# Web of Data : Past Research, Current State, Clear Opportunities
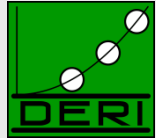
Giovanni Tummarello, DERI Institute (Galway)
FBK Institute (Trento)

National University of Ireland, Galway
*Ollscoil na hÉireann, Gaillimh*

science foundation ireland
fondúireacht eolaíochta éireann

# The Web of Data

A Web where **pages** have **elements** that machines can **interpret automatically**

Quite a significant slice of the web, already thanks to **Metadata Standards**

☐ RDF, RDFa

☐ Microformats

☐ Schema.org (Microdata)

# For example:



```html
<h1 id="name">
<span class="fn n">
        <span class="given-name">Giovanni
        </span>
        <span class="family-
name">Tummarello
```

Enabling **networked** knowledge.

- Google Rich Snippets Program

Lenovo ThinkPad **X201** 3626 Review - Watch CNET's Video Review
★★★★☆ Review by Dan Ackerman - Apr 30, 2010 - Price range: $1,029.00
1 Feb 2010 ... Anyone looking for the power of a midsize laptop in a compact 12-inch body
has only a few choices, and none to date tops **Lenovo's** excellent ...
reviews.cnet.com/.../**lenovo**...**x201**.../4505-3121_7-33998451.html - Cached - Similar

- Facebook Opengraph Protocol

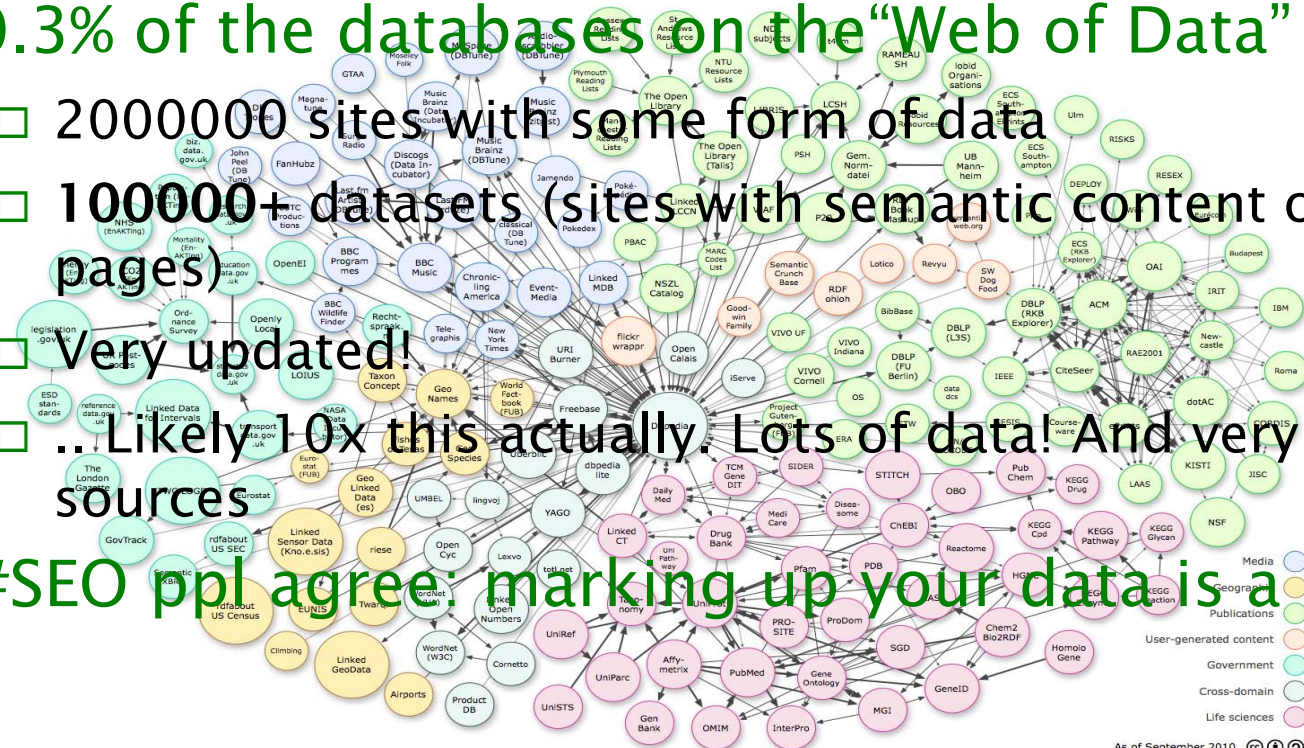👍 Like      ∫ 114,017 people like this. Be the first of your friends.
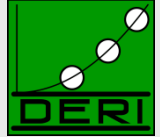
- Schema.org..

- **LOD Cloud: approx 300 datasets**
  - Relatively clean, with some mutual integration
  - Not very updated
- **0.3% of the databases on the "Web of Data"**
  - 2000000 sites with some form of data
  - **100000+** datasets (sites with semantic content on many pages)
  - Very updated!
  - ...Likely 10x this actually. Lots of data! And very important sources
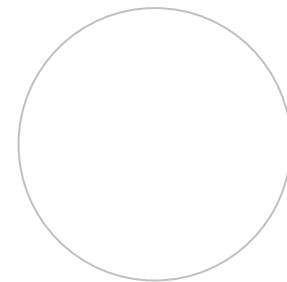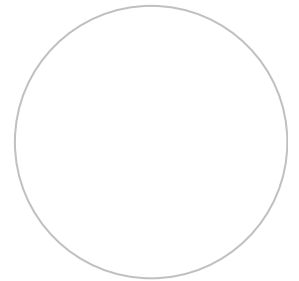- **#SEO ppl agree: marking up your data is a must**

As of September 2010

■ Questions

　□ How to find sources ?
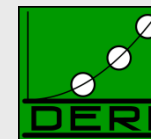
　□ How to clean and transform them for one's specific needs?

Sindice "WebStar" Cluster

Enabling **networked** knowledge.

# Sindice.com: Portal for Sindice Services

# Building a graph of web graphs

THE SEMANTIC WEB INDEX

Home    About    Search    Submit    Forum    Dev

Search interface type: **Simple** Advanced Guru    Query Language Documentation

keyword(s) stefan decker

**SEARCH**    Group By Dataset: ☑

**Quick filters** (All options)

Time range:
*Any date*
Today Yesterday Last week
Last month Last year

Format:
*Any format*
RDF RDFA MICRODATA
MICROFORMAT XFN HCARD
HCALENDAR HLISTING
HRESUME LICENSE GEO ADR

Predicate: ⑦
[                    ]

Class: ⑦
[person OR publication]

Ontology: ⑦
[                    ]

Domain: ⑦
[                    ]

**Sindice search:stefan decker class:(person OR publication)** found 1,286 documents (in 0.04 seconds) showing first 10 datasets out of 100+

**www.bibsonomy.org** (294)
➕ See all results from: www.bibsonomy.org dataset

**dblp.l3s.de** (127)
➕ See all results from: dblp.l3s.de dataset

**dblp.rkbexplorer.com** (108)
➕ See all results from: dblp.rkbexplorer.com dataset

**www.aifb.kit.edu** (64)
➕ See all results from: www.aifb.kit.edu dataset

**oai.rkbexplorer.com** (56)
➕ See all results from: oai.rkbexplorer.com dataset

**sw.deri.org** (50)
➕ See all results from: sw.deri.org dataset

**www.sembase.at** (48)
➕ See all results from: www.sembase.at dataset

**data.semanticweb.org** (47)
➕ See all results from: data.semanticweb.org dataset

**twitter.com** (44)
➕ See all results from: twitter.com dataset

**www.deri.ie** (38)
➕ See all results from: www.deri.ie dataset

Searching on 298.04 million documents.

← Previous    **1** 2 3 4 5 6 7 8 9 10    Next →

# Analytics Example

- **Bibsonomy.org**
- **L3s.de**
- **Deri.ie**

National University of Ireland, Galway
Ollscoil na hÉireann, Gaillimh

sfi
science foundation ireland
fondúireacht eolaíochta éireann

Enabling **networked** knowledge.

# Building on Sindice (2): Site Services

Enabling **networked** knowledge.

Menu just demo purpose

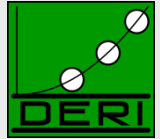The Thirteenth Floor 2012 10,000 B.C. The Day After Tomorrow Trade The Secret Life of Bees Alien Resurrection Aliens Alien 3 AVP - Alien Vs. Predator

## 10,000 B.C.

Menu just demo purpose
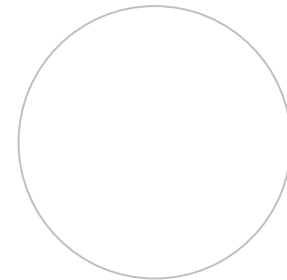
The Thirteenth Floor 2012 10,000 B.C. The Day After Tomorrow Trade The Secret Life of Bees Alien Resurrection Aliens Alien 3 AVP - Alien Vs. Predator
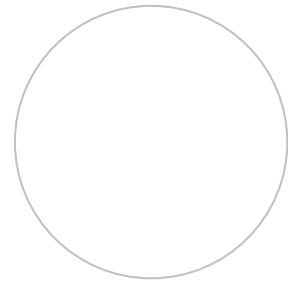
## 10,000 B.C.

👍 Like    124 people like this. Be the first of your friends.

Enabling **networked** knowledge.

# Just a small code addition! (?)

```
<meta property="og:url" content="http://demo.sindice.net/..." />
<meta property="og:title" content="10000 BC />
<meta property="og:type" content="movie" />
<meta property="og:image" content="http://demo.sindice.net/..."/>
```

# Just make it known to Sindice

## .. And keep itself in sync

# Pick a widget from our "Store"

# Copypaste Tag into your HTML

# Sindice.com/YOURSITE/movies/search

# Enjoy Services

## 10,000 B.C.

👍 Like    f 124 people like this. Be the first of your friends.

**More info:**

**Title:**
10,000 BC
**Starring:**
Steven Strait
Cliff Curtis
Camilla Belle
**Writer:**
Harald Kloser
Roland Emmerich
**Director:**
Roland Emmerich
**Budget:**
105000000
**Runtime:**
109 min.

powered by Sindice

National University of Ir
Ollscoil na Mireann.

# Experts create and $ from widgets

## Widget Composer

**Example URLs:**

From this dropdown you can choose a few example URLs
to be used for URL parameter

http://demo.sindice.net/sindice-widget/og-data/movie/2 ▼

Check that we have this URL in sindice: ☐

**SPARQL query name to be used by widget**

(pickup from list or create your own one inside the box)

get_more_info ▼

Use custom query: ☐

**Optional SPARQL query parameters:**

Enter value for: URL:string
Optional use only if you want to force widget to use different URL

http://demo.sindice.net/sindice-widget/og-data/movie/5

[ Show ]

**Required SPARQL query parameters:**

**Other optional widget parameters:**

Widget container selector (.class or #id). Possible values:
**div.foo** for <div class="foo"/>
**#foo** for <div id="foo"/>

#my_widget

Text for widget's header: (empty to hide header)

My header

Text for widget's footer: (empty to hide footer)

powered by Sindice

Widget width in pixels:

200

Widget height (without header and footer) in pixels:
(empty for auto)

**Widget styles editor:**

Use this editor to change widget styles on the fly.

**Rules for writing a custom query:(click to expand)**

**Sparql query description:**

This widget uses the film title embedded in the OpenGraph markup to look up
DBpedia for the name of the director of the movie, then looks for titles of
other movies by the same director, orders them by grossing and joins them
with the titles that are in the embedding website: this effectively creates a
reccomendation box which will link to other "popular" movies by the same
director available on the same website. Also retrieves rating from
Rottentomatoes.com

**Custom SPARQL query editor:**

Sparql query (if you want to try widget with your own query thick the box "Use custom query"):

```
PREFIX foaf:<http://xmlns.com/foaf/0.1/>
PREFIX og:<http://opengraphprotocol.org/schema/>
PREFIX dbpedia:<http://dbpedia.org/property/>
PREFIX rdfs:<http://www.w3.org/2000/01/rdf-schema#>

SELECT DISTINCT  ?Title ?RImage ?Starring ?Writer ?Director ?Budget ?Runtime
?_dburl
WHERE {
   <##URL##> og:type ?type.
   <##URL##> og:title ?Title.
   <##URL##> og:image ?image.
   <##URL##> og:url ?url.

   OPTIONAL {
      ?_rottenurl og:type ?type.
      ?_rottenurl og:title ?Title.
      ?_rottenurl og:image ?RImage.
      FILTER( !regex( STR(?RImage), "##DOMAIN##", "i")).
   }
OPTIONAL{
      {
        ?_dburl dbpedia:name ?Title.
      }
      UNION
      {
       ?_dburl foaf:name ?Title.
      }
      UNION
      {
       ?_dburl rdfs:label ?Title.
```
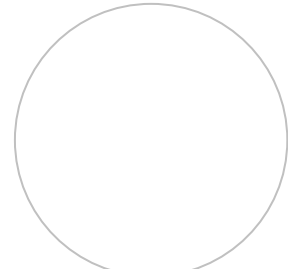
Enabling **networked** knowledge.

# Crowdsourcing rich annotations

FP7-SME-2010-1
SEMLIB 262301

Part of
## SemLib Project
## Semantic tools for digital libraries

SemLib - Annotation Demo - YouTube.flv

Enabling **networked** knowledge.

# Observations

- **RDF is very handy not per se a magic bullet**
  - ☐ Allows "late data integration efforts" but at a high cost: SPARQL esily gets prohibitively expensive
  - ☐ Used naively or too directly will allow bloated, unusable data structures
- **Large scale Data transformation capabilities come to the rescue.**
  **Models:**
  - ☐ Cloud Powered data pipelines allow free experiementation and data trasformations while the system keeps on running and serving users.
  - ☐ An extra dimension of flexibility and preprocessing which eases the task of SPARQL late
- **Now in trials with several private entities large enterprises**
  - ☐ Pharmaceuticals, Publishing

National University of Ireland, Galway
Ollscoil na hÉireann, Gaillimh

sfi
science foundation ireland
fondúireacht eolaíochta éireann

Enabling **networked** knowledge.

# Thanks

Enabling **networked** knowledge.