# HPI Hasso Plattner Institut

IT Systems Engineering | Universität Potsdam

# Extreme Web Data Integration
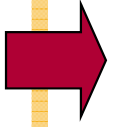
# Keynote @ ICWE 2012

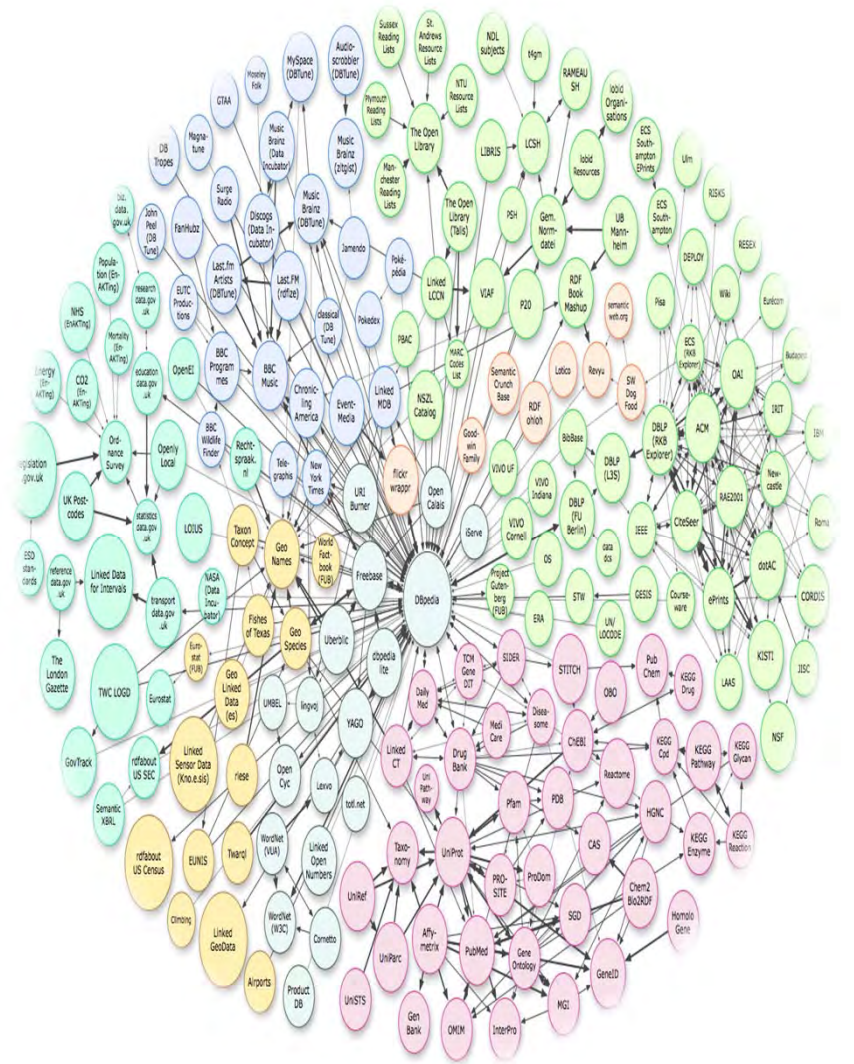26.7.2012

Felix Naumann
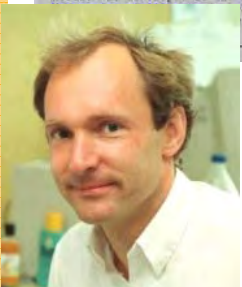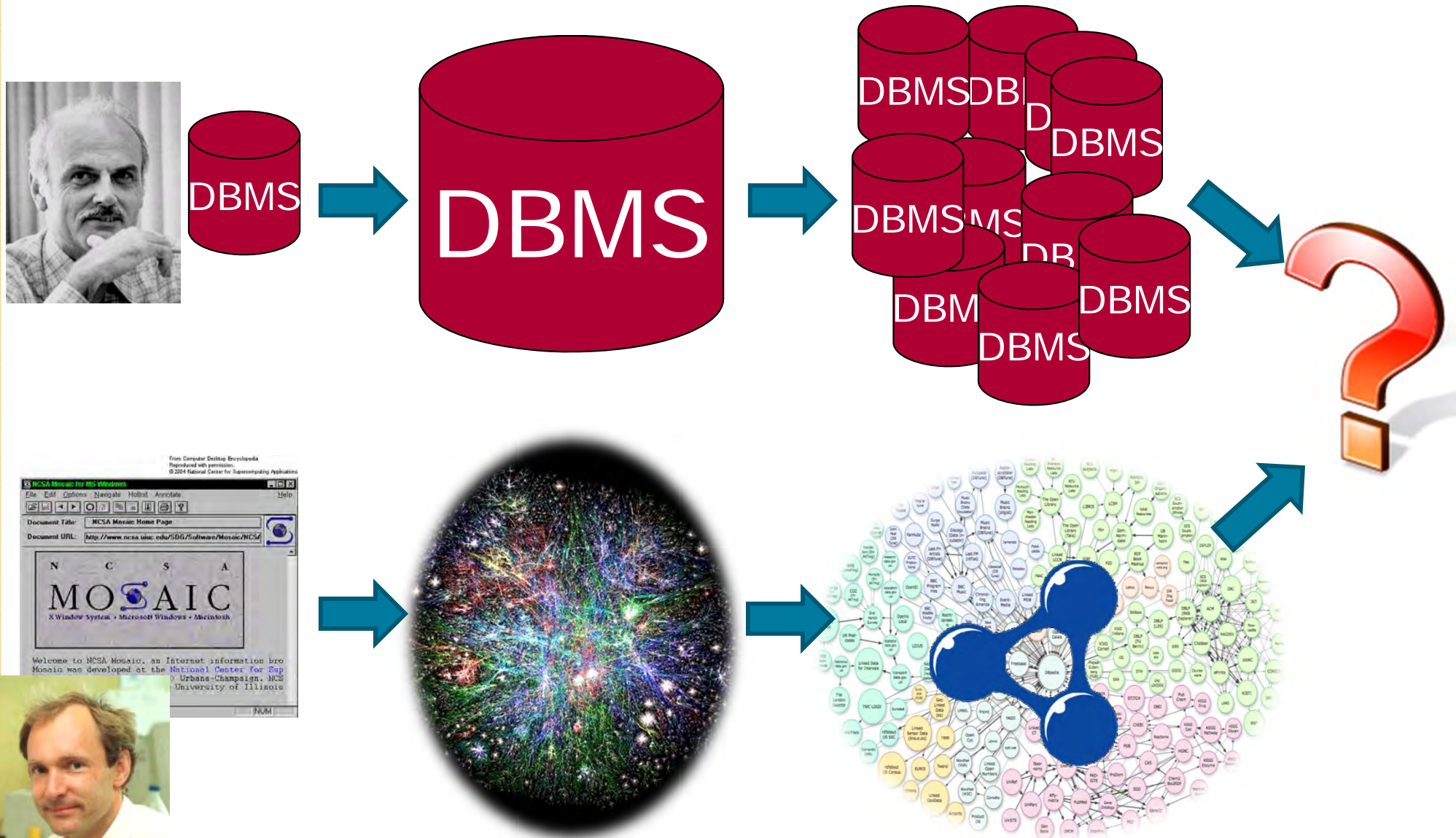
# Overview

- **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

# A brief history of data

# Linked Data & Data Spaces:
## A database guy's point-of-view

**Dataspaces / Data integration**

- Some schema
- Integrated
- Ad-hoc
- Data quality
- High accessibility

**Linked data**

**Relational databases**

**Semantic Web**

# Linked data – 4 Principles, 7 Properties

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://.../Berlin
     P: location
     O: http://.../Germany
     ```
   - ☐ Often user generated
   - ☐ Nice domains
   - ☐ Free
- ■ The Bad
   - ☐ Voluminous
   - ☐ Heterogeneous
- ■ The Ugly
   - ☐ Dirty, inconsistent, sparse

Linking Open Data cloud diagram, by Richard Cyganiak and Anja Jentzsch. http://lod-cloud.net/

Media
Geographic
Publications
User-generated content
Government
Cross-domain
Life sciences

As of September 2011

# DBpedia – Extraction

```
{{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

## 1. Core Datasets

| Dataset | en | de | fr | es | it | pl | nl | pt | sv | ja | | zh | fi | bp |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| Titles ( preview ) | nt csv | nt csv | nt csv | nt csv | nt csv | nt csv | -- | nt csv | nt csv | nt csv | nt csv | nt csv | nt csv | nt csv |
| Short Abstracts ( preview ) | nt - | nt - | nt - | nt - | nt - | nt - | -- | nt - | nt - | | | | | |
| Extended Abstracts ( preview ) | nt - | nt - | nt - | nt - | nt - | nt - | -- | nt - | nt - | nt - | nt - | nt - | nt - | nt - |
| Images ( preview ) | nt csv | -- | -- | -- | -- | -- | -- | | | | | | | |
| Links to Wikipedia Article ( preview ) | nt csv | nt csv | nt csv | nt csv | nt csv | nt csv | -- | nt csv | | | | | | |
| Articles Categories ( preview ) | nt csv | -- | -- | -- | -- | -- | -- | | | | | | | |
| External Links ( preview ) | nt csv | -- | -- | -- | -- | -- | -- | | | | | | | |
| Infoboxes ( preview ) | nt csv | nt csv | nt csv | nt csv | nt csv | nt csv | -- | nt csv | nt csv | nt csv | nt csv | nt csv | | |
| Properties ( preview ) | nt csv | nt csv | nt csv | nt csv | nt csv | nt csv | | nt csv | nt csv | nt csv | nt csv | nt csv | nt csv | |
| DBpedia Ontology ( preview ) | owl | -- | -- | -- | -- | -- | -- | | | | | | | |
| Ontology Infoboxes ( preview ) | nt | -- | -- | -- | -- | -- | | | | | | | | |
| Ontology Types ( preview ) | nt | -- | -- | -- | -- | -- | -- | | | | | | | |
| Homepages ( preview ) | nt csv | nt csv | nt csv | -- | -- | -- | -- | | | | | | | |
| Geographic Coordinates ( preview ) | nt csv | nt csv | nt csv | nt csv | nt csv | nt csv | | nt csv | nt csv | nt csv | nt csv | nt csv | nt csv | nt csv |
| Pagelinks ( preview ) | nt csv | nt csv | nt csv | nt csv | nt csv | nt csv | | nt csv | nt csv | nt csv | nt csv | nt csv | nt csv | nt csv |
| Persondata ( preview ) | nt csv | nt csv | -- | -- | -- | -- | -- | | | | | | | |
| Redirects ( preview ) | nt csv | -- | -- | -- | -- | -- | -- | | | | | | | |
| Disambiguation Links ( preview ) | nt | | | | | | | | | | | | | |

- 1 billion triples
  - 385 million English
- From 97 languages of Wikipedia
- 3.6 million things
  - 416,000 persons
  - 526,000 places
  - 106,000 music albums
  - 60,000 films
  - 17,500 video games
  - …
- http://wiki.dbpedia.org/Datasets

# And more sources

- Government data
  - □ www.data.gov
    450k data sets
  - □ 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, …

Browse Raw Datasets

| | Name | | Popularity | Type |
|---|---|---|---|---|
| 1. | **Worldwide M1+ Earthquakes, Past 7 Days** Geography and Environment ANSS, geologist, plate, real time, environment, … Real-time, worldwide earthquake list for the past 7 days | | 167,711 views | |
| 2. | **U.S. Overseas Loans and Grants (Greenbook)** Foreign Commerce and Aid foreign assistance, economic assistance, Greenbook, … These data are U.S economic and military assistance by country from 1946 to 2010. | | 62,348 views | |
| 3. | **CMS Medicare and Medicaid EHR Incentive Program, electronic health record products used for attestation** Science and Technology electronic health record, … Data set merges information about the Centers for Medicare and Medicaid Services, | | 34,285 views | |
| 4. | **Federal Data Center Consolidation Initiative (FDCCI) Data Center Closings 2010-2013** Federal Government Finances and Employment fddci, … Federal Data Center Consolidation Initiative (FDCCI) Data Center Closings 2010-2013 | | 32,648 views | |
| 5. | **TSCA Inventory** Geography and Environment new chemicals, manufactured chemicals, … This dataset consists of the non confidential identities of chemical substances | | 27,007 views | |
| 6. | **Data.gov Catalog** Other dataset, metadata, catalog, data extraction tool, … An interactive dataset containing the metadata for the Data.gov raw datasets and tools | | 23,117 views | |
| 7. | **US DOE/NNSA Response to 2011 Fukushima Incident: Radiological Air Samples** Geography and Environment radiation, Japan, nuclear, Tohoku, … Field Samples are physical media collected during the response which are | | 22,458 views | |
| 8. | **US DOE/NNSA Response to 2011 Fukushima Incident: Field Team Radiological Measurements** Geography and Environment Japan, nuclear, Tohoku, radiation, … Field Measurements describe &alpha; and &beta; activity and &gamma; exposure rate. | | 20,940 views | |
| 9. | **Federal Executive Branch Internet Domains** Federal Government Finances and Employment .gov, domains, agencies, federal, registered Listing of Federal Agency Internet Domains (This list is updated bi-weekly to reflect the | | 17,267 views | |

Killer app?

# Overview

- **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

# 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 "thoughtcrime", and "Newspeak" have entered th itself has come to refer to anything reminiscent generally considered to be George Orwell's mag

iPopulator

## Contents [hide]

**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** | 978-0452284234 |

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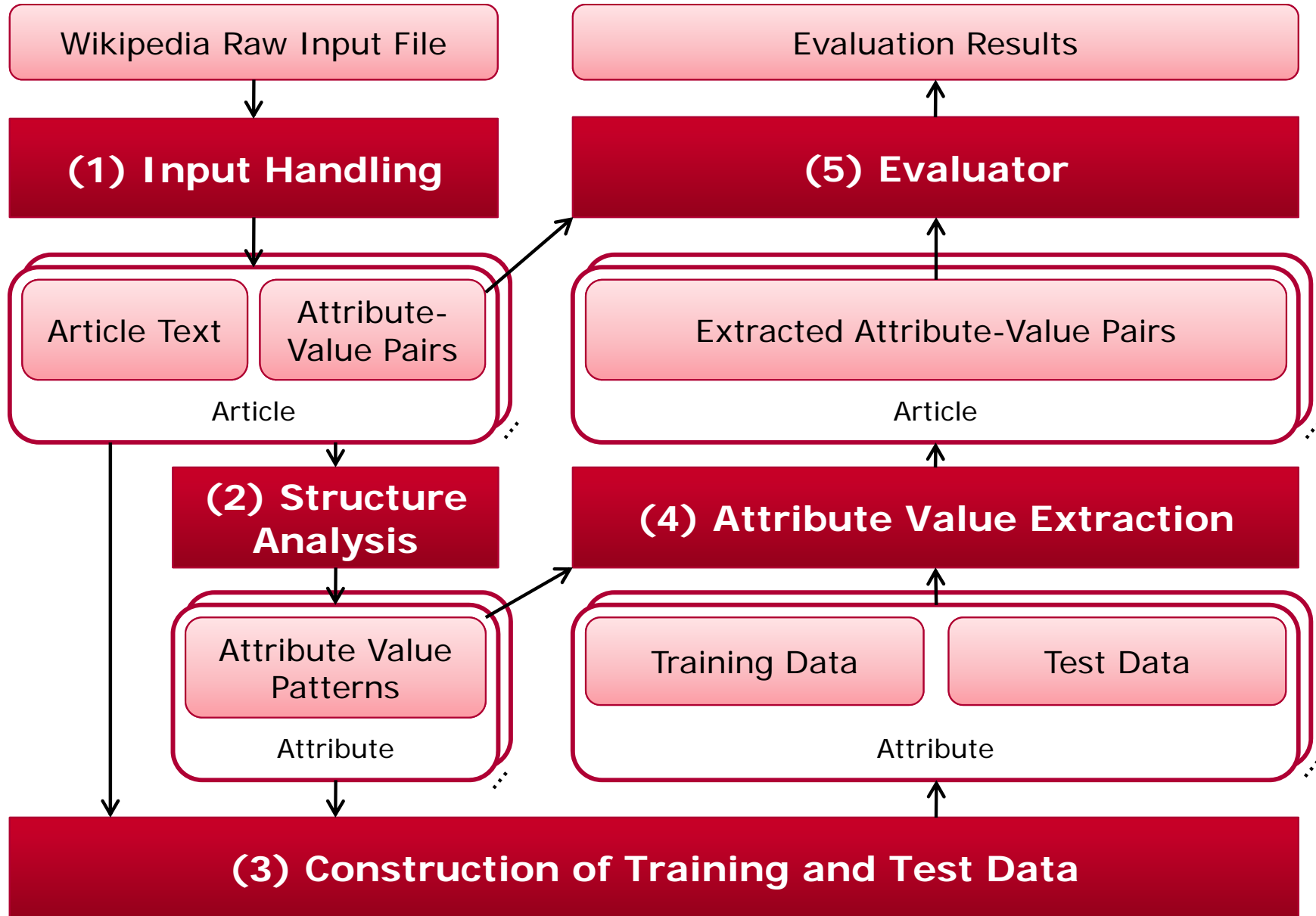## 12 most frequent attributes in infobox_book



**72.0 %** of the book articles specifying a `series` in the infobox also contain the `series` in the article text. **8.7 %** of these occurrences could only be found by separately searching for parts of these values.

Attributes (x-axis): author, name, publisher, country, isbn, release_date, followed_by, preceded_by, series, pub_date, cover_artist, subject

Legend:
- Complete match (exact)
- Complete match (similar)
- Part match (similar, average)

Y-axis: Occurrence rate
X-axis: Attribute

# 20 most frequent templates



On average, **42.2 %** of the `infobox_album` attribute values can be found in the article text. **38.2** % of these occurrences could only be found by separately searching for parts of these values.

Occurrence rate

Infobox template

infobox_indian_jurisdiction, infobox_single, infobox_album, football_player_infobox, infobox_football_biography, infobox_actor, infobox_radio_station, infobox_military_person, infobox_musical_artist_, infobox_musical_artist, infobox_settlement, infobox_television, infobox_australian_place, infobox_company, infobox_film, infobox_cityit, infobox_mlb_player, infobox_planet, infobox_nrhp, football_club_infobox

# Architecture of iPopulator

This is a presentation slide with a title, charts, and a Wikipedia excerpt. It's largely image-dominant but has text elements.

# Evaluation: infobox_planet
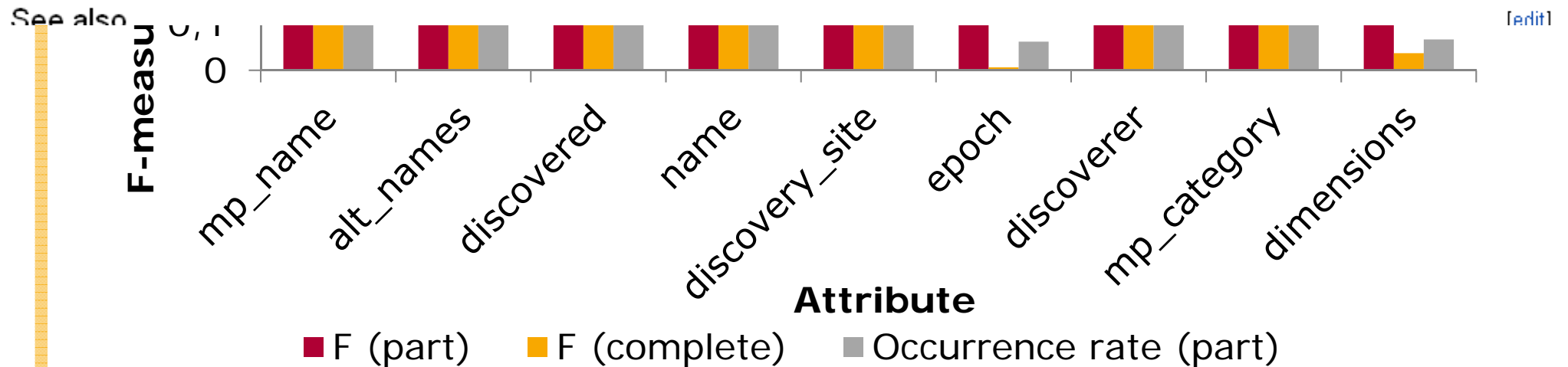
HPI Hasso Plattner Institut
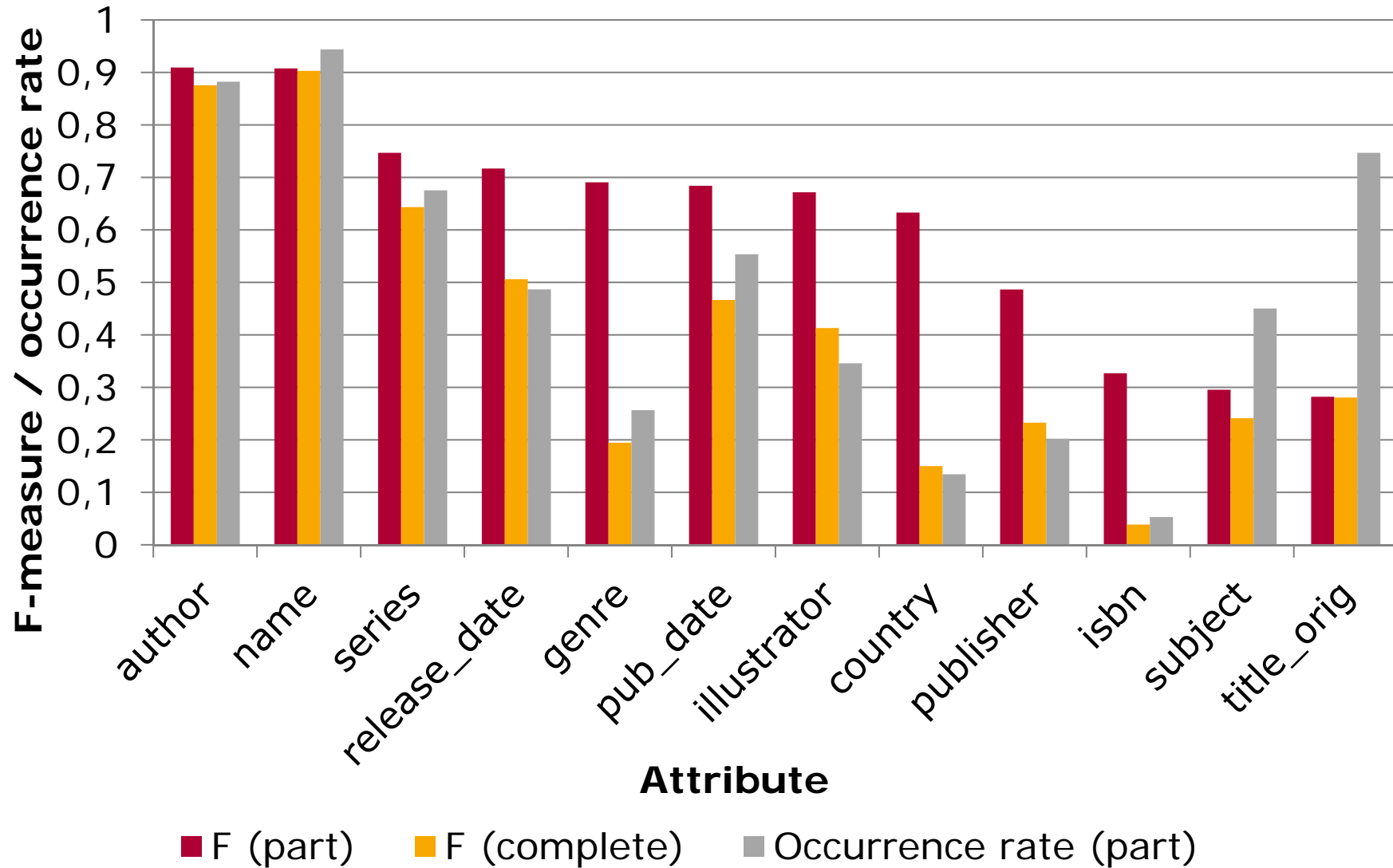
15

**ce rate**

1
0,9
0,8
0,7

## 22032 Mikekoop

From Wikipedia, the free encyclopedia

**22032 Mikekoop** (provisional designation: **1999 $XB_{151}$**) 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]

**F-measu** 0,1

0

mp_name  alt_names  discovered  name  discovery_site  epoch  discoverer  mp_category  dimensions

**Attribute**

■ F (part)   ■ F (complete)   ■ Occurrence rate (part)

# Evaluation: infobox_planet

HPI Hasso Plattner Institut

**ce rate**

1
0,9
0,8
0,7

## 22032 Mikekoop

From Wikipedia, the free encyclopedia

**22032 Mikekoop** (provisional designation: **1999 $XB_{151}$**) 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]

**F-measu** 0,1

0

mp_name | alt_names | discovered | name | discovery_site | epoch | discoverer | mp_category | dimensions

**Attribute**

■ F (part)   ■ F (complete)   ■ Occurrence rate (part)

# Evaluation: infobox_book



**F-measure / occurrence rate** (y-axis: 0, 0,1, 0,2, 0,3, 0,4, 0,5, 0,6, 0,7, 0,8, 0,9, 1)

**Attribute** (x-axis: author, name, series, release_date, genre, pub_date, illustrator, country, publisher, isbn, subject, title_orig)

Legend: ■ F (part)  ■ F (complete)  ■ Occurrence rate (part)

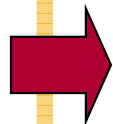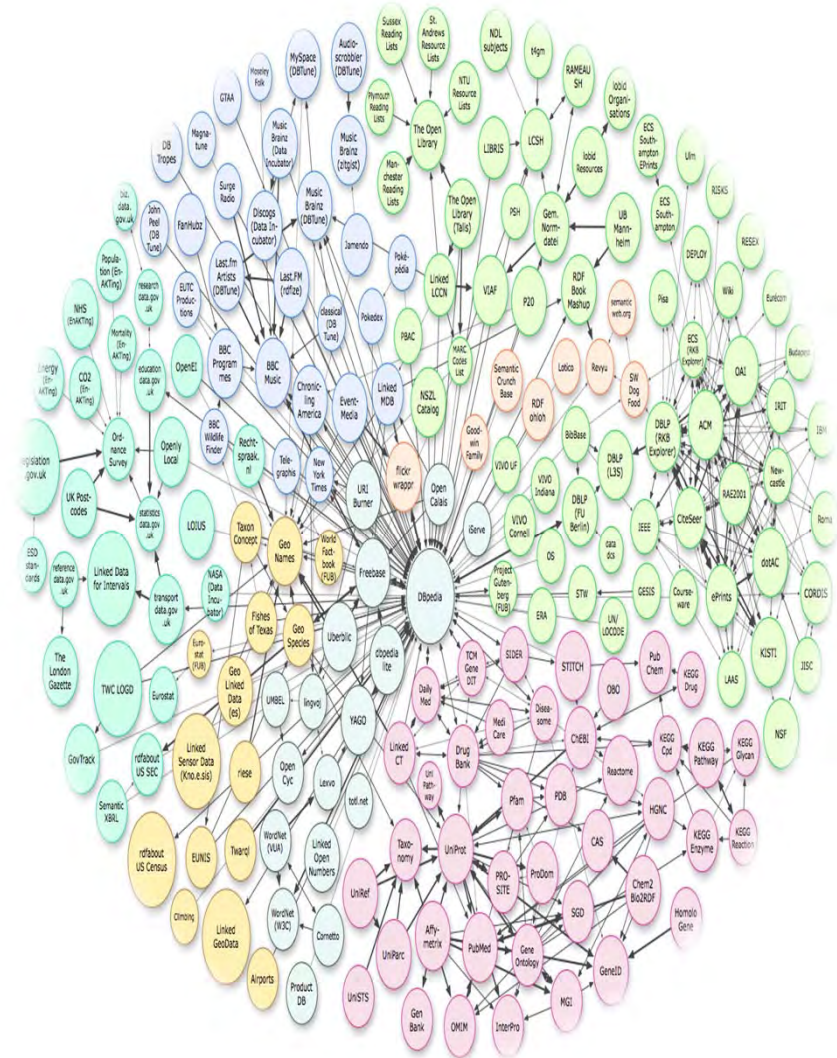# Evaluation on all attributes (>4000) of all infobox templates (>800)



http://www.hpi.uni-potsdam.de/naumann/projekte/completed_projects/ipopulator.html

# Overview

- 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

- Source
  - Formats ⟷ File converters
  - Domain ⟷ Clustering, topic mining
  - Bandwidth ⟷ Patience
- Schema
  - Structure ⟷ Schema Mapping
  - Semantics ⟷ Domain knowledge
- Data
  - Formatting ⟷ Scrubbing
  - Duplicates ⟷ Entity Matching

# 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*REPUBLIQUE TCHEQUECZECH REPUBLIC | | |
| Address: | | Country / Territory: | Czech Republic |

# The problem – a domain mess

- 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

**1.companyName**

135

438

12

22305

4739

7511

**2.company class**

448

**3.company template**

# Company Template

```
{{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 |
|---|---|
| `{{Infobox Company` | |
| `| name             =` | REQUIRED |
| `| logo             =` | |
| `| type             =` | REQUIRED |
| `| genre            =` | |
| `| fate             =` | |
| `| predecessor      =` | |
| `| successor        =` | |
| `| foundation       =` | REQUIRED |
| `| founder          =` | |
| `| defunct          =` | |
| `| location_city    =` | REQUIRED |
| `| location_country =` | REQUIRED |
| `| 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             =` | |
| `}}` | |

# The problem – a schema mess

- 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`

# The **foundation** attribute

# Infoboxes in Company class

25

- **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
- industry   28720
- foundation   26875
- products   26486
- homepage   25982
- location   24094
- companyName   23297
- companyType   19591
- companyLogo   14644
- numEmployees   11395
- locationCity   9210
- name   8700
- locationCountry   7985
- founder   7867
- revenue   7391
- parent   6468
- type   6358
- areaServed   5842
- logo   5434
- founded   4107
- companySlogan   4053
- netIncome   3528
- genre   3369
- subsid   3288

- headquarters   3191
- airline   2686
- services   2568
- callsign   2391
- icao   2386
- iata   2363
- owner   2303
- fleetSize   2246
- operatingIncome   2246
- hubs   2244
- website   2104
- intl   1996
- defunct   1987
- fate   1944
- slogan   1807
- country   1734
- destinations   1712
- assets   1591
- url   1505
- locations   1384
- divisions   1227
- logoSize   1217
- successor   1211
- distributor   1125

| 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 Econom |
| 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 | | | | | | | | |

Phew!

1pt font!

# The problem – a data mess

- Poor schemata: No types, no constraints
- Sloppy data entry:
  - Data value are neither standardized nor normalized
- `Revenue` attribute 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
  - And lots of other stuff:

?

Wal-Mart

Undisclosed

Assets exceed £4 billion GBP

http://www.credit-suisse.com/investors/en/reports/2007_results_q4.jsp

Image:green_up.png ⟶ △

€ bn (as of 2004)
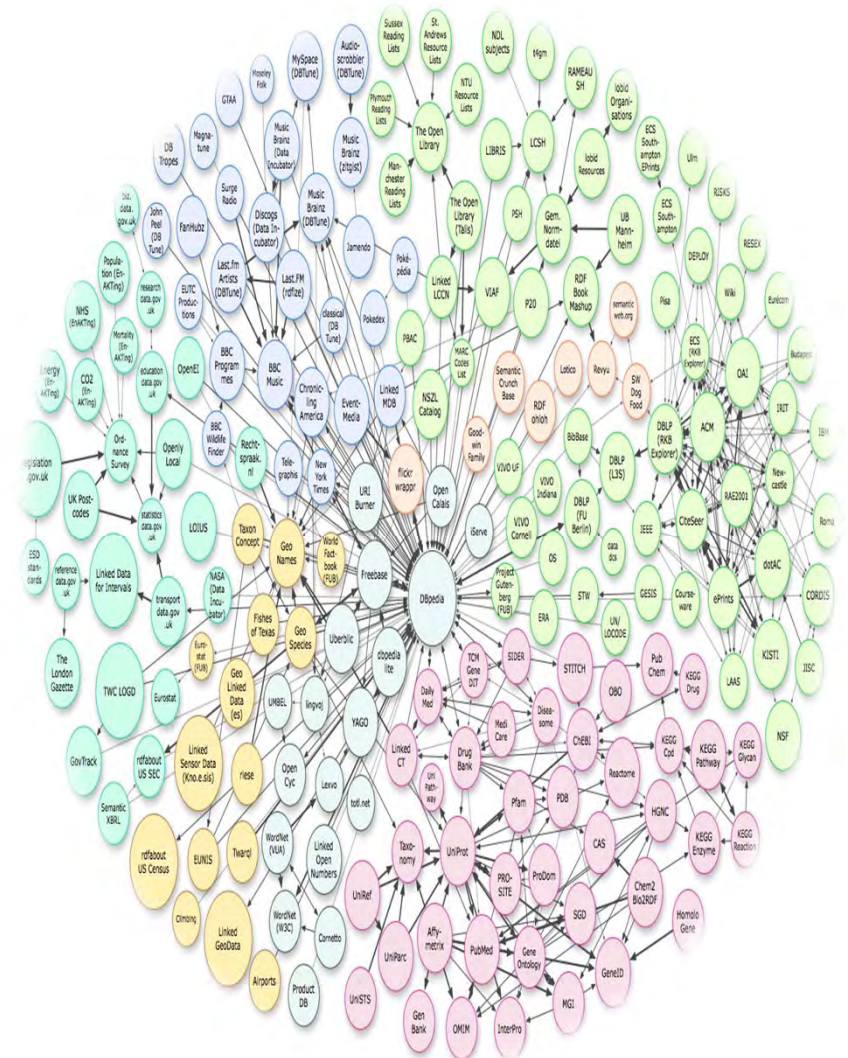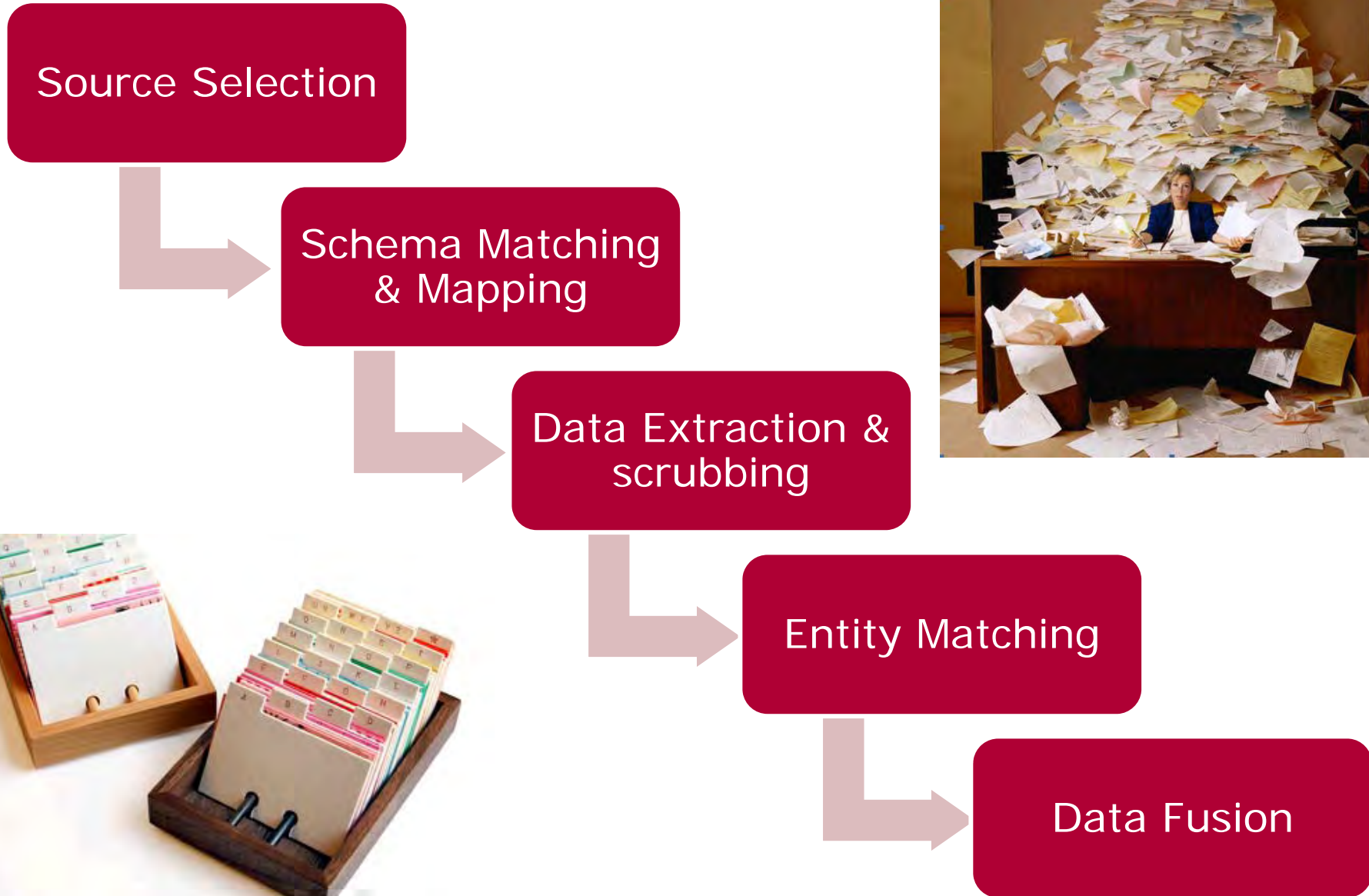
# Overview

- 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

# ProLOD profiling tasks

- **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

# Overview

- **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

# Five steps for integration

Source Selection

Schema Matching & Mapping

Data Extraction & scrubbing

Entity Matching

Data Fusion

# Step 1: Source selection

- 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

# Step 2: Schema matching and mapping

- **Semi-automated matching**
  - □ Label-based and instance-based
- **Challenges:**
  - □ Multi-lingual
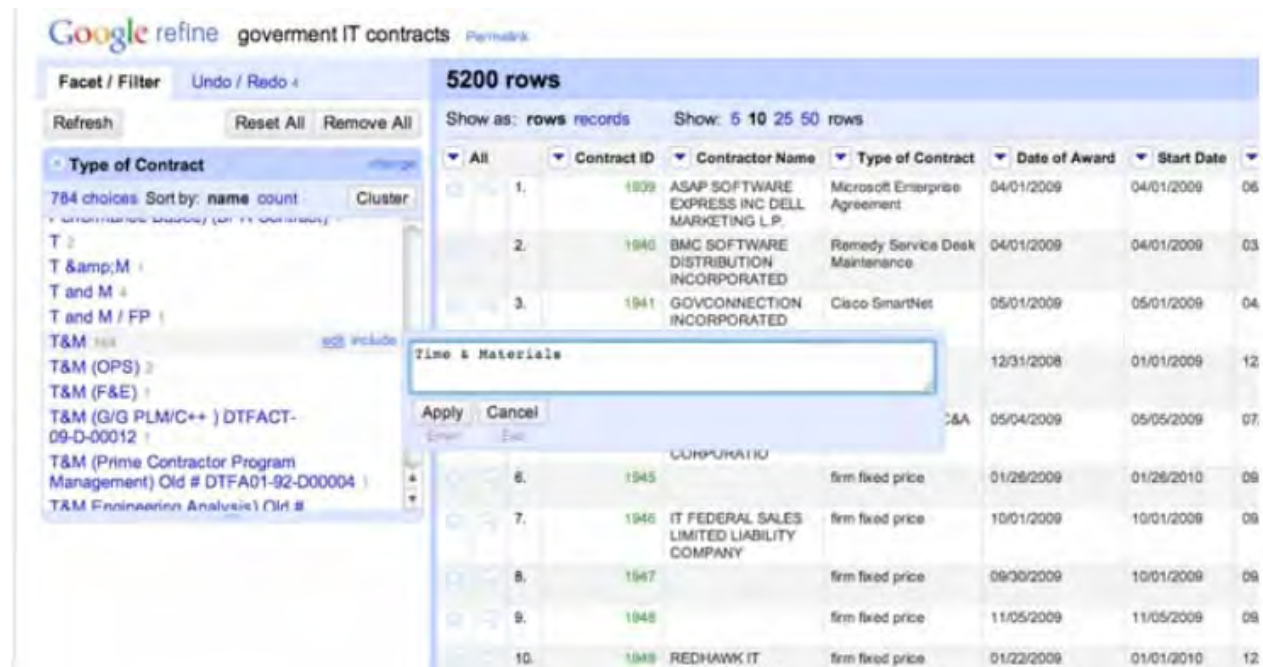  - □ Homonyms and Synonyms
  - □ 1:1, 1:n, n:m
- **Complex data transformation**

| Final Schema | DBPedia | SEC | Freebase |
|---|---|---|---|
| dbpediaURI | | | /type/object/key |
| cik | secCik | CIK | |
| irsnumber | | | |
| companyName | companyName, name, nonProfitName | name | /type/object/name, /common/ |
| address | | BusinessAddress, MailingAddress | /location/mailing_address/stre /location/mailing_address/pos |
| locationCity | locationCity, location | BusinessAddress, MailingAddress | /location/mailing_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/employees /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 | | |

- Recognize data types

- Regular expressions for multi-valued strings

- Remove spurious values (layout, formatting, …)

- Standardize formats
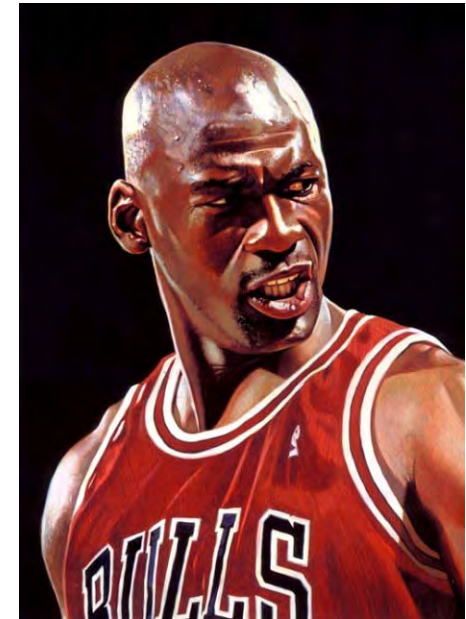
- Translate from foreign languages

- Many tools

# Step 4: Entity matching

- Duplicate entities

- Linking between entities

- Challenges

  □ Fuzzy matching: Similarity measures

  □ Data volume: Partitioning algorithms

  □ Sparse data

    ◇ **`Michael Jordan born_in Miami`**

**Find People**

First Name | *Last Name | City, State or ZIP | Find People
Michael | Jordan | CA

**Whoa!** Over 100 Results Found

Michael Jordan (disambiguation)

From Wikipedia, the free encyclopedia

**Michael Jordan** is an American basketball player.

**Michael Jordan** may also refer to:

- Michael Jordan (mycologist), English mycologist
- Michael Jordan (footballer) (born 1986), English goalkeeper (Ars
- Michael B. Jordan (born 1987), American actor
- Michael I. Jordan (born 1957), American researcher in machine
- Michael H. Jordan (d. 2010), American executive for CBS, Peps
- Michael-Hakim Jordan (born 1977), American professional bask
- Michael Jordan (Irish politician), Irish Farmers' Party TD from W

# Step 5: Data fusion

- Combine multiple representations of real-world entities

  □ Survivorship, consolidation, etc.

- Resolve data conflicts

  □ Conflict resolution functions
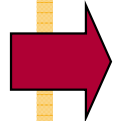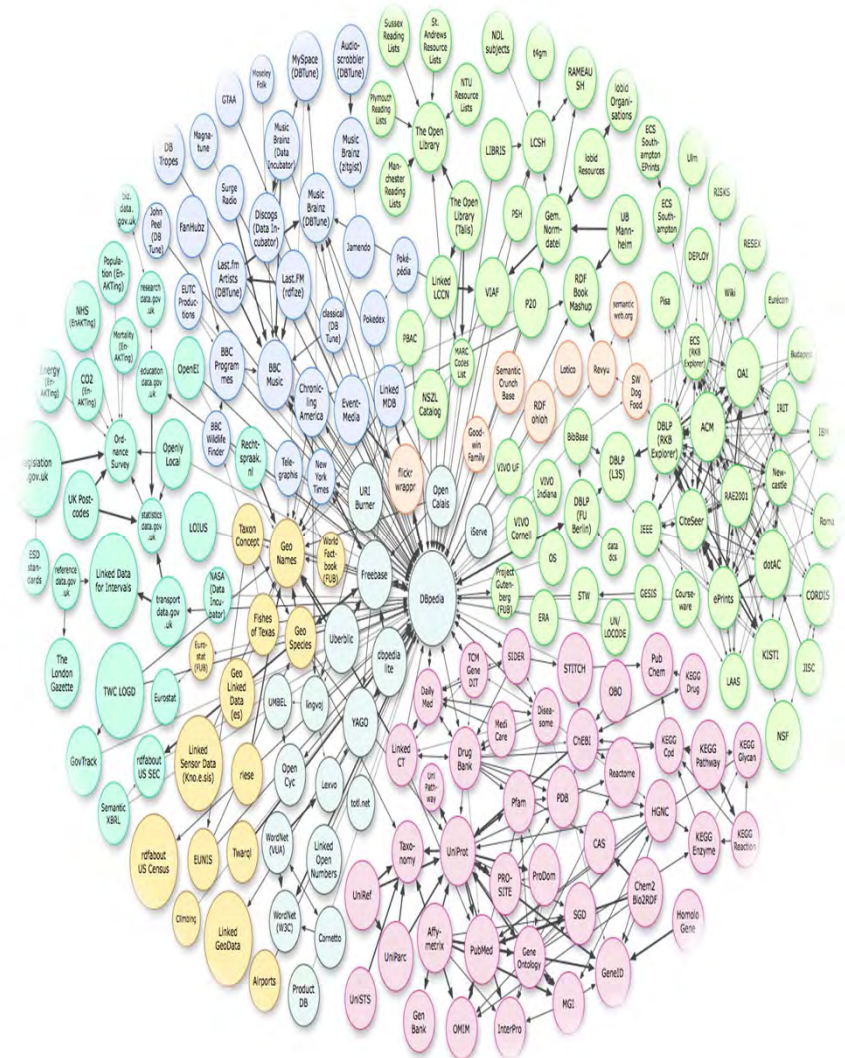
  □ Reputation / accuracy / freshness -> "truth discovery"

| 0766607194 | H. Melville | | $3.98 | 📄 |
|---|---|---|---|---|

ID       max length       MIN   CONCAT

| 0766607194 | Herman Melville | Moby Dick | $5.99 | 📄 📄 |
|---|---|---|---|---|

- Retain data lineage

# Overview

- **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

# Multi-Lingual Wikipedia

- **Goal:** Schema matching across languages
    - Complement infobox data
    - Autocomplete for authors
    - Detect errors or inconsistencies
    - Keep values up to date
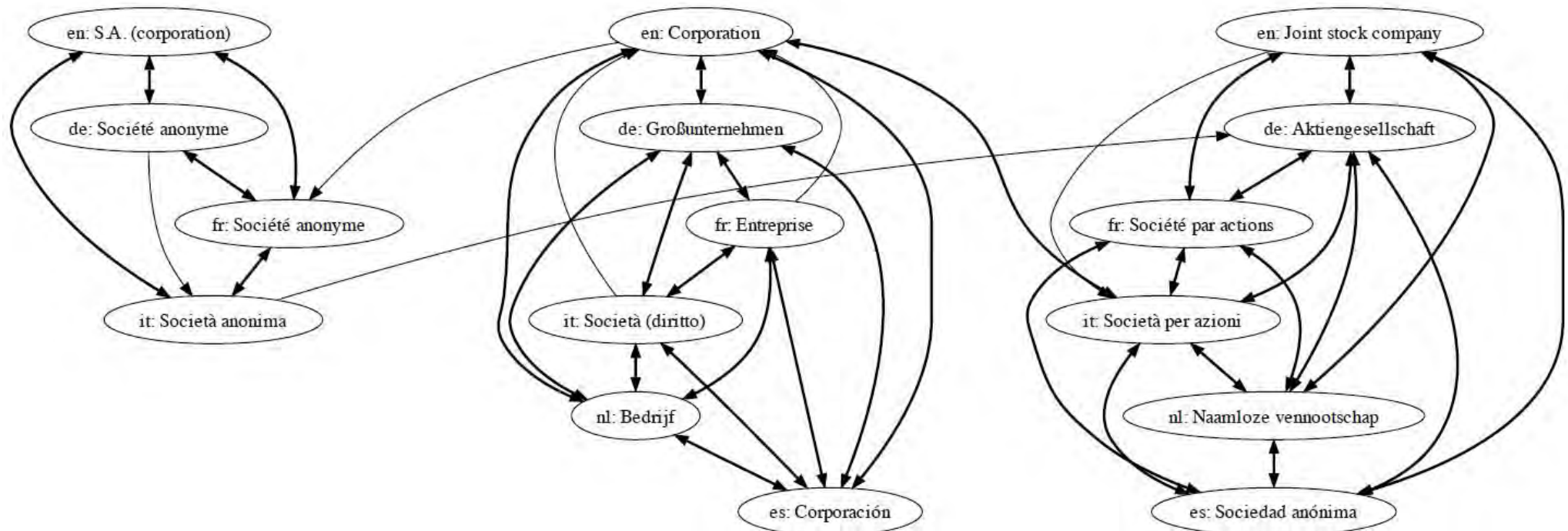- **Idea:** Use cross-language links across all 285 languages
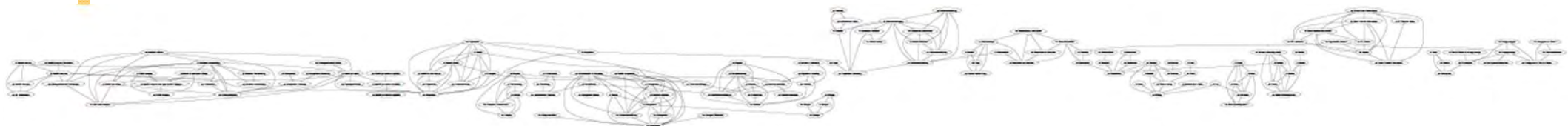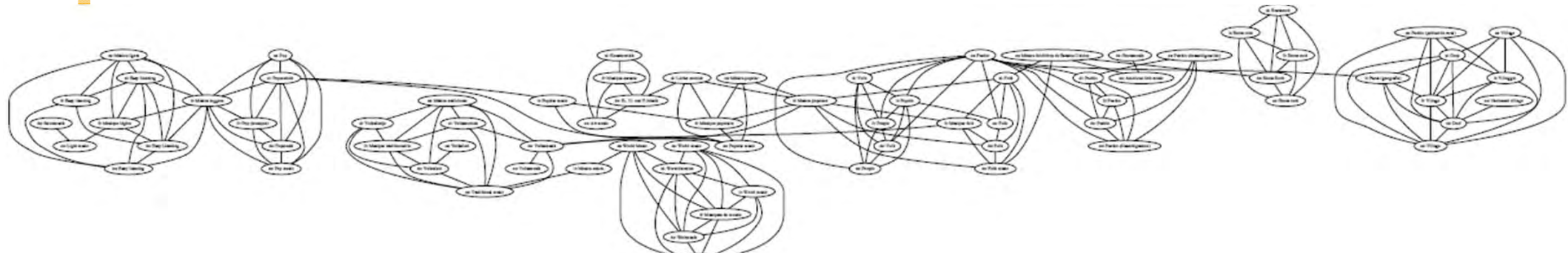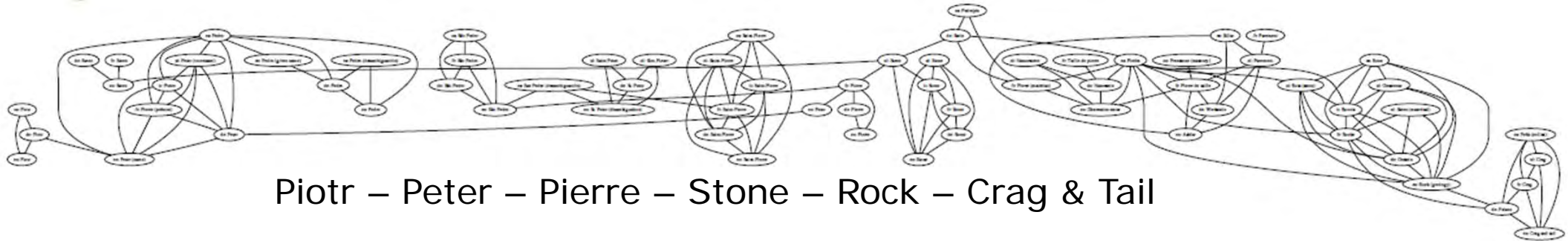
# Interlanguage links (ILLs)

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

# Other large components

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

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



- **SCC**
  - Strongly connected components (SCC)
  - Each node is reachable from each other node
  - 1,067,753 SCCs of which 3,469 are incoherent

- **BCC**
  - Bidirectionally connected components (BCC)
  - Undirected graph of bidirectional components is connected
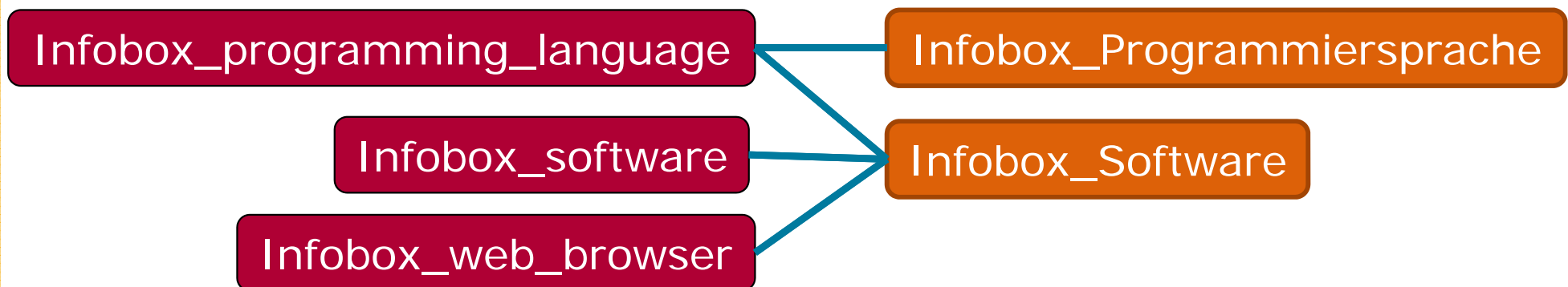  - 4,241 BCCs of which 2,980 are incoherent

- **2CC**
  - Bi-connected components (2CC)
  - Each pair of vertices is connected via **two** vertex-independent paths.
  - 8,828 2CCs of which 4,770 are vertex-disjoint

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

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

| Infobox_programming_language | Infobox_Programmiersprache |
|---|---|
| Infobox_software | Infobox_Software |
| Infobox_web_browser | |

# Duplicate-based Schema Matching

- 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

47

| Coordinates: | 52°30'2"N 13°23'56"E |
|---|---|
| Country | Germany |
| Government | |
| - Governing Mayor | Klaus Wowereit (SPD) |
| - Governing parties | SPD / Die Linke |
| - Votes in Bundesrat | 4 (of 69) |
| Area | |
| - City | 891.85 km$^2$ (344.3 sq mi) |
| Elevation | 34 - 115 m (-343 ft) |
| Population (31 March 2010)[1] | |
| - City | 3,440,441 |
| - Density | 3,857.6/km$^2$ (9,991.3/sq mi) |
| - Metro | 4,429,847 |
| Time zone | CET (UTC+1) |
| - Summer (DST) | CEST (UTC+2) |
| Postal code(s) | 10001–14199 |
| Area code(s) | 030 |
| ISO 3166 code | DE-BE |
| Vehicle registration | B |
| GDP / Nominal | € 90.1[2] billion (2009) [citation needed] |
| NUTS Region | DE3 |
| Website | berlin.de |

| Basisdaten | |
|---|---|
| Fläche: | 891,85 km$^2$ (14.) |
| Einwohner: | 3.456.264[1] (8.) (31. Oktober 2010) |
| Bevölkerungsdichte: | 3.875 Einw. je km$^2$ (1.) als Bundesland, (2.) als Gemeinde |
| BIP: | 90,1 Mrd. € (2009) |
| Höhe: | 34–115 m ü. NN |
| Geografische Lage: | 52° 31' N, 13° 24' O |
| Zeitzone: | Mitteleuropäische Zeit (MEZ) UTC+1 |
| Postleitzahlen: | 10115–14199 |
| Vorwahl: | 030 |
| Kfz-Kennzeichen: | B |
| Gemeindeschlüssel: | 11 0 00 000 |
| ISO 3166-2: | DE-BE |
| UN/LOCODE: | DE BER |
| Website: | www.berlin.de |
| Politik | |
| Reg. Bürgermeister: | Klaus Wowereit (SPD) |
| Reg. Parteien: | SPD und Die Linke |
| Sitzverteilung im Abgeordnetenhaus | SPD 54 CDU 36 |

# Evaluation

- 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 |
| $F_1$ Score | 93.06 | 94.50 | 94.97 | 91.42 | 92.44 | 93.33 | 93.42 |

# Next step by community: Wikidata

- Free knowledge base about the world

- Read and edited by humans and machines

- Data in all the languages of the Wikimedia projects

  □ In particular: Wikipedia pages

- Central access to data
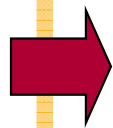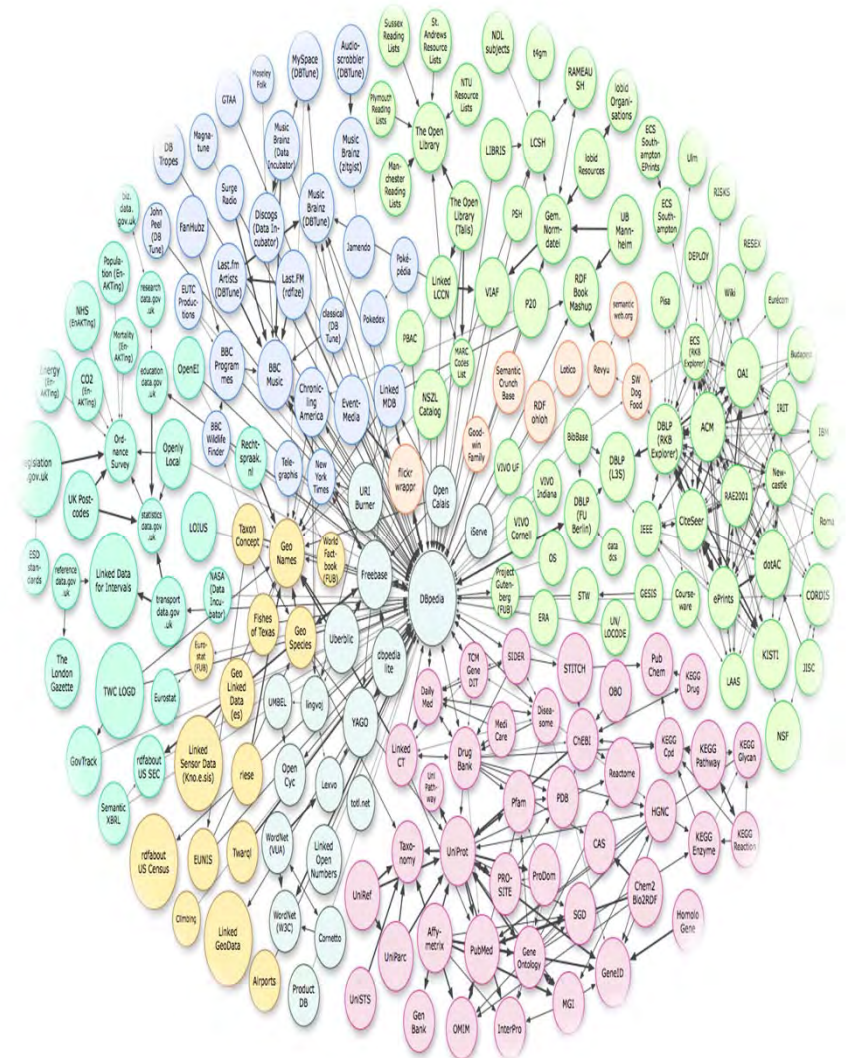

- Begin April 2012 – much to do
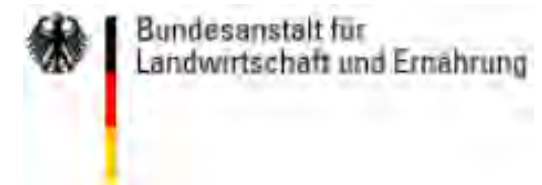
- http://meta.wikimedia.org/wiki/Wikidata/de
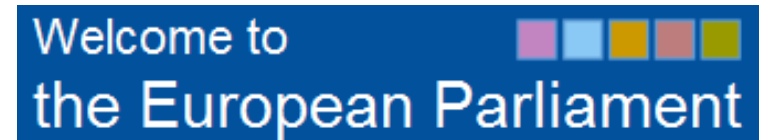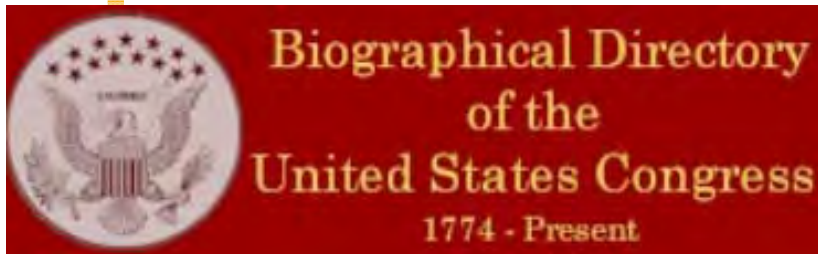
50

- 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**
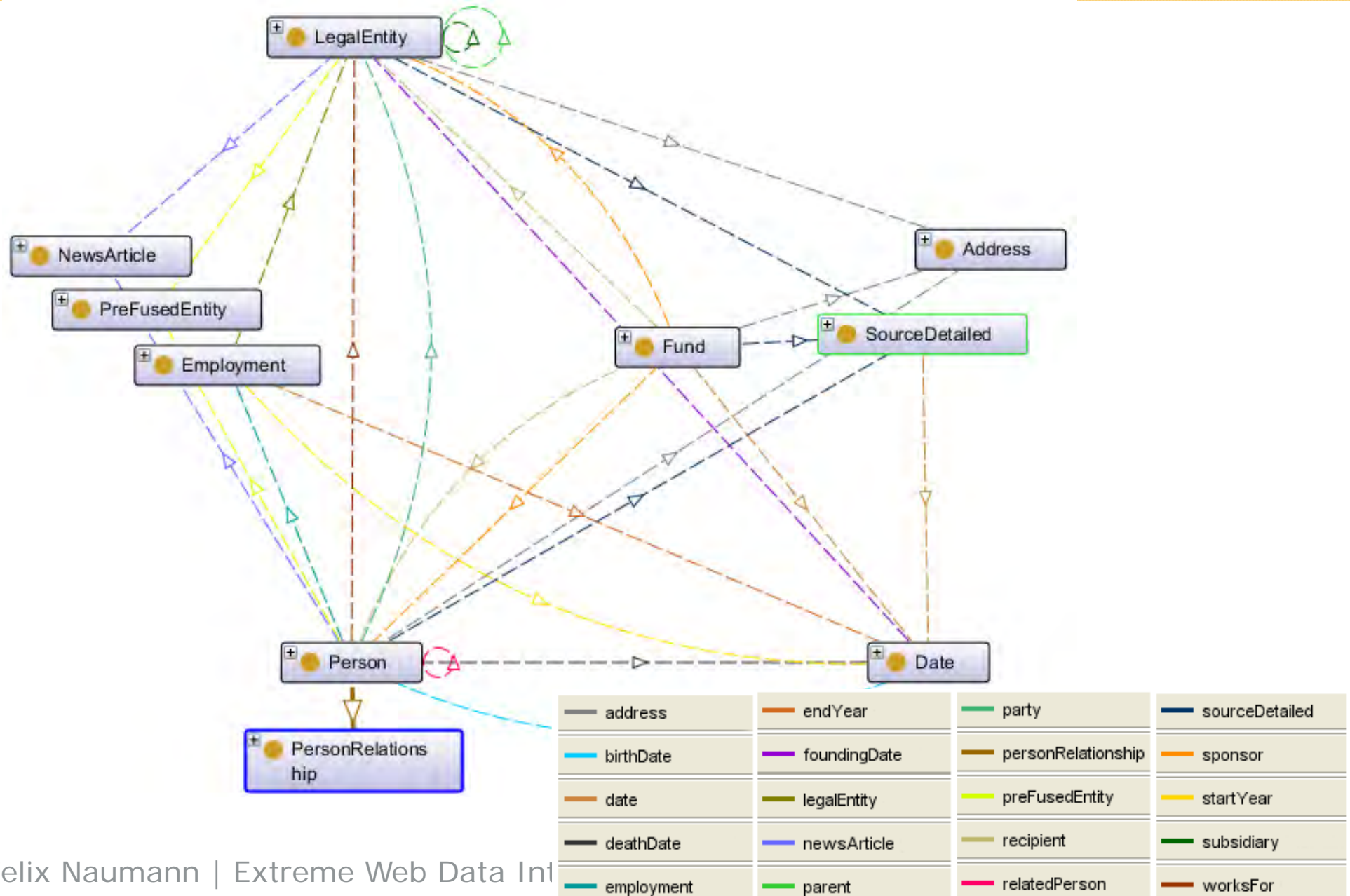  - Politicians, friends, and funds
  - The GovWILD experience

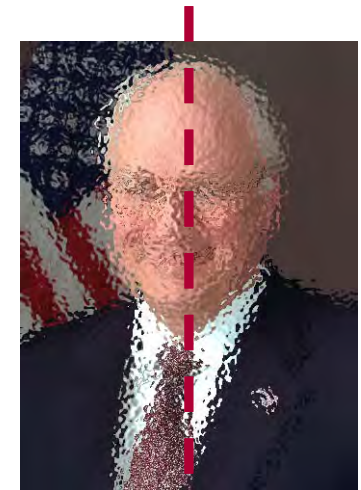# Motivation – Wealth of Open Gov Data

# Companies, Agencies, and People
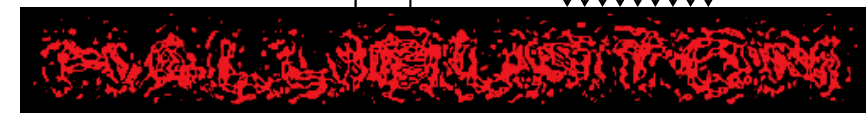
# Interesting queries

- 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
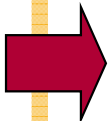  - *Sets*: "for each … find all …"
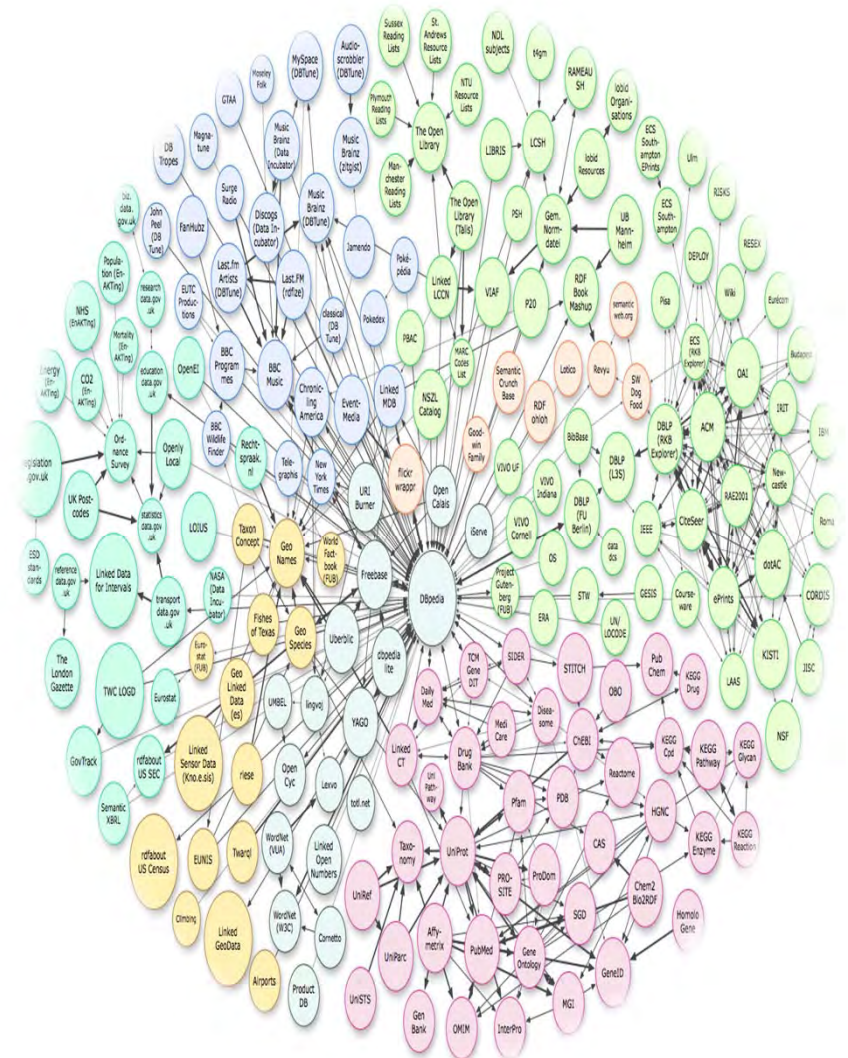
Chairman
of the board

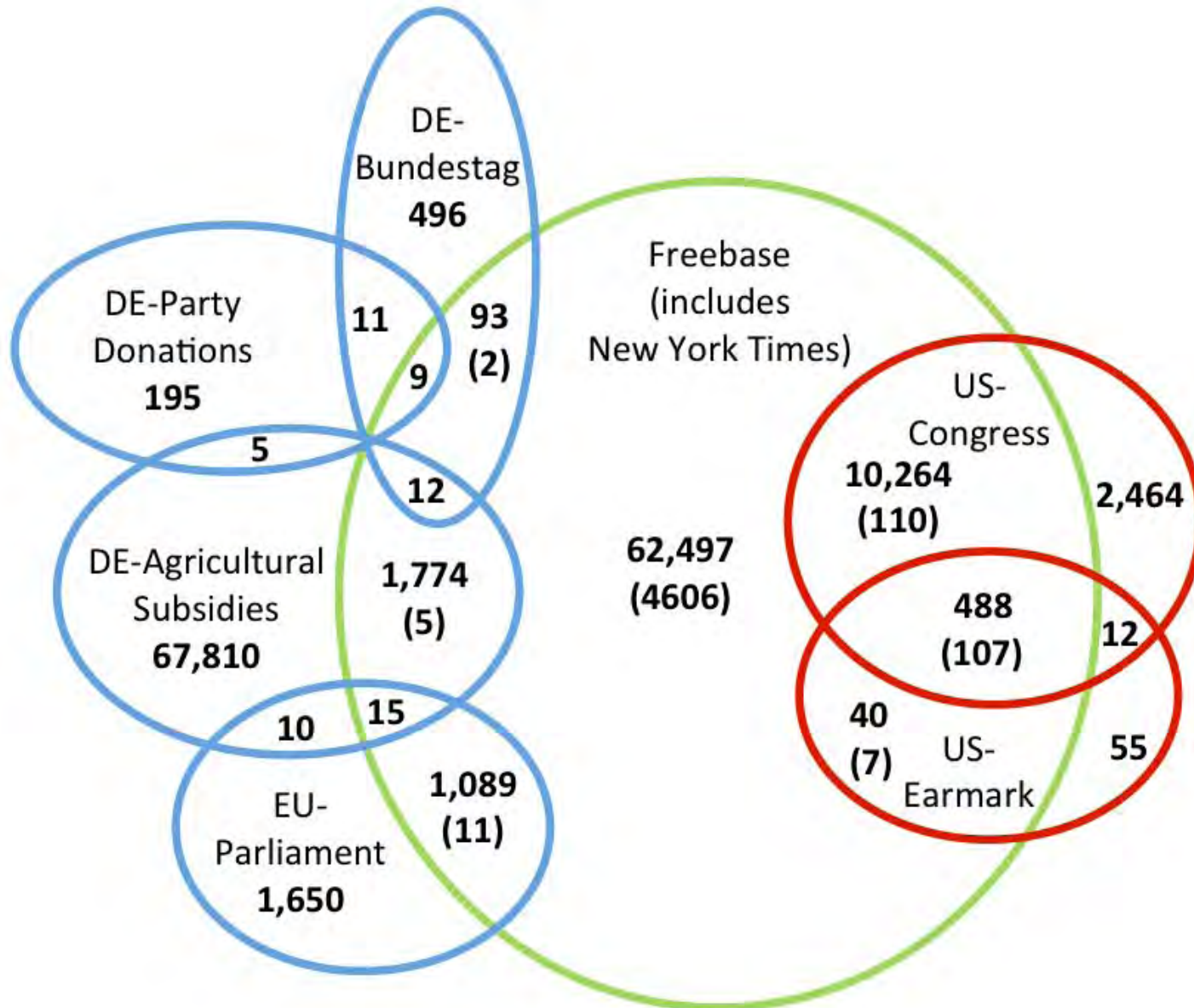Funds

CEO

# Overview

- **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

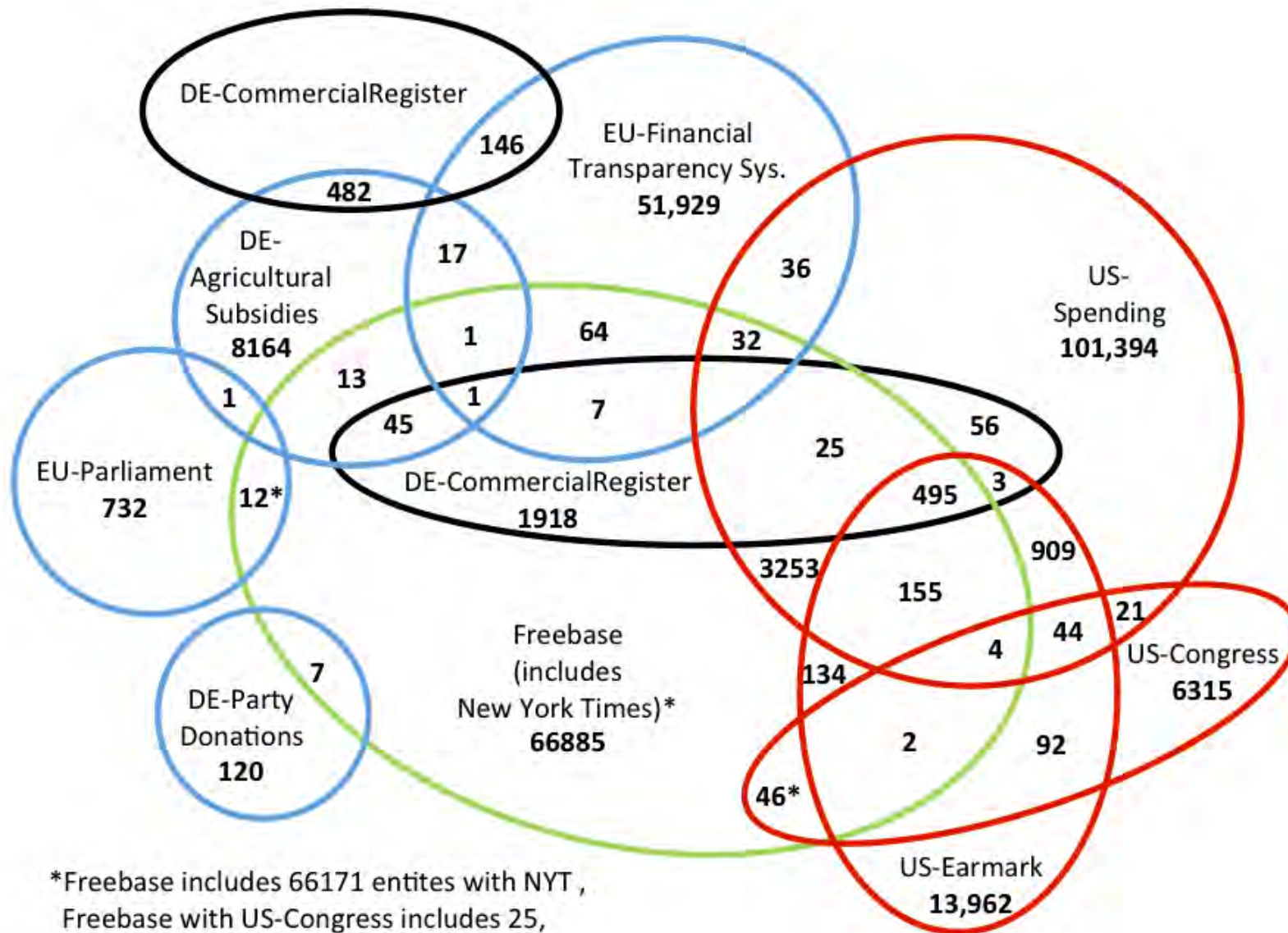| Step | Time | | Input size on master node and element count | | | Details |
|---|---|---|---|---|---|---|
| | Jaql 0.4 | Jaql 0.5 | LegalEntity.json | Person.json | Fund.jon | |
| **Scrubbing** | | | | | | |
| Scrub and map 15 files | 1h 15min | 1h 15min | Start with 11 GB size | | | - map and normalize attributes<br>- set references within a source (includes many joins)<br>- group entities / match entities of the same source<br>- use dictionaries for enrichment |
| Merging Scrubbed Files | 3 min | 3 min | | | | - concatenate files in HDFS to achieve 3 files containing persons, legal entities, and funds |
| | | | 162 MB - 217 087 entities | 544 MB - 1 357 810 entities | 471 MB - 998 150 ent. | |
| **Matching of LegalEntity** | | | | | | |
| Write from HDFS to master | 6 min | 7 sec | -||- | | | |
| Find similar entities on workstation | 30 min | 24 min | -||- | | | - computes duplicates in pairs of 2, non-parallel |
| Write back to HDFS | 7 sec | 6 sec | 44 MB – 7530 pairs | | | |
| Fuse similar objects | 10 min | 10 min | -||- | | | - compute transitive closure of IDs (transform and combine with UDF) |
| | | | Join 5 402 fused IDs with 7530 Leg.E. | | | - join clustered IDs with objects (2 minutes)<br>- group by cluster_ID |
| | | | -||- | | | - split large clusters (transform with UDF) |
| | | | -||- | | | - fuse these clusters (transform with UDF) |
| Update fused Ids in all files (merge new IDs from Legal Entity into Person, Fund and LegalEntity) | 10 min | 10 min | 211 362 entities | 1 357 810 entities | 998 150 ent. | - transform on source file to find all ID changes<br>- transform on target file to find all possibly old references<br>- join both<br>- group by target ID<br>- join this with target file (3 min for merging from LegalEntity to Person)<br>- transform this to set new IDs |
| **Matching of Person** | | | | | | |
| Write from HDFS to master | 18 min | 20 sec | | 544 MB | | |
| Find similar entities on workstation | 44 min | 48 min | | -||- | | - as above, non-parallel |
| Write back to HDFS | 8 sec | 12 sec | | 79 MB – 51 634 pairs | | |
| Fuse similar objects | 11 min | 10 min | | Join 35 744 fused with all Persons | | - as above |
| Remove irrelevant Freebase Persons | 1 min | 1 min | | filter 328 889 out of 1 323 112 | | - remove freebase persons without references (filter) |
| Update fused Ids in all files | 10 min | 9 min | 211 362 entities | 328 889 entities | 998 150 ent. | - as above, from Person file to all others |
| **Finalize data** | | | | | | |
| Precanned Query for US states | 9 min | 10 min | -||- | -||- | -||- | - for every object create stateEntities array with connected state names (transform on LegalEntity, Person, Fund)<br>- filter US states from legal entities to create US states file<br>- replace state names with state IDs (similar to updating IDs before) |
| Clean up attributes | 2 min | 1.5 min | -||- | -||- | -||- | - remove empty arrays |
| Write JSON from HDFS to master | 40 min | 1min | 175 MB | 428 MB | 524 MB | |
| **Prepare for RDF export** | | | | | | |
| Add attributes | 1 min | 2 min | -||- | -||- | -||- | - add „label" and „uri" fields (transform with UDF) |
| Replace ID references by URI references | 23 min | 19 min | -||- | -||- | -||- | - as update IDs above, for most combinations of LegalEntity, Person, and Fund (Funds are never referenced) |
| Write from HDFS to master | 46 min | 1 min | 185 MB - 211 362 entities | 453 MB - 328 889 entities | 689 MB - 998 150 ent. | |
| | sum: 5h 39 min | Sum: 3h 45min | | | | |

*Freebase includes 66171 entites with NYT ,
Freebase with US-Congress includes 25,
Freebase with EU-Parliament includes 0.

# http://govwild.org

- 150,000 persons
- 270,000 legal entities
- 1,100,000 funds
- 43,000,000 triples

- Keyword Queries

- Linked Data Interface (dereference URIs)

- Exploration of entities mentioned in New York Times articles

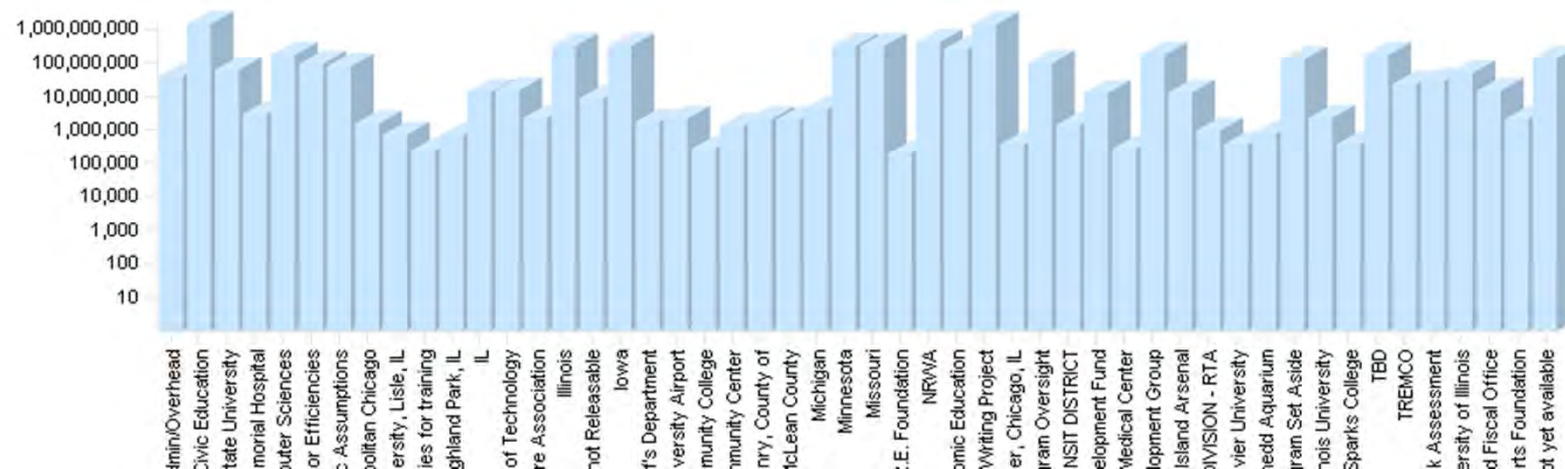- Data Download (RDF, SQL Dump, JSON files)

## 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 Fe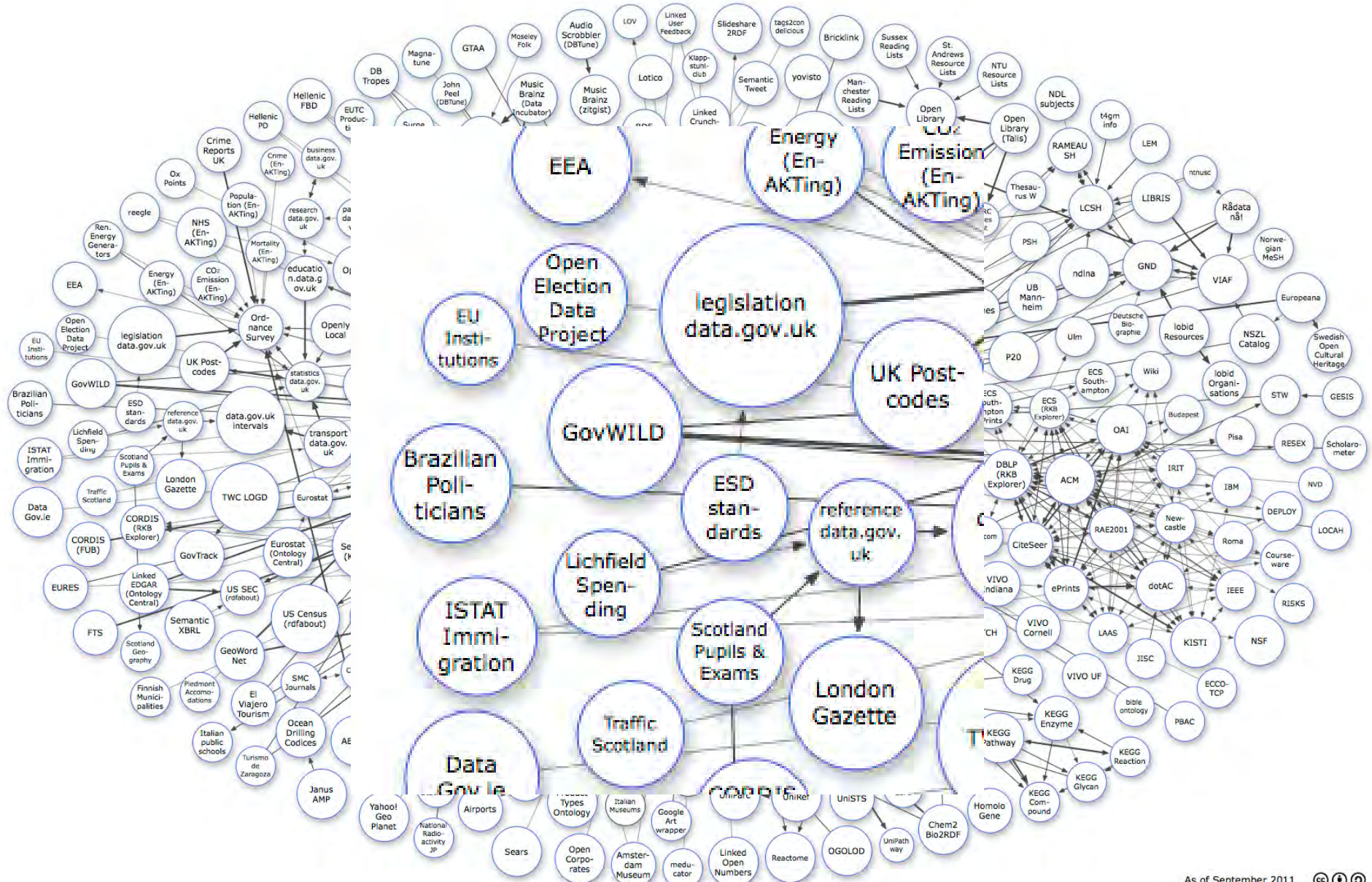bruary 2007, and after a close campaign in the 2008 Democratic Party presidential primaries against Hillary Rodham Clinton, he won his party's nomination. In the 2008 general election, he defeated Republican nominee John McCain and was inaugurated as president on January 20, 2009.4

## Earmarks

# Allianz insurance
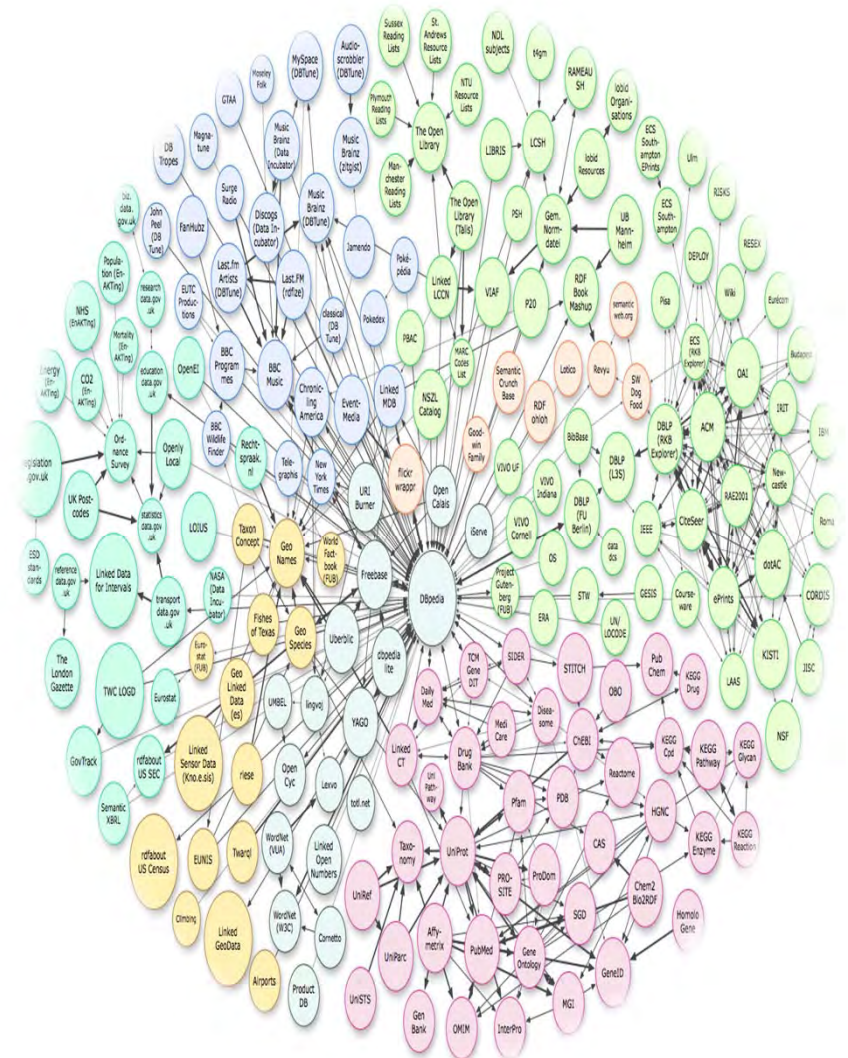
# Summary

- **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

- [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