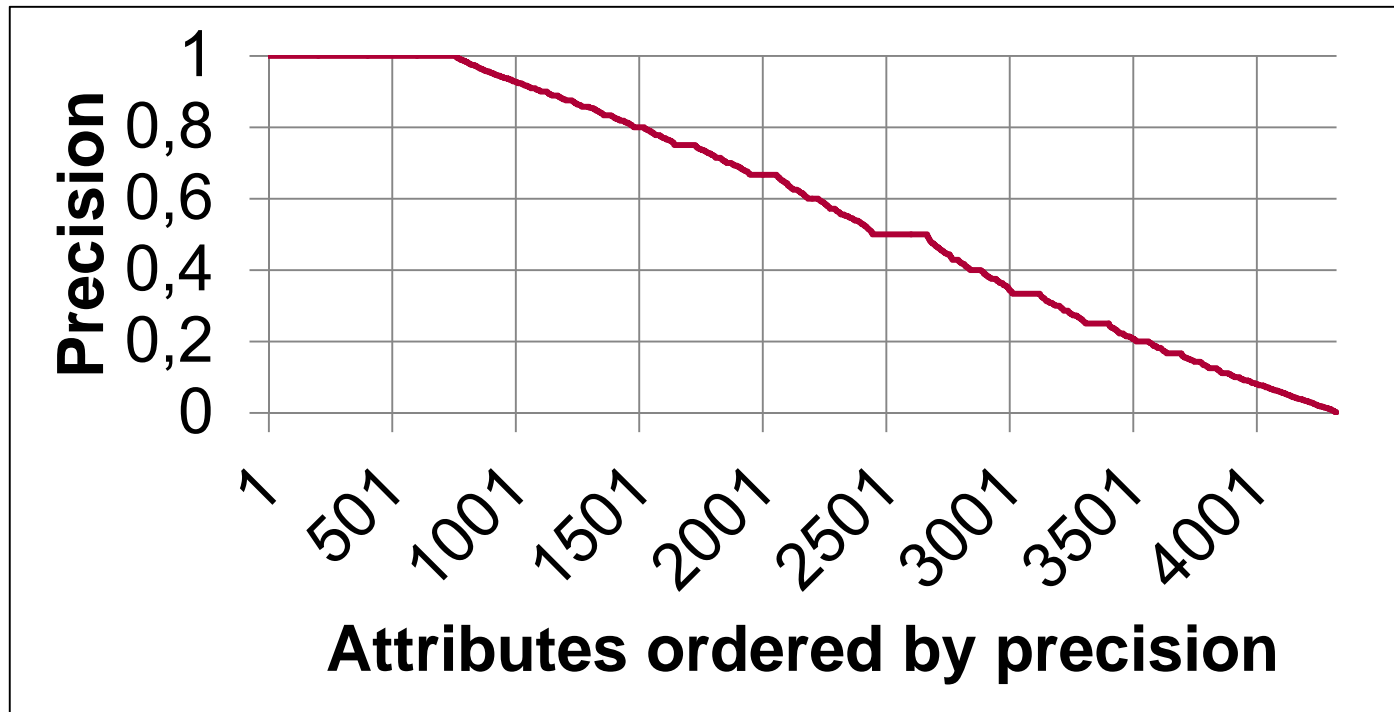


Evaluation: All Infobox-Templates



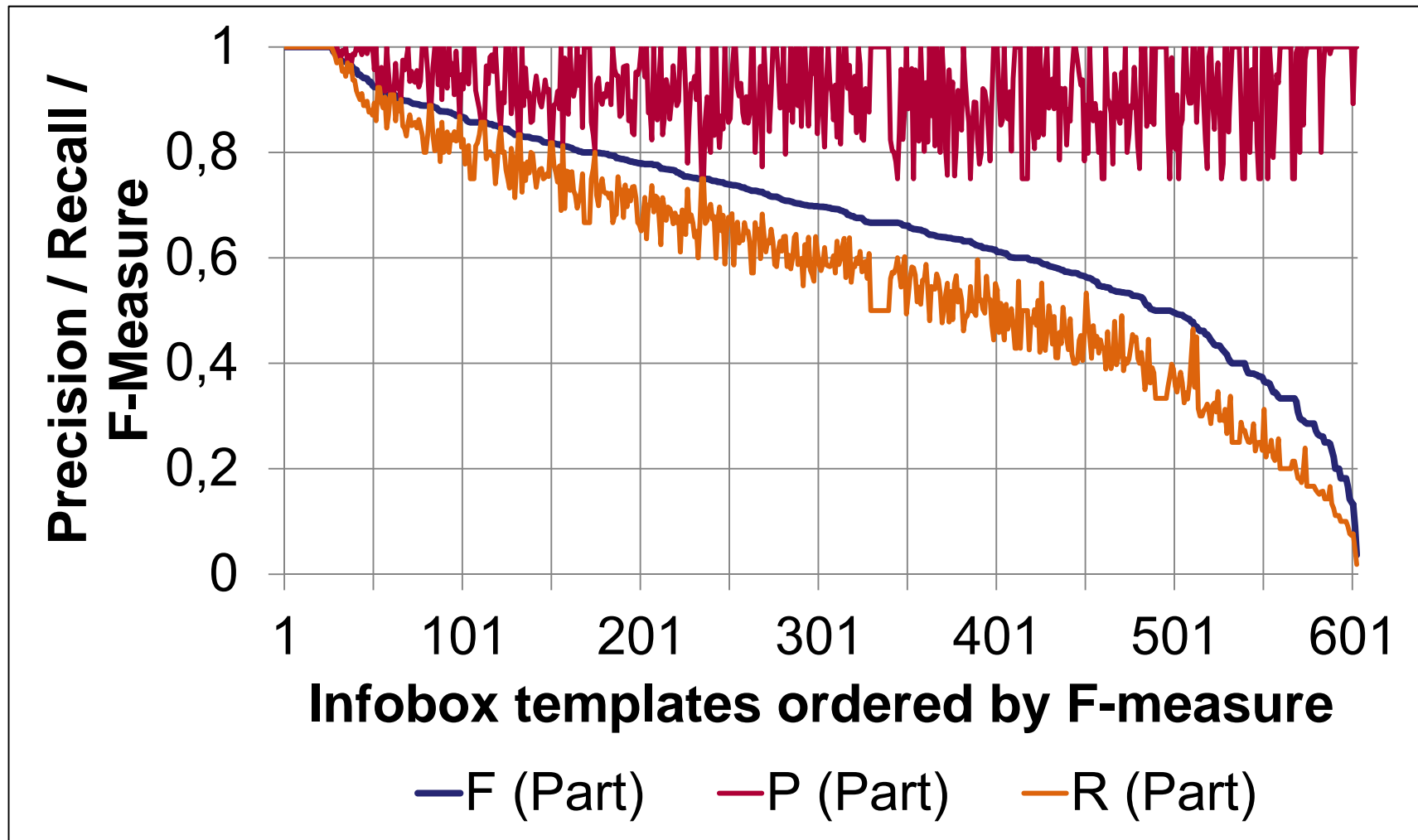
Evaluation on all attributes (>4000) of all infobox templates (>800)

23



Evaluation on all attributes with Precision > 0.75

24



Challenges: Heterogeneity at all levels

26

■ Source

- Formats ↔ □ File converters
- Domain ↔ □ Clustering, rules
- Bandwidth ↔ □ Patience

■ Schema

- Structure ↔ □ Schema Mapping
- Semantics ↔ □ Domain knowledge

■ Data

- Formatting ↔ □ Scrubbing
- Duplicates ↔ □ Entity Matching

Now: Examples for each

The problem – a format mess

Commitment position key: SI2.514875.1

Year:	2008	Amount €:	99.965.021,40
Subject of grant or contract: 2007-EU-50010-P EasyWay * - K(2008) 8479			
Responsible Department:	Trans-European Transport Network Executive Agency	Budget line name and number:	Financial support for projects of common interest in the trans-European transport network (06.03.03)
Programme:	TEN Transport	Co-financing rate:	100,00 %

Beneficiary

Name:	ANONYMI ETAIREIA EKMETALLEFSIS KAIDIACHEIRISIS ELLINIKON AFTOKINITODROMON*TEO AE SOCIETE ANONYME OF HELLENIC MOTORWAYS		
Address:	14342 ATHINA, VITNIS STREET 14-18	Country / Territory:	Greece
Name:	BUNDESREPUBLIK DEUTSCHLAND*REPUBLIQUE FEDERALE D ALLEMAGNE FEDERAL REPUBLIC OF GERMANY		
Address:		Country / Territory:	Germany

Name:	CESKA REPUBLIKA*REPUBLIC OF CZECHIA
Address:	

```

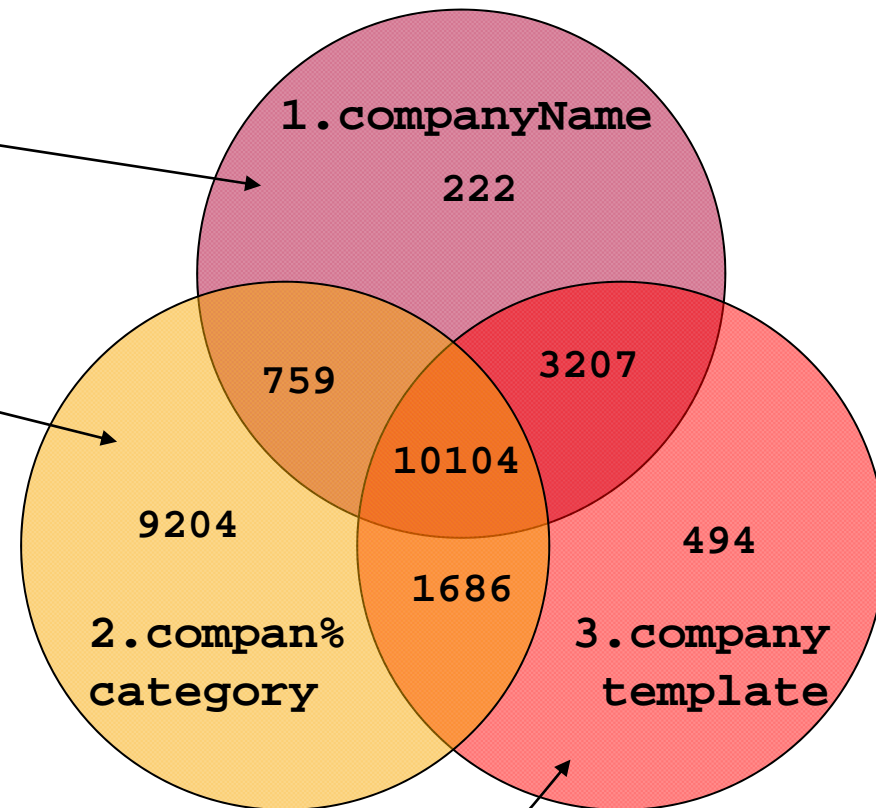
{
  "_id" : "euFinance#28994",
  "year" : 2008,
  "nameOfBeneficiary" : "ROBERT BOSCH GMBH*",
  "coordinator" : false,
  "countryTerritory" : "Germany 70049 STUTTGART",
  "coFinancingRate" : "67,51 %",
  "amount" : 3199959.00,
  "commitmentPositionKey" : "F13.A22622.1",
  "subjectOfGrantOrContract" : "MULTISPECTAL TERAHERTZ, INFRARED ...",
  "responsibleDepartment" : "Information Society and Media",
  "budgetLineNameAndNumber" : "Support for research ..."
}

```

The problem – a domain mess 2008

28

- What is a company?
- Def. 1: Entities having a `companyName`
 - 14292 companies
- Def. 2: Entities in a category that starts with `'compan%'`
 - 21753
- Def. 3: Entities having a `wikiPageUsesTemplate` with value `Template:infobox_company`
 - 15491



The problem – a domain mess 2011

29

- What is a company? 35,588 candidates

- Def. 1: Entities having a %companyName%

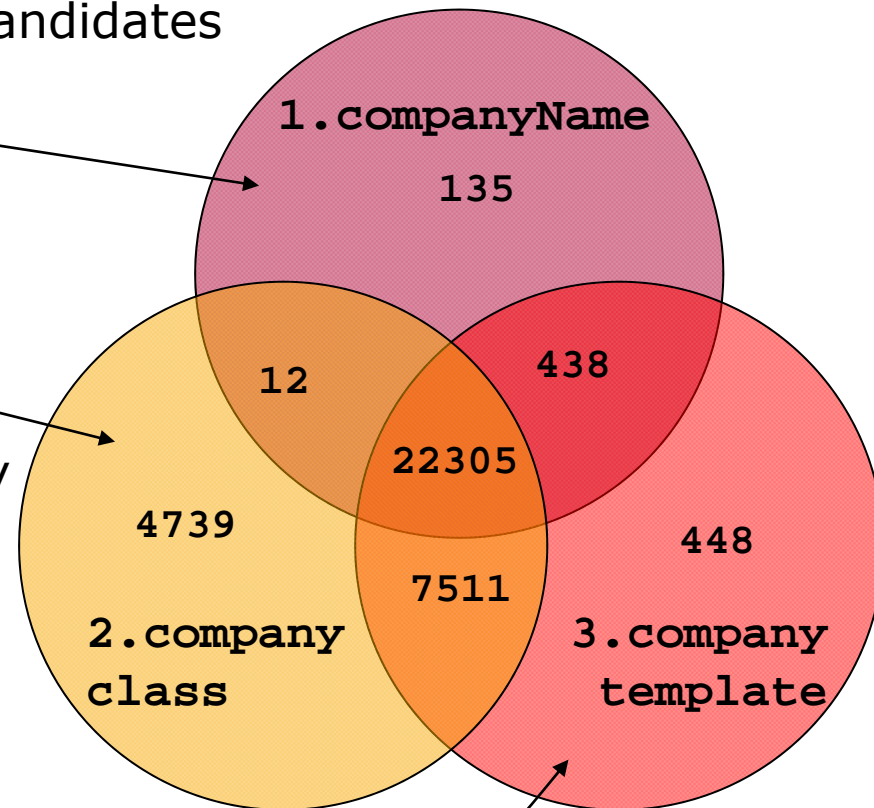
- 22,890

- Def. 2: "Company" according to DBpedia ontology

- 34,567

- Def. 3: Entities having a wikiPageUsesTemplate with value %compan%

- 30,702



Company Template

30

```

{{Infobox Company
| name           = The Corporation Company
| logo           = [[Image:Example.png|160px]]
| type           = [[Public company|Public]] {{{nyse|TCC1}}, {{{tyo|TCC1}}}
| genre          = Corporate histories
| predecessor    = The Wikitory Company
| foundation     = [[New York City]], [[United States|U.S.]] {{{Start date|1900}}}
| founder       = Wikiped Wikiad
| location_city  = [[Seattle]], [[Washington]]
| location_country = [[United States|U.S.]]
| location       =
| locations      = 300 stores (2000) at [[2000-12-31]]
| area_served   = [[North America]]
| key_people     = Wikiped Wikiad <small>[[Entrepreneur|Founder]]</small> <br />
                 Waldo Wikiad <small>[[Chief executive officer|CEO]]</small>
| industry      = [[Publishing]]
| products      = [[Book]]s, [[magazine]]s
| services      = Literary restoration, literary archiving
| revenue       = US$500,000,000 (2000), {{{increase}} 5% from 1999
| operating_income = US$350,000,000 (2000) {{{steady}} from 1999
| net_income    = US$50,000,000 (2000) {{{decrease}} 12% from 1999
| assets        = US$1,500,000,000 at [[2000-12-31]] {{{decrease}} 9% from year earlier
| equity        = US$950,000,000 at [[2000-12-31]] {{{increase}} 6% from year earlier
| owner         = Wikiped Wikiad
| num_employees = 1,500 (2000)
| parent        = Mega Corporation Inc.
| divisions     = TCC Company Histories, TCC Magazine Services
| subsid        = Restored Book Company, Super Archives, Ltd.
| homepage      = [http://www.thecorporationcompany.com/ TheCorporationCompany.com]
| footnotes    =
| intl         =
}}

```

Vertical list	Requirements
<pre> {{Infobox Company name = logo = type = genre = fate = predecessor = successor = foundation = founder = defunct = location_city = location_country = location = locations = area_served = key_people = industry = products = services = revenue = operating_income = net_income = aum = assets = equity = owner = num_employees = parent = divisions = subsid = homepage = footnotes = intl = }} </pre>	<p>REQUIRED</p> <p>REQUIRED</p> <p>REQUIRED</p> <p>REQUIRED</p> <p>REQUIRED</p>

The problem – a schema mess

31

- Wikipedia/DBpedia: Triples and ill-defined templates invite disaster.
 - Schema chaos: Many attribute synonyms
 - Hundreds of different attributes
 - Schema misuse: Many attribute homonyms
 - Foundation attribute in DBpedia may contain
 - ◇ Person who founded the company
 - ◇ Year/Date company was founded
 - ◇ Location where the company was found
- `_percent_27_percent_27_percent_27companyName`
 - `_percent_3Cbr/_percent_3ECompanyName`
 - `automatedImagingAssociationCompanyName`
 - `bTcgvuvCompanyName`
 - `bellFoundryCompanyName`
 - `companyNameLocal`
 - `companyNameZh`
 - `companyName_percent_E3_percent_80_percent_80`
 - `companyNames`
 - `dvdEuroCompanyName`
 - `europeanTradeAssociationCompanyName`
 - `iceCreamCompanyName`
 - `itIsExpensiveCompanyName`
 - `publicCompanyName`
 - `companyNameEn`
 - `companyNamesBigBum`
 - `companyName`

Infoboxes with CompanyTemplate

32

- 1083 different attributes
 - 499 appear only once
- Of the 1083 attributes, 39 distinct ones contain 'name' as substring
- 273 companies without any name attribute

location	20617	companyName	13355
products	18176	name	2036
wikiPageUsesTemplate	18048	surname	25
keyPeople	17836	railroadName	8
industry	16822	companyNickname	4
foundation	15826	pastNames	4
homepage	14476	absNameProperty	3
companyType	13433	dnvNameProperty	3
companyName	13355	labelName	3
companyLogo	9006	logoFilename	3
numEmployees	6207	dvdEuroCompanyName	2
revenue	5030	filename	2
locationCity	4098	longName	2
locationCountry	3212	websitename	2
companySlogan	2815	alternativeNames	1
areaServed	2557	birthname	1
relatedInstance	2284	brandName	1
type	2152	bTcgvuvCompanyName	1
parent	2054	companyNameLocal	1
name	2036	companyNamesBigBum	1
netIncome	1663	europeanTradeAssociationCompanyName	1
founder	1597	familyCorporationCompanyName	1
subsidi	1232	formerNames	1
nihongoProperty	1141	fukCompanyName	1
slogan	1087	golfFacilityName	1
coorTitleDmsProperty	960	hangulName	1
logo	925	iceCreamCompanyName	1
services	904	nativeName	1
operatingIncome	896	nickname	1
owner	680	officialName	1
otheruses4Property	510	oldName	1
intl	503	organisationName	1
forProperty	467	publicCompanyName	1
divisions	429	renamed	1
date	422	shortName	1
locations	419	wineryName	1

Infoboxes in Company class 2011

33

- 34567 companies with 455821 triples

- 1729 different attributes

- 894 appear only once

- After cleansing by DBpedia

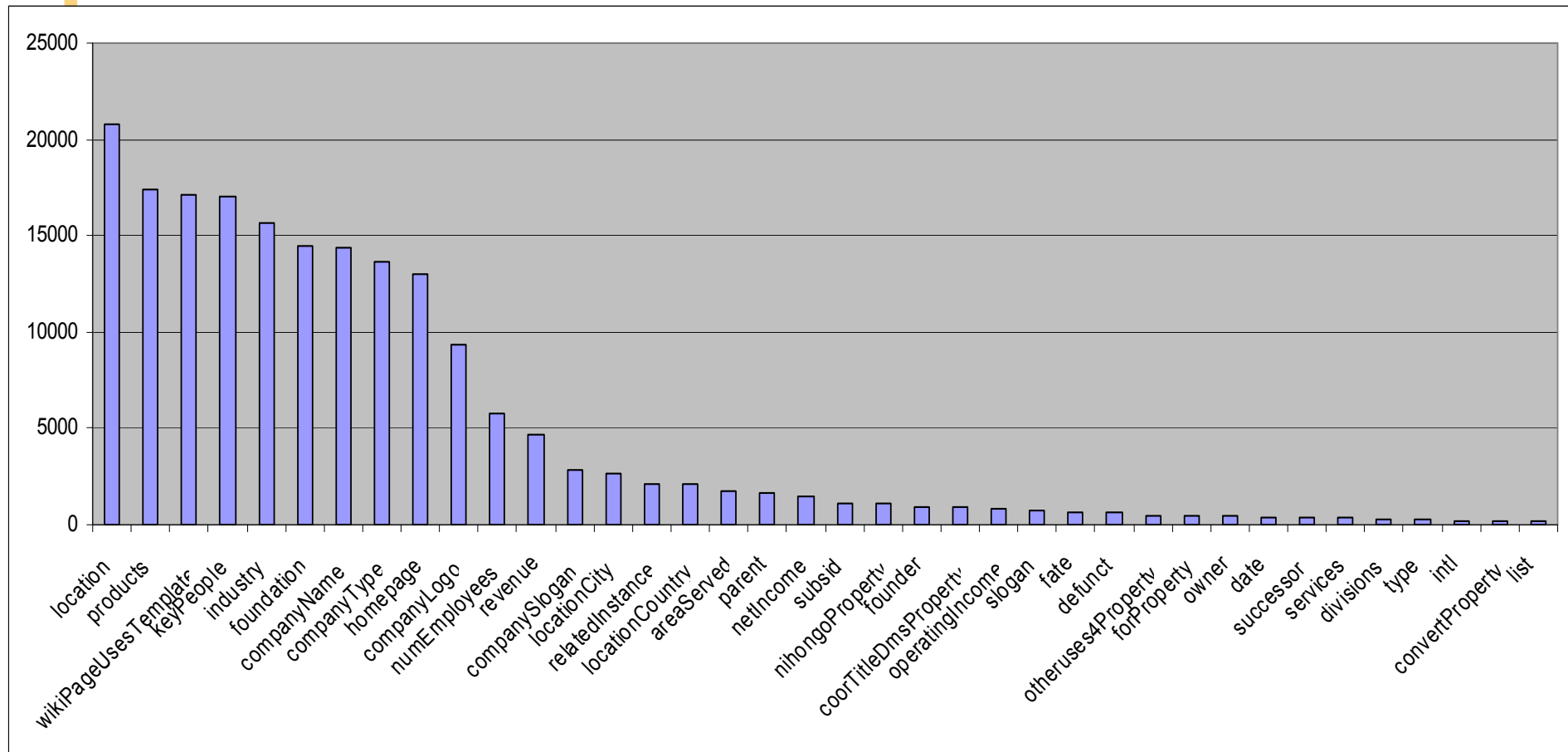
- 34711 companies with 368185 triples

- Only 50 different attributes

- | | |
|------------------------|------------------------|
| ■ keyPeople 34100 | ■ headquarters 3191 |
| ■ industry 28720 | ■ airline 2686 |
| ■ foundation 26875 | ■ services 2568 |
| ■ products 26486 | ■ callsign 2391 |
| ■ homepage 25982 | ■ icao 2386 |
| ■ location 24094 | ■ iata 2363 |
| ■ companyName 23297 | ■ owner 2303 |
| ■ companyType 19591 | ■ fleetSize 2246 |
| ■ companyLogo 14644 | ■ operatingIncome 2246 |
| ■ numEmployees 11395 | ■ hubs 2244 |
| ■ locationCity 9210 | ■ website 2104 |
| ■ name 8700 | ■ intl 1996 |
| ■ locationCountry 7985 | ■ defunct 1987 |
| ■ founder 7867 | ■ fate 1944 |
| ■ revenue 7391 | ■ slogan 1807 |
| ■ parent 6468 | ■ country 1734 |
| ■ type 6358 | ■ destinations 1712 |
| ■ areaServed 5842 | ■ assets 1591 |
| ■ logo 5434 | ■ url 1505 |
| ■ founded 4107 | ■ locations 1384 |
| ■ companySlogan 4053 | ■ divisions 1227 |
| ■ netIncome 3528 | ■ logoSize 1217 |
| ■ genre 3369 | ■ successor 1211 |
| ■ subsid 3288 | ■ distributor 1125 |

Profiling Companies

34



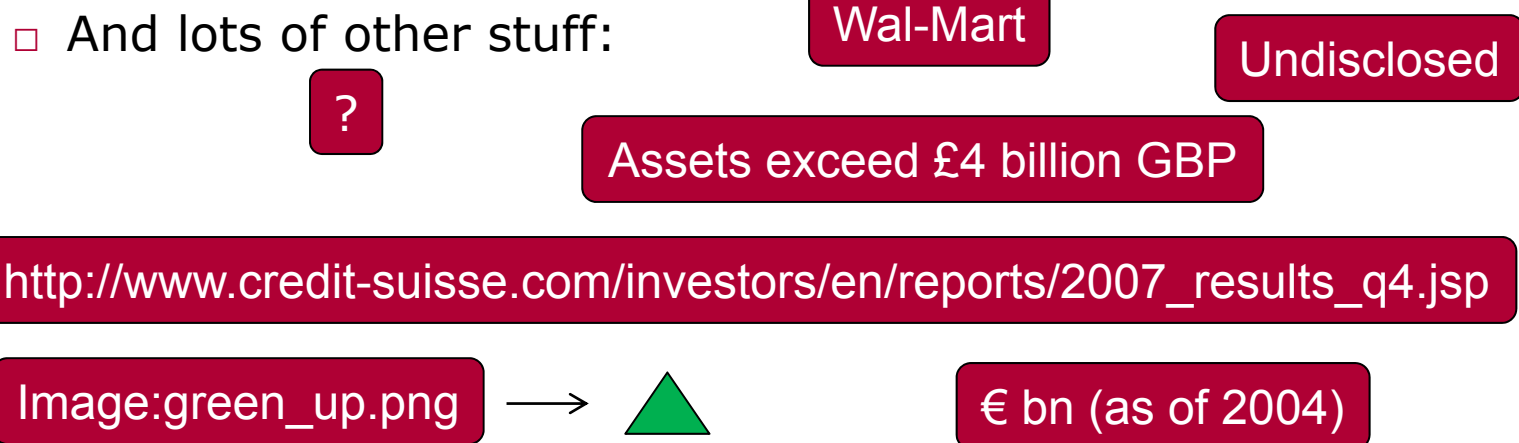
fieldName	<info>	Dollars Obligated	Current Contract Value	Ultimate Contract Value	Major Agency	Modified Contracting Agency	Contracting Agency	Contracting Office	Program / Funding Agency	Program / Funding Office	Reason For Purchase For DoD
example1		\$220,989,132	\$220,989,132	\$220,989,132	Dept. of Defense	97AS: Defense Logistics Agency	Defense Logistics Agency	SP0600	Defense Logistics Agency	SP0600	Invalid code
example2		\$33,710,000	\$33,710,000	\$33,710,000	Dept. of Defense	1700: NAVY, Department of the	NAVY, Department of the	N00024	NAVY, Department of the	N00024	Convenience and Economy
info		add?			kind of category for subagency						
info2		never null	never null	never null	never null, use standardized from modified	never null			Contracting Agency, one contract might have several funding agencies		
scrubbing						split			use Contracting Agency if left blank		
map to LegalEntity as recipient											
map to LegalEntity as Parent recipient											
	subject = "USSpending",		amount.curr	amount.ulti							

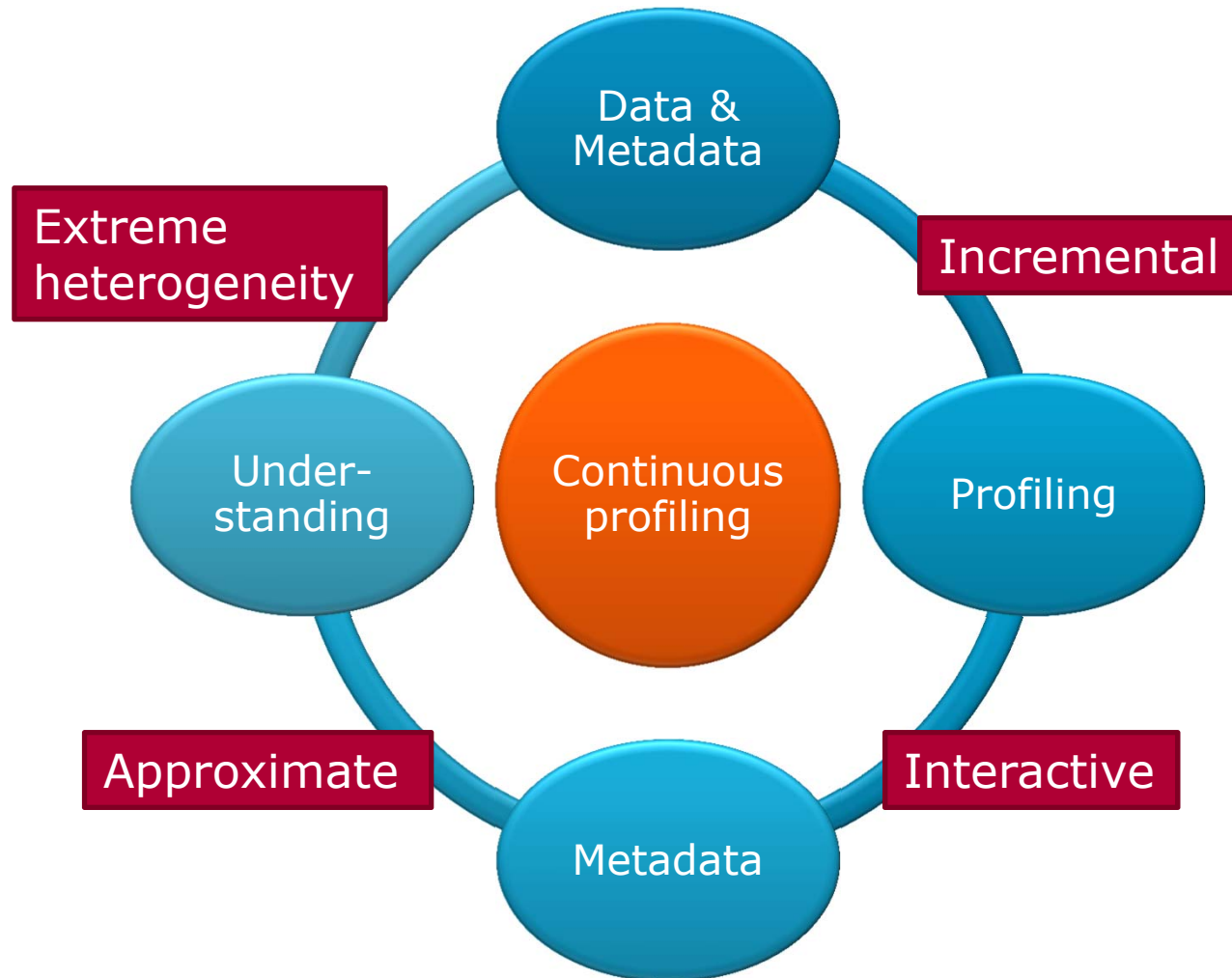


The problem – a data mess

36

- Poor schemata: No types, no constraints
- Sloppy data entry:
 - Data value are neither standardized nor normalized
- Revenue attribute in DBpedia may contain different units, different currencies, and different number-formats.
 - 1.64 billion USD vs. \$1640 m vs. 1,6 vs. more than one million Euro in 2006





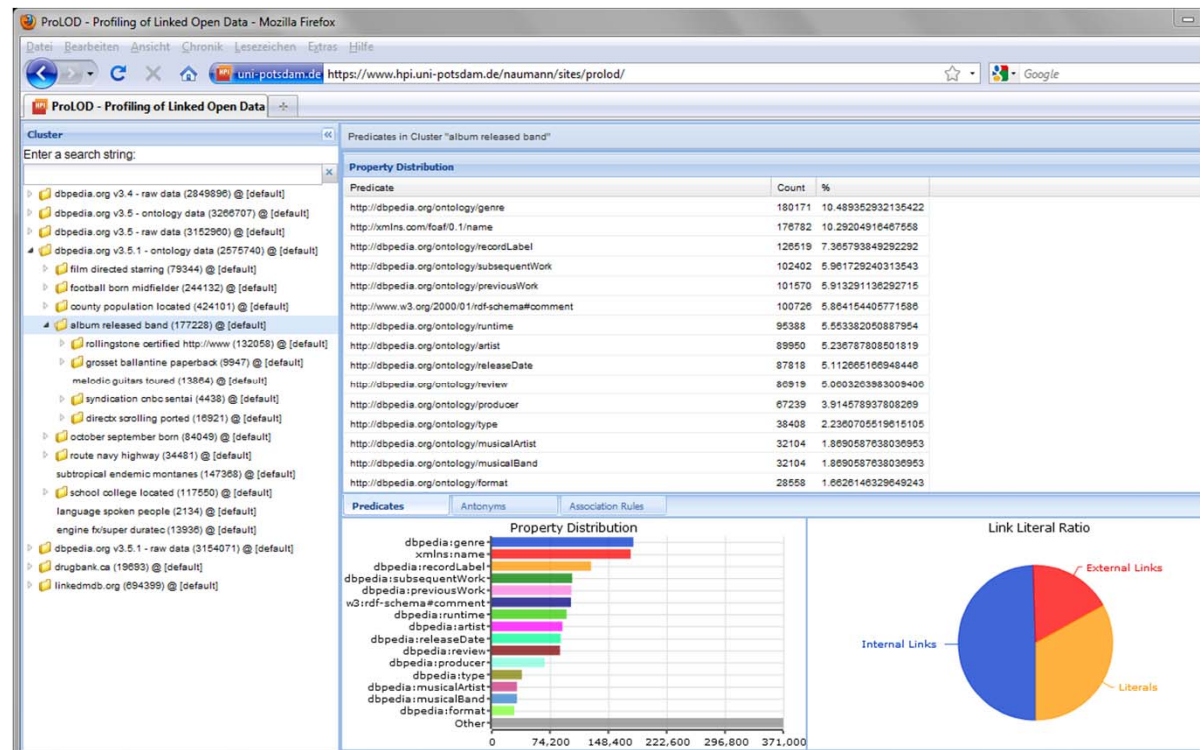
Prototype: ProLOD

39

- Platform for ongoing and future work
 - <https://www.hpi.uni-potsdam.de/naumann/sites/prolod/>

■ Steps:

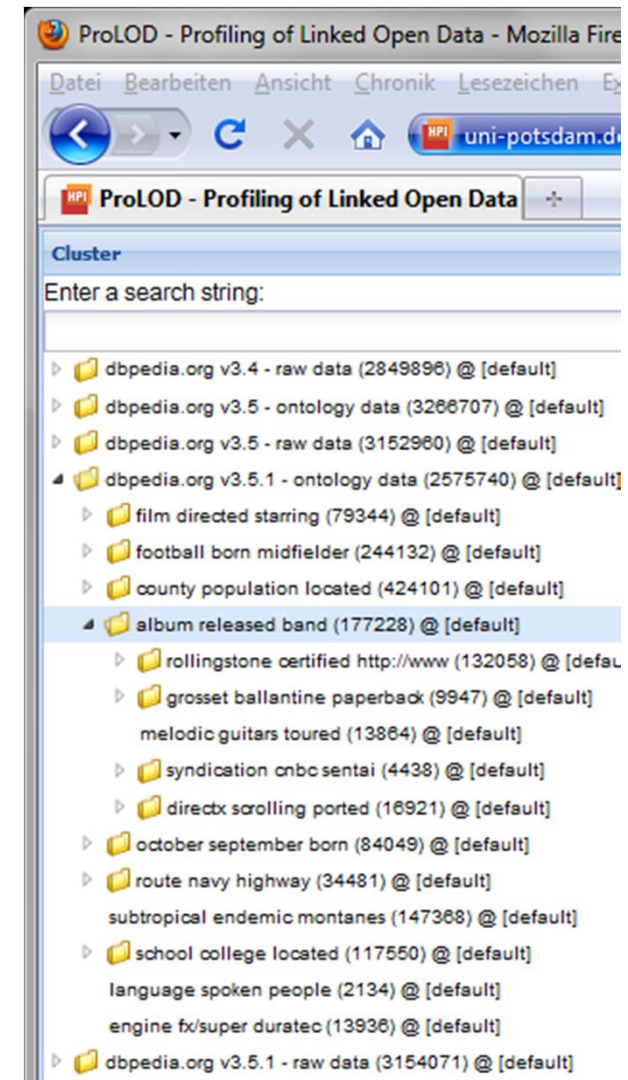
- Data upload
- Preprocessing
- Visualization



ProLOD profiling tasks

40

- Clustering
 - Hierarchical, based on schema
 - Labeling
- Predicate statistics
 - State-of-the-art profiling for attribute values
 - Value types: literals, internal and external links
 - Data types (String, Text, Integer, Decimal, Date)
 - Strings → determine (normalized) patterns
 - Integers, Decimals → display value ranges



ProLOD – Profiling Linked Open Data

ProLOD - Profiling of Linked Open Data - Mozilla Firefox

uni-potsdam.de https://www.hpi.uni-potsdam.de/naumann/sites/prolod/

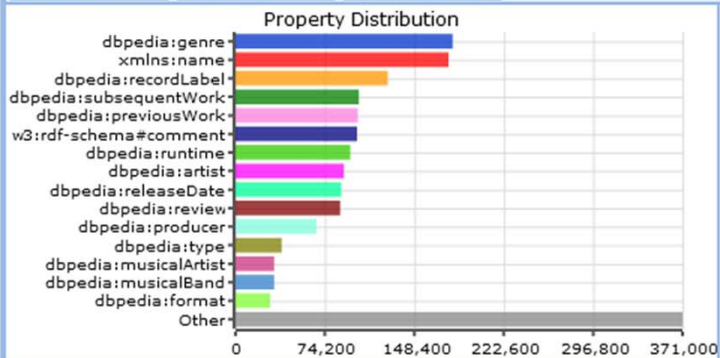
Cluster: album released band (177228) @ [default]

Enter a search string:

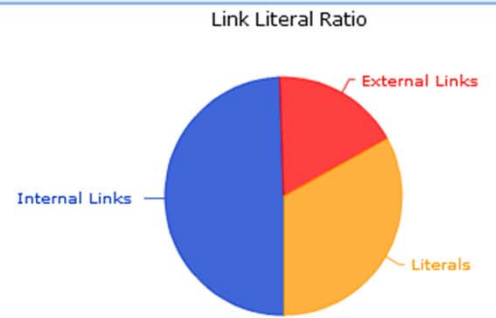
Property Distribution

Predicate	Count	%
http://dbpedia.org/ontology/genre	180171	10.489352932135422
http://xmlns.com/foaf/0.1/name	176782	10.29204916467558
http://dbpedia.org/ontology/recordLabel	126519	7.365793849292292
http://dbpedia.org/ontology/subsequentWork	102402	5.961729240313543
http://dbpedia.org/ontology/previousWork	101570	5.913291136292715
http://www.w3.org/2000/01/rdf-schema#comment	100728	5.864154405771586
http://dbpedia.org/ontology/runtime	95388	5.553382050887954
http://dbpedia.org/ontology/artist	89950	5.236787808501819
http://dbpedia.org/ontology/releaseDate	87818	5.112665166948446
http://dbpedia.org/ontology/review	86919	5.0603263983009406
http://dbpedia.org/ontology/producer	67239	3.914578937808269
http://dbpedia.org/ontology/type	38408	2.260705519615105
http://dbpedia.org/ontology/musicalArtist	32104	1.8690587638036953
http://dbpedia.org/ontology/musicalBand	32104	1.8690587638036953
http://dbpedia.org/ontology/format	28558	1.6626146329649243

Property Distribution (Bar Chart)



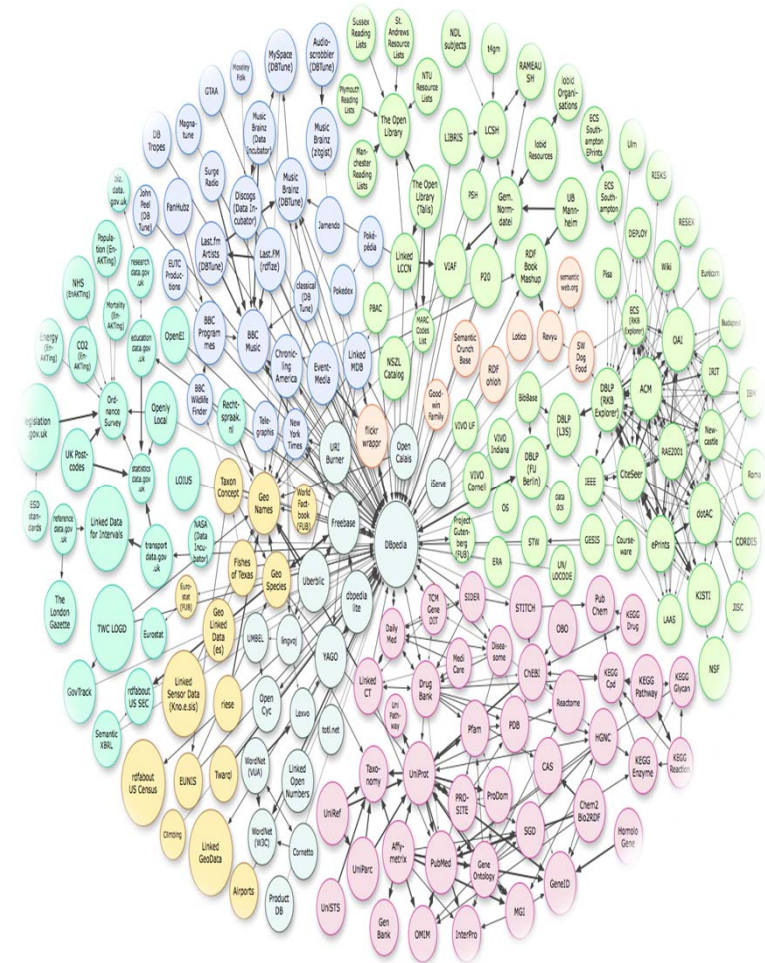
Link Literal Ratio (Pie Chart)



Overview

42

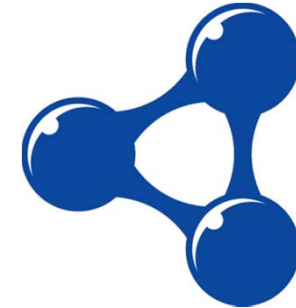
- Web Data abounds
 - Linked, open, and otherwise
 - iPopulator
- Web Data stinks
 - Dirt, grime, and some surprises
 - ProLOD – Profiling LOD
- Cleansing and Integration
 - ...of mops and brooms
 - Cross-Language Integration
- Government data
 - Politicians, friends, and funds
 - The GovWILD experience



Midas – Integration project with IBM Almaden Research Center

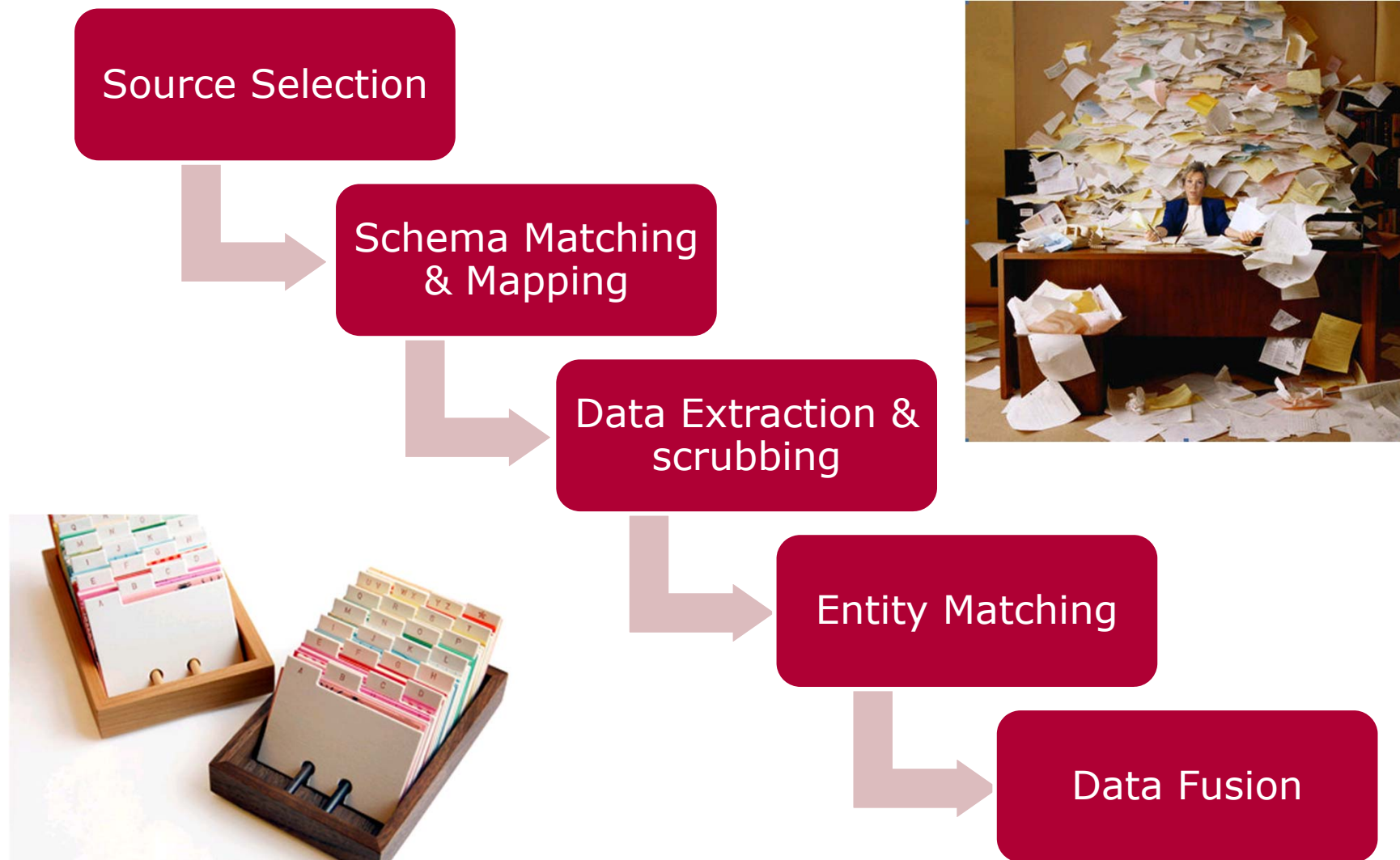
43

- Linked Open Data (Midas, LOD)
 - Integrating DBpedia, Freebase, SEC and FDIC at the level company entities
- Regulatory sources (Midas.Finance)
 - Integrating unstructured/semi-structured data sources containing information about a wide range of entities (e.g., SEC and FDIC)
- Government (Midas.Gov)
 - Integrating structured data from government data sources like usaspending.gov, senate.gov, etc.
 - Persons, legal entities, funding



Five steps for integration

44



Five steps – Source selection

45

- Performed by domain experts
- Criteria
 - Availability and downloadability
 - Coverage of domain (completeness)
 - Complementation with other sources
 - Reputation of source
 - Accuracy of data
 - Cost
 - Other data quality criteria...

Top: Health (57,758)

- [Animal](#) (5,432)

• Alternative (4,700)	• Medicine (10,070)
• Conditions and Diseases (14,289)	• Mental Health (4,577)
• Healthcare Industry@ (5,652)	• Regional (0)

• Addictions (2,302)	• Nutrition (550)
• Aging (77)	• Occupational Health and Safety (423)
• Beauty (432)	• Organizations (132)
• Child Health (433)	• Pharmacy (2,573)
• Conferences (0)	• Products and Shopping (0)
• Dentistry (533)	• Professions (1,337)
• Directories (6)	• Public Health and Safety (3,064)
• Disabilities@ (881)	• Publications@ (131)
• Education (165)	• Reproductive Health (1,812)
• Employment@ (361)	• Resources (106)
• Environmental Health@ (279)	• Search Engines (11)
• Fitness (305)	• Senior Health (647)
• History@ (8)	• Senses (297)
• Home Health (245)	• Services (37)
• Insurance@ (131)	• Specific Substances (581)
• Issues@ (2,003)	• Support Groups (280)
• Medical Tourism@ (67)	• Teen Health (49)
• Men's Health (178)	• Travel Health@ (67)
• News and Media (202)	• Weight Loss (286)
• Nursing (1,109)	• Women's Health (513)

dmoz.org

Five steps – Schema matching and schema mapping

46

- Semi-automated matching
 - Label-based and instance-based

- Challenges:

- Multi-lingual
- Homonyms and Synonyms
- 1:1, 1:n, n:m

- Complex data transformation

Final Schema	DBPedia	SEC	Freebase
dbpediaURI			/type/object/key
cik	secCik	CIK	
irsnumber			
companyName	companyName, name, nonProfitName	name	/type/object/name, /common/ /location/mailling_address/stre
address		BusinessAddress, MailingAddress	/location/mailling_address/pos
locationCity	locationCity, location	BusinessAddress, MailingAddress	/location/mailling_address/city
locationCountry	locationCountry, location, showflag	BusinessAddress, MailingAddress	
telephone		BusinessAddress	
symbol	symbol	Symbol	/business/company/ticker_syn
homepage	homepage, url		
keyPeople (name,title)	keyPeople	KeyPeople	/business/employer/employee /business/company/board_me
industry	industry		industry
products	products, services, genre		
companyType	companyType, type, nonProfitType		company_type
numEmployees	numEmployees, employees		
revenue	revenue		
netIncome	netIncome, grossProfit, earnings, operatingIncome		
foundingYear	foundation, ageProperty		/business/company/founded
fate	fate, currentStatus, end, dissolved, defunct, successor, origins		
companySlogan	companySlogan, motto, slogan		

Five steps – Data extraction & scrubbing

47

- Recognize data types
- Regular expressions for multi-valued strings
- Remove spurious values (layout, formatting, ...)
- Standardize formats
- Translate from foreign languages

Five steps – Entity matching

48

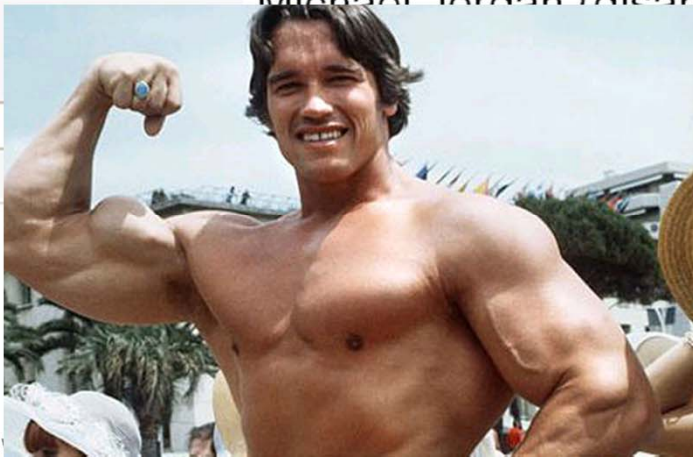
- Duplicate entries
- Linking between entries
- Challenges
 - Fuzzy matching: Similarity measures
 - Data volume: Partitioning algorithms
 - Sparse data
 - ◇ “Michael Jordan visited Indianapolis”



Find People

First Name	*Last Name	City,
<input type="text" value="Michael"/>	<input type="text" value="Jordan"/>	<input type="text" value="CA"/>

Whoa! Over 100 Results Found



Michael Jordan (disambiguation)

...l player.

...ychologist

...English goalkeeper (Ars

...n actor

...researcher in machine

...executive for CBS, Peps

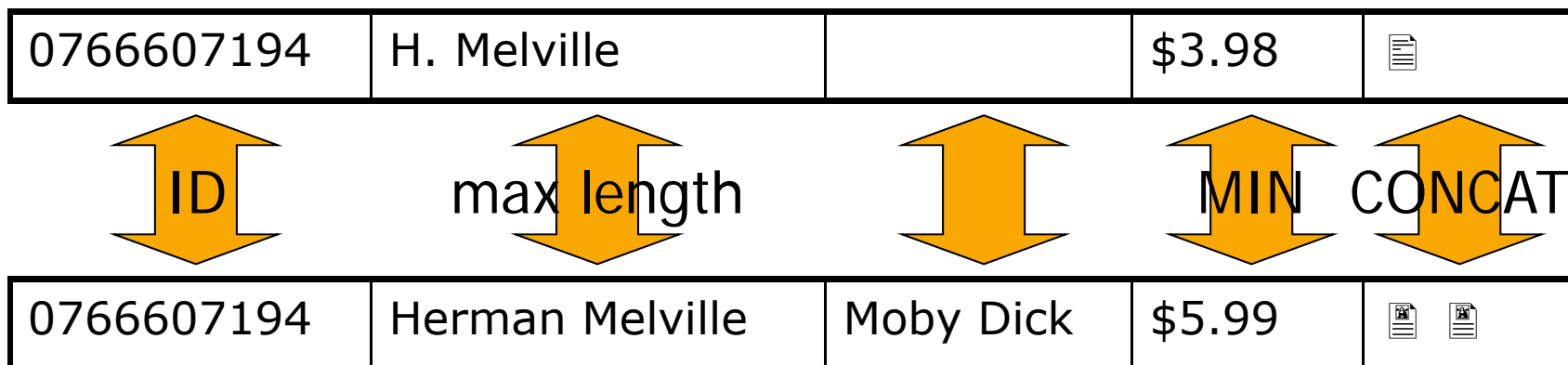
...merican professional bask

...rners' Party TD from Wv

Five steps – Data fusion

49

- Combine multiple representations of real-world entities
 - Survivorship, consolidation, etc.
- Resolve data conflicts
 - Conflict resolution functions
 - Reputation / accuracy / freshness -> "truth discovery"



- Retain data lineage

Multi-Lingual Wikipedia

51

- Goal: Schema matching across languages
 - Complement infobox data
 - Autocomplete for authors
 - Detect errors or inconsistencies
 - Keep values up to date
- Idea: Use cross-language links across 281 languages (Mar 2011)

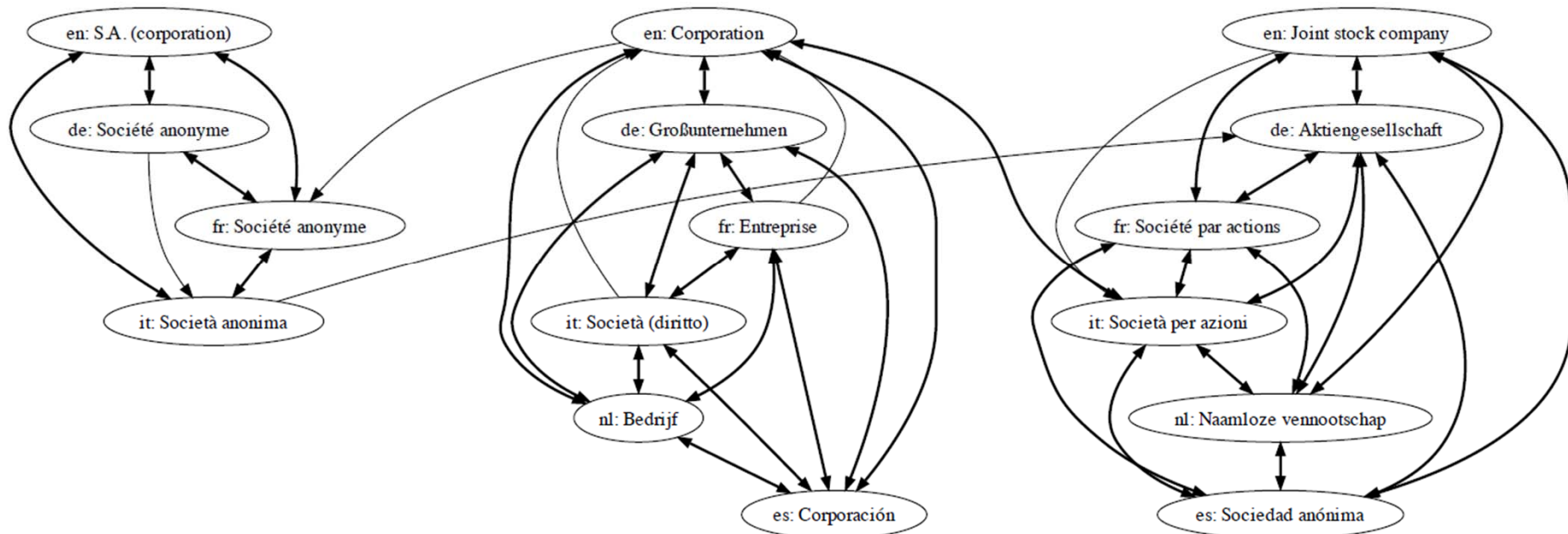


- ▼ Languages
 - العربية
 - Български
 - Català
 - Česky
 - Dansk
 - Deutsch
 - Eesti
 - Español
 - Euskara
 - فارسی
 - Français
 - 한국어
 - हिन्दी
 - Bahasa Indonesia
 - Italiano
 - עברית
 - ಕನ್ನಡ
 - Latviešu
 - Lietuvių
 - Magyar
 - Nederlands
 - 日本語
 - Norsk (bokmål)
 - Polski
 - Português
 - Română

Interlanguage links (ILLs)

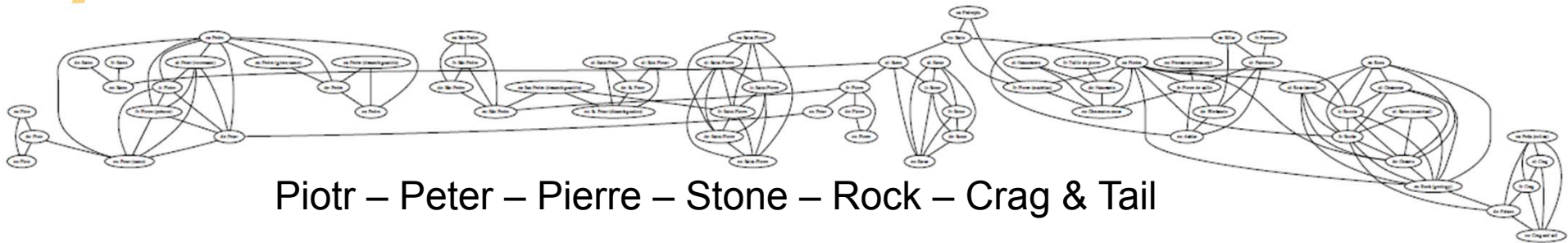
52

- First, evaluate quality of ILLs and build duplicate clusters
 - Build connected components using cross-language links (restricted to the six largest languages)
- But, largest weakly connected component has 108 articles
 - 26 English, 26 German, 21 French, 13 Italian, 13 Dutch, 9 Spanish articles

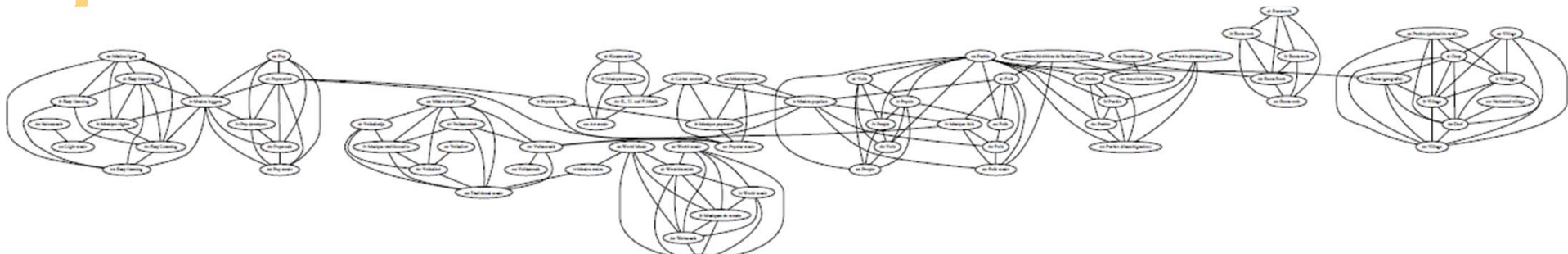


Other large components

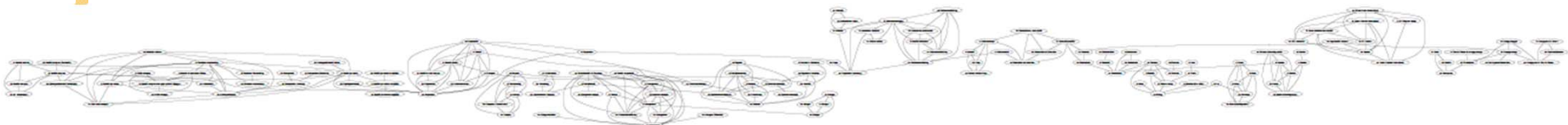
53



Piotr – Peter – Pierre – Stone – Rock – Crag & Tail



Easy Listening – Pop music – World music – Musique folk – Folk – Pueblo - Village

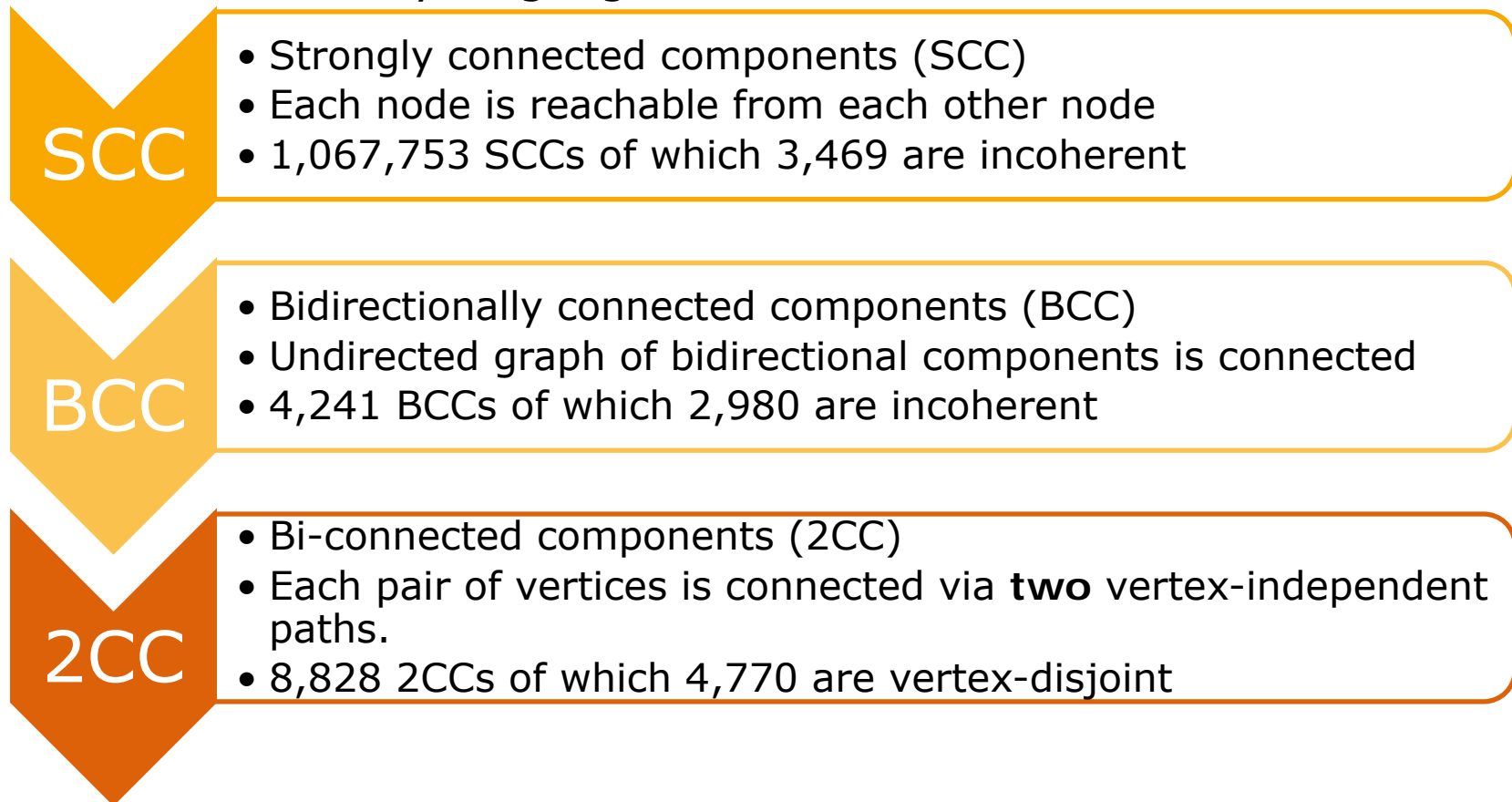


Joint Stock Company – ... – Brother

Whittling down the ILL set

54

- A connected component is **incoherent** if it contains more than one node for any language.

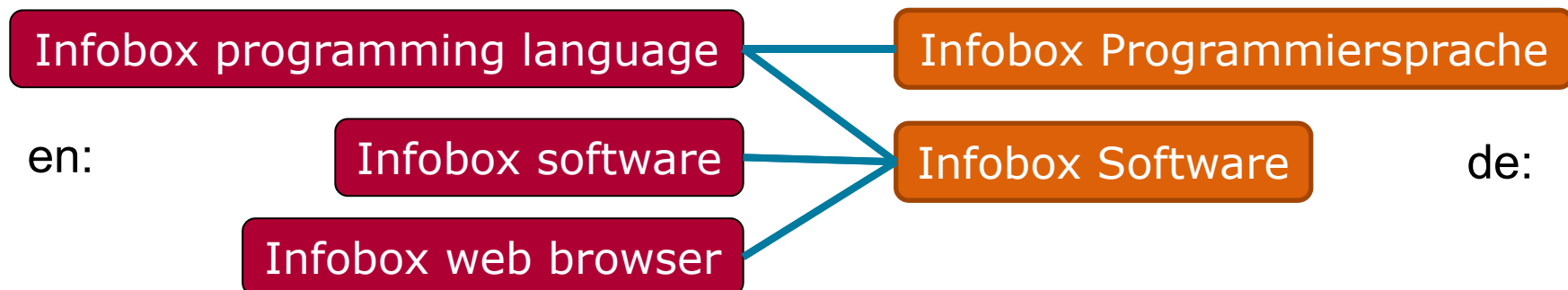


- Result: 1,069,948 coherent, connected components

Infobox Template Mapping

55

- Match schemas of **corresponding** infobox templates only.
- Different granularities in templates => n:m mapping
- **Idea:** Count co-occurrences of infobox templates in terms of connected components
- Apply thresholds:
 - **Absolute:** at least 5 co-occurrences
 - **Relative:** co-occurrence frequency at least 20% of individual occurrences of the templates



Duplicate-based Schema Matching




56

- General technique of data is available under both schemas
- Idea: If data coincides for attributes of two schemata, they probably match.

- For each infobox template pair
 - For each article pair
 - ◇ For each attribute value pair
 - Determine similarity of values (edit-distance)
 - Store in matrix
 - Aggregate similarities across all articles
 - Perform global matching: bipartite assignment

Duplicate-based Schema Matching

57

Coordinates:  52°30'2"N 13°23'56"E		Basisdaten	
Country	Germany	Fläche:	891,85 km ² (14.)
Government		Einwohner:	3.456.264 ^[1] (8.) (31. Oktober 2010)
- Governing Mayor	Klaus Wowereit (SPD)	Bevölkerungsdichte:	3.875 Einw. je km ² (1.) als Bundesland, (2.) als Gemeinde
- Governing parties	SPD / Die Linke	BIP:	90,1 Mrd. € (2009)
- Votes in Bundesrat	4 (of 69)	Höhe:	34–115 m ü. NN
Area		Geografische Lage:	52° 31' N, 13° 24' O
- City	891.85 km ² (344.3 sq mi)	Zeitzone:	Mitteuropäische Zeit (MEZ) UTC+1
Elevation	34 - 115 m (-343 ft)	Postleitzahlen:	10115–14199
Population (31 March 2010) ^[1]		Vorwahl:	030
- City	3,440,441	Kfz-Kennzeichen:	B
- Density	3,857.6/km ² (9,991.3/sq mi)	Gemeindeschlüssel:	11 0 00 000
- Metro	4,429,847	ISO 3166-2:	DE-BE
Time zone	CET (UTC+1)	UN/LOCODE:	DE BER
- Summer (DST)	CEST (UTC+2)	Website:	www.berlin.de 
Postal code(s)	10001–14199	Politik	
Area code(s)	030	Reg. Bürgermeister:	Klaus Wowereit (SPD)
ISO 3166 code	DE-BE	Reg. Parteien:	SPD und Die Linke
Vehicle registration	B	Sitzverteilung im Abgeordnetenhaus	SPD 54 CDU 36
GDP/ Nominal	€ 90.1 ^[2] billion (2009) <small>[citation needed]</small>		
NUTS Region	DE3		
Website	berlin.de 		

Evaluation

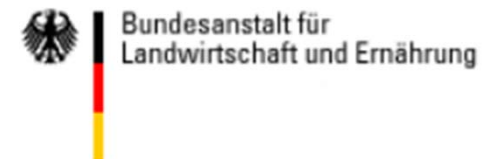
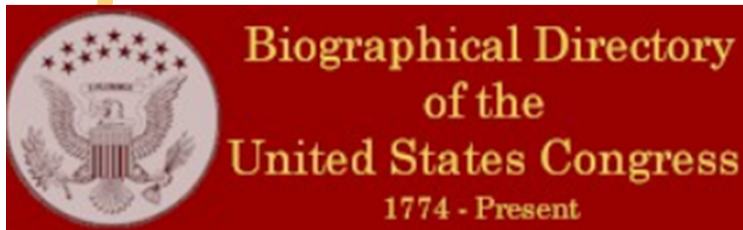
58

- Qualitative evaluation via hand-crafted attribute mappings
 - 96 infobox template pairs
 - 1,417 expected attribute pairs

%	en de	en fr	en nl	de fr	de nl	fr nl	Overall
Precision	91.97	92.28	95.15	90.78	91.67	93.85	92.64
Recall	94.17	96.83	94.80	92.06	93.22	92.82	94.21
F₁ Score	93.06	94.50	94.97	91.42	92.44	93.33	93.42

Motivation – Wealth of Open Gov Data

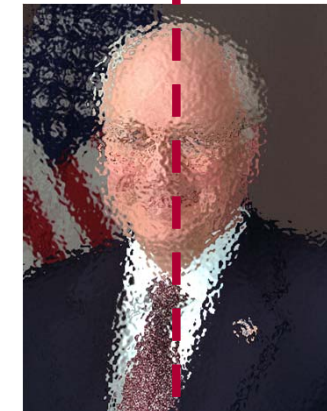
60



Interesting queries

62

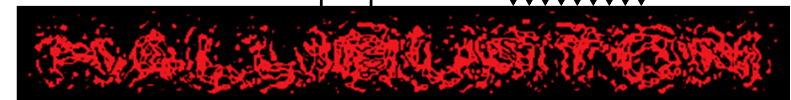
- Find all *classmates* of George W. Bush who, during his term, have worked at a company that has received government funding.
- For each member of congress, find all earmarks awarded to organizations that have *employed a relative* of that member of congress.
- For each government employees, find all companies that have received funding supported by that member and have *employed him after/before their term in congress*.
- Goal: Demonstrate the power of
 - *Joins*: Find unknown connections
`<person - university|company|fund - person>`
 - *Grouping and aggregation*: Combine data about parties, companies, and persons; calculate sums.
 - *Sorting*: Order results by funding amount
 - *Sets*: "for each ... find all ..."



Chairman
of the board

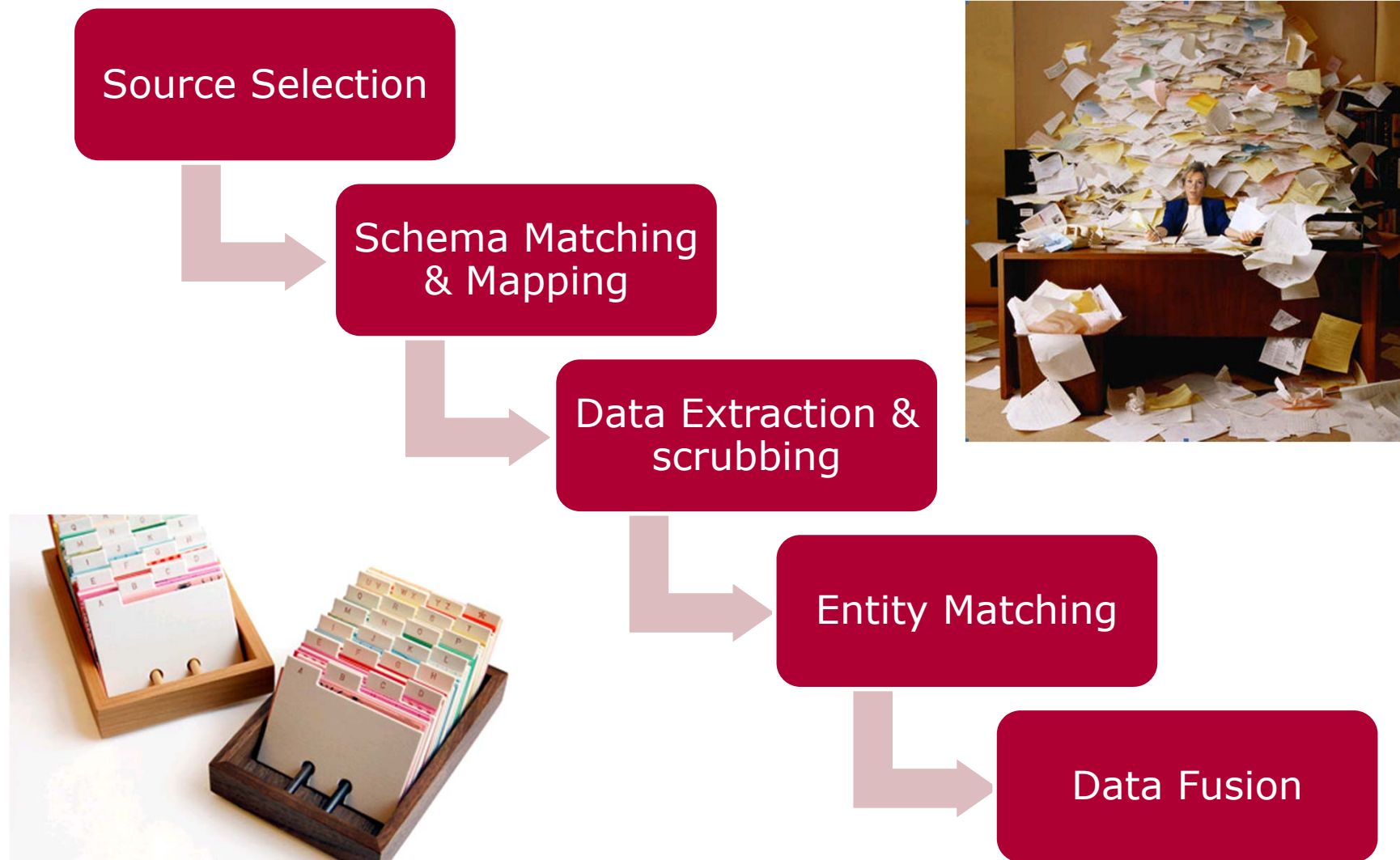
Funds

CEO



Five steps for integration

63



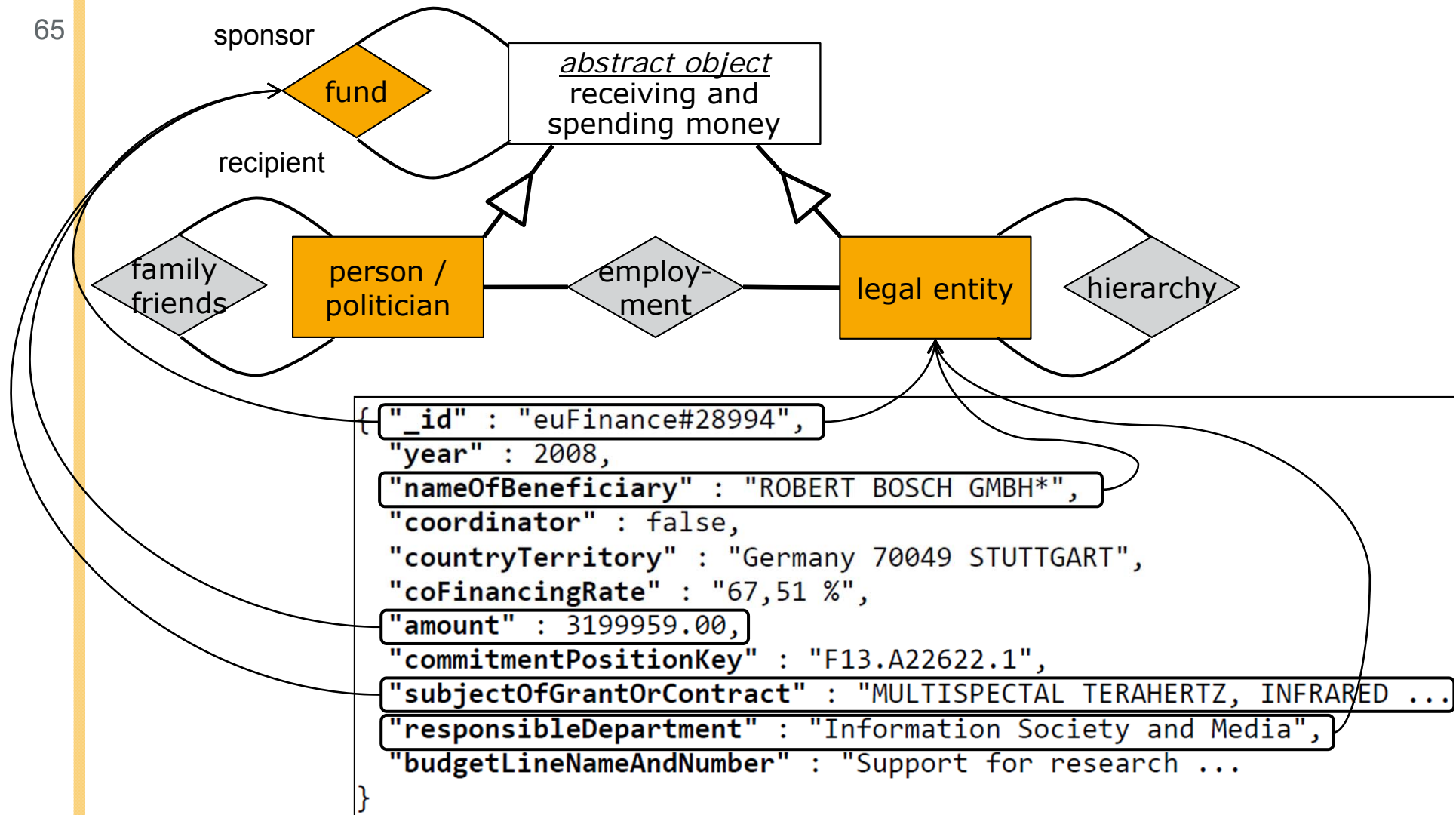
Data sources so far

64

Source	Num. of entities	Num. of attributes	Format	Content
US Spending	1.7m	122	XML	all gov spending
US Earmarks	20,000	37	CSV	anonymous garrantees
US Congress	12,000	8	HTML	members of congress since 1744, incl. bio
DE Party Donations	1,500	4	HTML	Donations > 20,000 €
EU Finance	122,000	11	HTML	EU spending
EU Agric. Subventions	207,000	8	HTML	EU spending
EU Parliam. Data	900	14	HTML	members of parliament
Freebase Person Data	1,8m	32	TSV	person data

Data – Mapping and Scrubbing

65




Data – Cleansing

66

- Deduplication / Entity Matching
 - Intra Source Consolidation
 - Intra Source Duplicate Detection
 - ◇ Duplicate Detection Toolkit – DuDe
 - ◇ Hundreds of duplicates within original sources
 - Entity Matching across Sources
 - ◇ Augment discovered Person Data with Freebase Info
 - ◇ Jaro-Winkler and Monge-Elkan distance
- Entity Fusion
 - ◇ Dempster-Shafer-Theory

<http://govwild.org>

68

- 200,000 persons
 - 248,000 legal entities
 - 1,000,000 funds
- 
- The image shows the GovWILD logo in a stylized, golden, serif font with a shadow effect. Below the logo is a white search input field with a black 'Search' button to its right.
- Keyword Queries
 - Linked Data Interface (dereference URIs)
 - Exploration of entities mentioned in New York Times articles
 - Data Download (RDF, SQL Dump, JSON files)



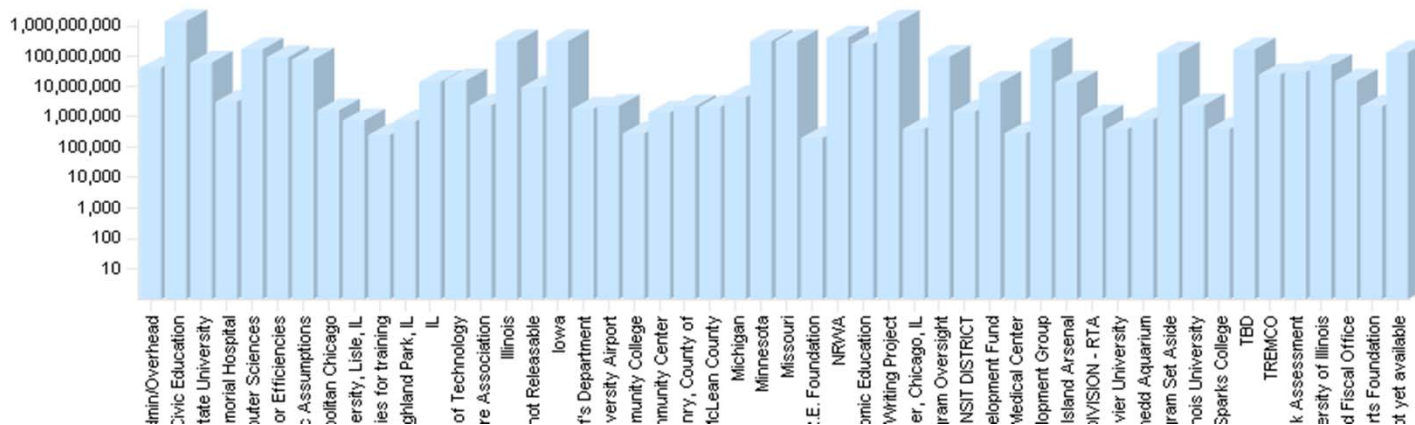
Barack Obama

























Barack Obama

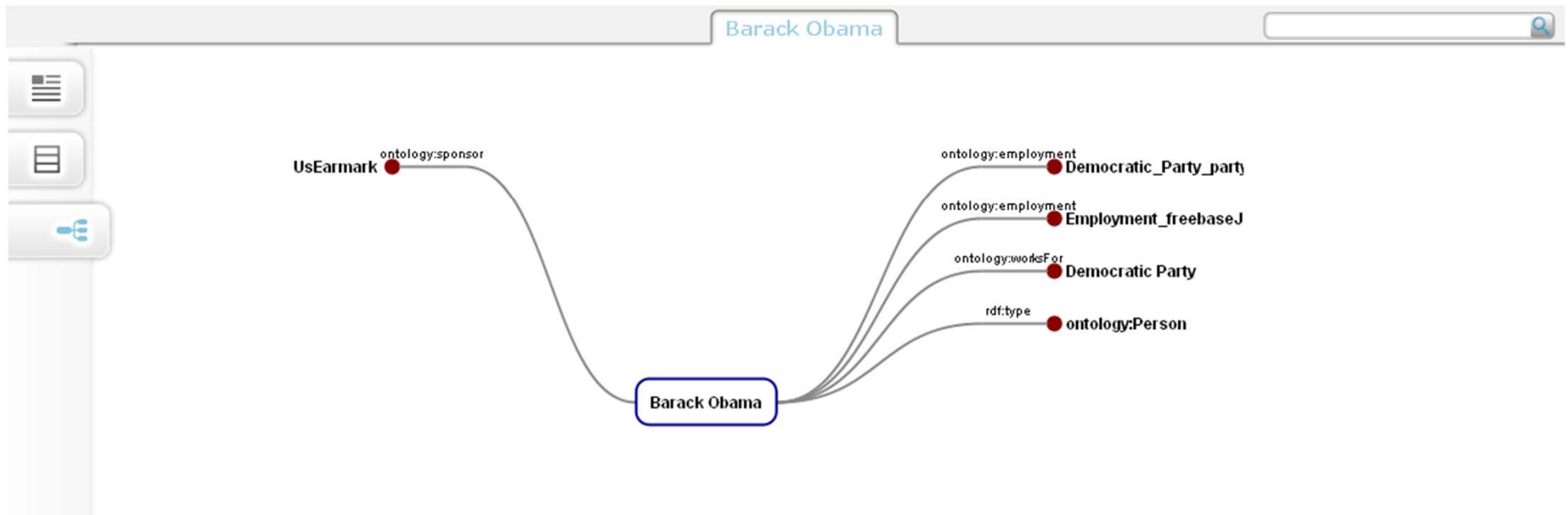
Barack Hussein Obama II (born in 1961) is the 44th and current President of the United States. He is the first African American to hold the office. Obama previously served as a United States Senator from Illinois, from January 2005 until he resigned after his election to the presidency in November 2008. A native of Honolulu, Hawaii, Obama is a graduate of Columbia University and Harvard Law School, where he was the president of the Harvard Law Review. He was a community organizer in Chicago before earning his law degree. He worked as a civil rights attorney in Chicago and taught constitutional law at the University of Chicago Law School from 1992 to 2004. Obama served three terms in the Illinois Senate from 1997 to 2004. Following an unsuccessful bid against a Democratic incumbent for a seat in the U.S. House of Representatives in 2000, he ran for United States Senate in 2004.[1] Several events brought him to national attention during the campaign, including his victory in the March 2004 Democratic primary and his keynote address at the Democratic National Convention in July 2004. He won election to the U.S. Senate in November 2004. His presidential campaign began in February 2007, and after a close campaign in the

2008 Democratic Party presidential primaries against Hillary Rodham Clinton, he won his party's nomination. In the 2008 general election, he defeated Republican nominee John McCain and was inaugurated as president on January 20, 2009.4

Earmarks



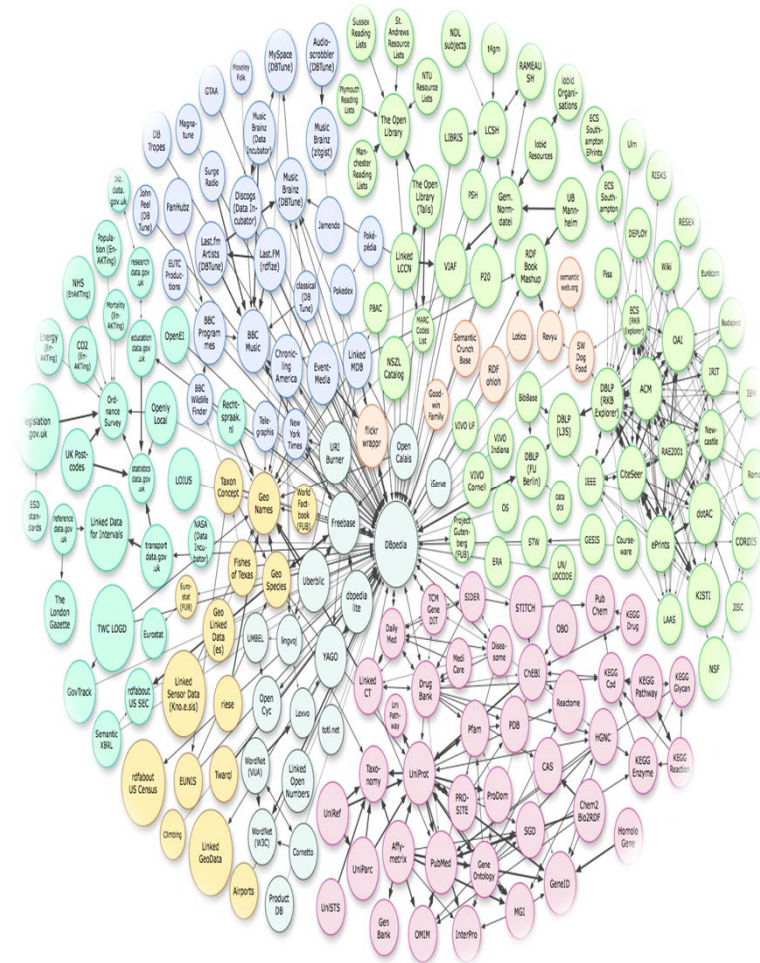
Predicate	Value	Info
ontology:placeOfBirth	Honolulu, Hawaii	
ontology:birthYear	1961	
ontology:lastName	Obama	
rdf:type	ontology:Person	
ontology:employment	Democratic Party party 269	
ontology:employment	Employment freebaseJoinedPersons 478 usStates 99	
rdfs:label	Barack Obama	
ontology:worksFor	Democratic Party	
ontology:birthDay	4	
ontology:birthMonth	8	
ontology:sponsor of	UsEarmark	
ontology:sponsor of	UsEarmark	
ontology:sponsor of	UsEarmark	
ontology:sponsor of	UsEarmark	
ontology:sponsor of	UsEarmark	
ontology:sponsor of	UsEarmark	
ontology:sponsor of	UsEarmark	
ontology:sponsor of	UsEarmark	
ontology:sponsor of	UsEarmark	
ontology:sponsor of	UsEarmark	
ontology:sponsor of	UsEarmark	
ontology:sponsor of	UsEarmark	
ontology:sponsor of	UsEarmark	
ontology:sponsor of	UsEarmark	



Summary

72

- Web Data abounds
 - Linked, open, and otherwise
 - iPopulator
- Web Data stinks
 - Dirt, grime, and some surprises
 - ProLOD – Profiling LOD
- Cleansing and Integration
 - ...of mops and brooms
 - Cross-Language Integration
- Government data
 - Politicians, friends, and funds
 - The GovWILD experience



References

73

- [Extracting Structured Information from Wikipedia Articles to Populate Infoboxes](#)
Dustin Lange, Christoph Böhm, and Felix Naumann
Proceedings of the 19th Conference on Information and Knowledge Management (CIKM) 2010, Toronto, Canada
(Extended version available as [technical report](#))
- [Profiling Linked Open Data with ProLOD](#)
Christoph Böhm, Felix Naumann, Ziawasch Abedjan, Dandy Fenz, Toni Grütze, Daniel Hefenbrock, Matthias Pohl, David Sonnabend
Workshop *New Trends in Information Integration (NTII) 2010, Long Beach, USA*
- [Linking Open Government Data: What Journalists Wish They Had Known](#)
Christoph Böhm, Felix Naumann, Markus Freitag, Stefan George, Norman Höfler, Martin Köppelmann, Claudia Lehmann, Andrina Mascher, and Tobias Schmidt.
[Honorable Mention](#) at Linked Data Triplification Challenge 2010 @ I-Semantics, Graz. (link to [GovWILD](#))
- [DuDe: The Duplicate Detection Toolkit](#)
Uwe Draisbach and Felix Naumann: QDB 2010 Workshop at VLDB, Singapore