# The Web as a
# **Semantic Source**

**Ricardo Baeza-Yates**

**VP, Yahoo! Research**

# and Applications

## Outline

- Internet moved from a curiosity to a substrate for life activity
- Content growing, changing, diversifying, fragmenting
- Semantics of content unlock the value of the data ecosystem
- Explicit and implicit semantic sources
- Applications
- A virtuous feedback cycle for enhancing semantics

# Content Growth and Trends

## Content trends

| Content type | Amount of content produced per day |
|---|---|
| Published content | 3-4 GB |
| Professional web content | $\sim$ 2 GB |
| User generated content | 8-10 GB |
| Private text content | $\sim$ 3 TB (300x more) |
| Upper bound on typed content | $\sim$700 TB ($\sim$200x more) |

[Ramakrishnan and Tomkins 2007]

# Metadata trends

| Metadata type | Amount of metadata produced per day |
| --- | --- |
| Anchortext | 100 MB |
| Tags | 40 MB |
| Pageviews | 180 GB |
| Reviews | Around 10 MB |

[Ramakrishnan and Tomkins 2007]

# Content ownership

- Content consumption is fragmenting – nobody owns more than 10% of the Web page views

- No single place will own all the content

- Best of breed processing will operate on the web version (?)

- Value transitions to ecosystem

# Content Consumption is fragmenting

1 to 3 | 0.5 | treats, catnips, daddy, mommy, purring, mice,

**del.icio.us / tag / jsr168**

All items tagged **jsr168** (**create tag description**) → view **popular**

« earlier | later »

**By topic**

**ONJava.com -- JSR168 portlet example** save this
by kwangomango to portlet jsr168 ... saved by 22 other people ... 4 hours ago

**Introducing Java Portlet Specifications: JSR 168 and JSR 286** save this
Portlets are Web-based components managed by portlet containers that supply dynamic content. Portals employ portlets as pluggable Portlet Specification achieves interoper
by pedavison to Portal Portlet portlets j2ee java jsr168 programming development ... saved by 71 other people ... 1 day ago

**portlet-container: OpenPortal Portlet Container Project supporting JSR 168 and JSR 286 portlets** save this
by johnalewis to java opensource portal portlet jsr168 ... saved by 36 other people ... 1 day ago

**FitzBlog : WSRP and JSR168 Are Two Completely Different Things...** save this
by bonifax to WSRP JSR168 portal SharePoint Microsoft portlet ... saved by 22 other people ... 1 day ago

**JSR 168 programming** save this
by nitingr to portlet tech jsr168 portlets ... saved by 1 other person ... 1 day ago

**Marina Sum's Blog: Weather Portlet in Portlet Repository** save this
by odalet to java portal portlet jsr168 ... saved by 1 other person ... 2 days ago

**Best Practices for Applying Ajax to JSR 168 Portlets** save this
by emcconne to ajax jsr168 portlets portal portlet ... saved by 94 other people ... 2 days ago

**light: Project Home Page** save this
by lhotari to portlet portal jsr168 ... saved by 75 other people ... 3 days ago

# Content access is fragmenting

**facebook**

Profile edit | Friends ▼ | Networks ▼ | Inbox (2) ▼ | home account privacy logout

Search ▼

Applications edit
- Photos
- Groups
- Events
- Marketplace
- Many Eyes Visualizer
▼ more

**Privacy Settings for Search**

You will show up in search results if anyone searches for "**andrew tomkins**" or any part of your name. Even though anyone can search for you, only your friends and everyone from Yahoo, MIT, Carnegie Mellon and Silicon Valley, CA, can see your profile. In addition, people in college networks, high school networks, company networks, regional networks and no networks can see you in search results. People who can't see your profile can see your profile picture, poke you, message you and send you a friend request from your search listing.

Back to Privacy Overview without saving changes.

**Who Can Find Me in Search and See My Public Se**

You can allow **everyone** on Facebook to find you in search resu[...] members list, or you can select **restricted** settings to allow only certain people from insid[...]o find you in search results. Your friends can always find you in search results.

**MIT**

Network Info
Members: 22,504
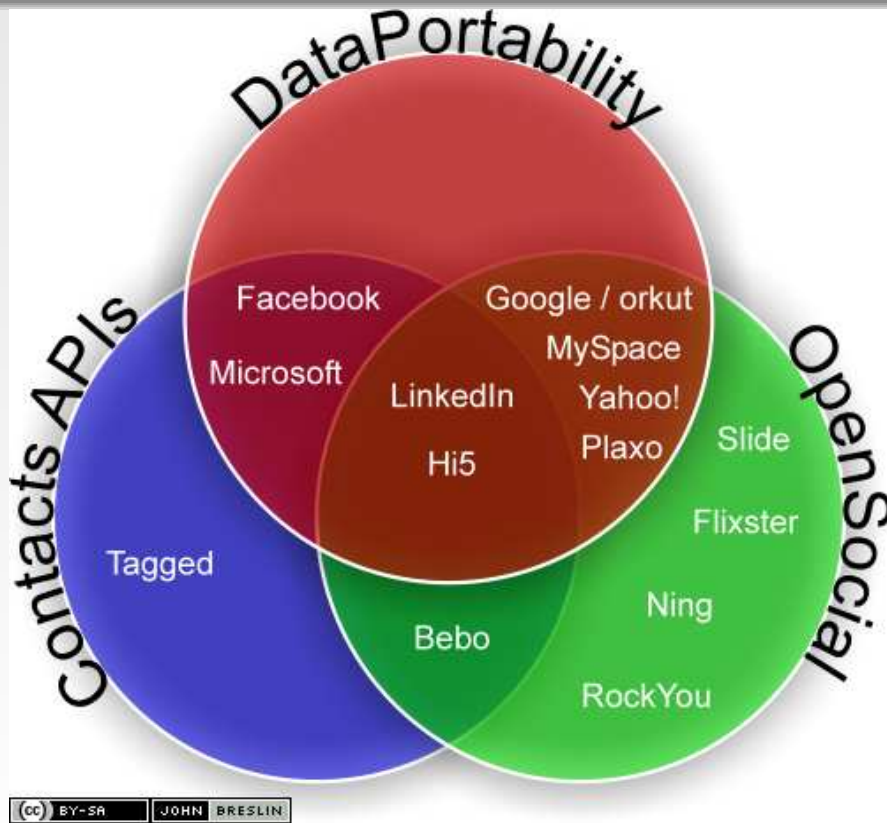Friends: 11
Type: College
Location: Cambridge, MA

**Which Facebook users can find me in search?**

Some of my networks and all my friends ▼

☑ Yahoo
☑ MIT
☑ Carnegie Mellon
☑ Silicon Valley, CA
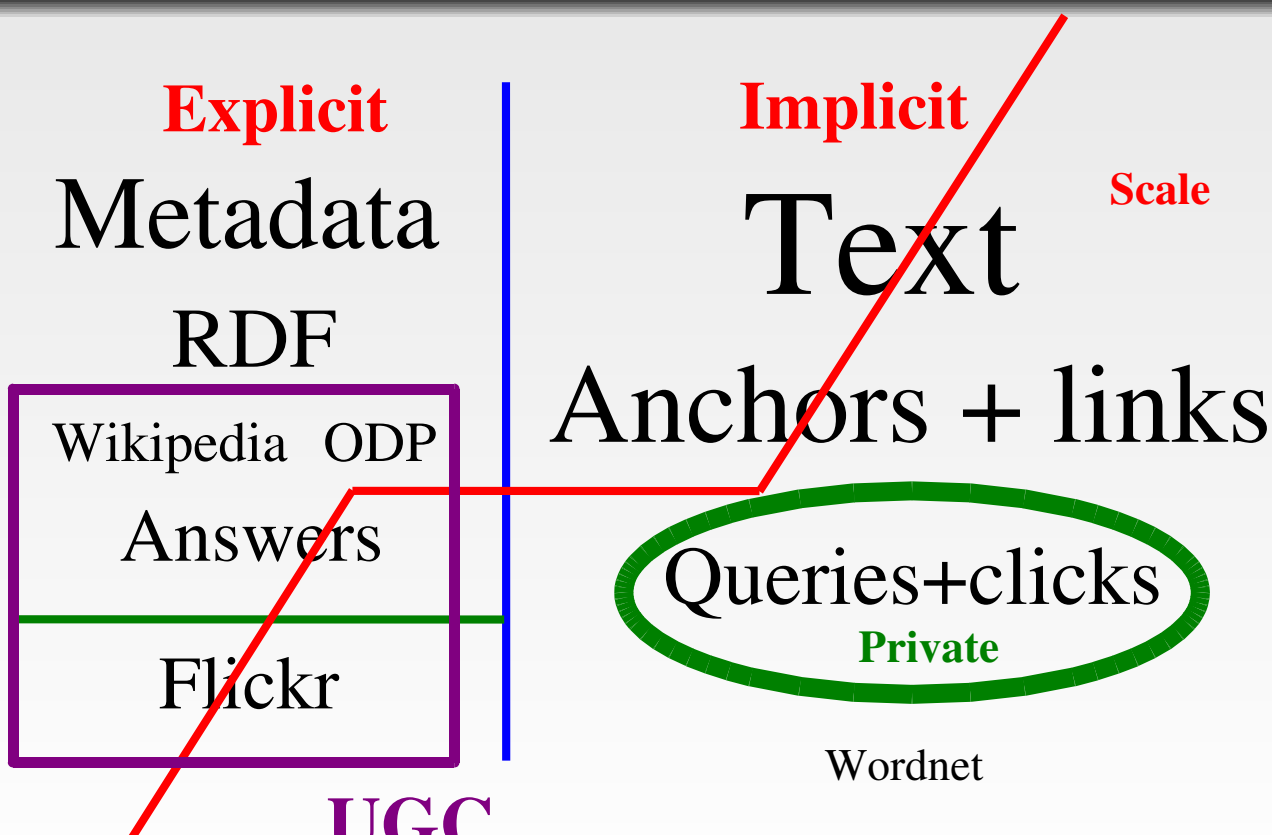
Semantic Sources

- Intrinsically Collaborative
- Explicit or Implicit
- Taxonomies vs. Folksonomies
- Different Size and Growth
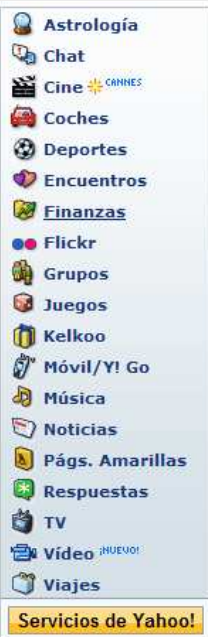- Diversity of Quality
- Public or Private

## Examples

**Explicit**

**Implicit**

Metadata

Text

**Scale**

RDF

Wikipedia  ODP

Answers

Anchors + links

Flickr

Queries+clicks

**Private**

Wordnet

**UGC**

# User Generated Content at Yahoo!

Social properties had 115M unique visitors worldwide, 56M "under 35".

| | | |
|---|---|---|
| • | **Yahoo! Grups** | 8 million groups, 1 in 10 Internet users |
| • | **Del.icio.us** | 2 million users |
| • | **Flickr** | 1 million photos uploaded daily |
| • | **Yahoo! Answers** | 90M unique users, 250M answers |
| • | **Messenger** | 85M unique users |

(2007 data)

- James Surowiecki, a **New Yorker** columnist, published this book in 2004
  - **"Under the right circumstances, groups are remarkably intelligent"**
- Importance of diversity, independence and decentralization **Aggregating data**

  *"large groups of people are smarter than an elite few, no matter how brilliant—they are better at solving problems, fostering innovation, coming to wise decisions, even predicting the future".*

- 23 -

- Crucial for Search Ranking
- Text: Web Writers & Editors
  - not only for the Web!
- Links: Web Publishers
- Tags: Web Taggers
- Queries: All Web Users!
  - Queries and actions (or no action!)

Metadata: Microformats

# Microformats

- microformats.org
- Originated by Tantek Celik and others
- Agreements on the way to encode certain kinds metadata in HTML
  - Reuse of semantic-bearing HTML elements
  - Based on existing standards
  - Community process
  - Persons, events, listings etc. but also syntactic metadata: licenses, tags
- Microformats have no shared syntax
  - Each microformat has a separate syntax tailored to the vocabulary
- Microformats are not ontologies
  - No formal descriptions of schema, only text
  - Limited reuse, extensibility of schemas
  - No datatypes
- No namespaces, unique identifiers (URIs)
  - no interlinking
  - mapping between instances is required
- Relationship to page context is unclear
- Widely used in millions of documents
  - User-generated as well as automatically generated

# Example: tags and machine tags

# Metadata is out there

- Question:
  - Just how much data is out there?
  - What is the quality?
- Idea: bring metadata to the surface of search
- How does it work?
  - User enters query
  - Metadata is extracted dynamically
  - Entity reconciliation
  - Metadata is used to display
    - rich abstracts,
    - related pages
    - spatial, temporal visualization
- Microsearch prototype
  - Play at http://www.yr-bcn.es/demos/microsearch/

# Example: ivan herman

# Example: peter site:flickr.com



Flickr users named "Peter" by geography

# Example: san francisco conference



Conferences in San Francisco by date

# Example: greater st. peter

# Lessons

- More metadata than we expected
  - 53% of unique queries have at least one metadata-enabled page in top 10 (n=7848)
- Performance is poor
  - Metadata needs to come from the index for performance
- Metacrap does exist
  - Users <u>have to</u> see metadata to spot mistakes in their markup, warn others
- RDF templating is hard
  - Adds extra complexity
- Scalability

# Exploting Metadata: SearchMonkey

## SearchMonkey

- Creating an ecosystem of publishers, developers and end-users
  - Motivating and helping publishers to implement semantic annotation
  - Providing tools for developers to create compelling applications
  - Focusing on end-user experience
- Rich abstracts as a first application
- Addressing the long tail of query and content production
- Standard Semantic Web technology
  - dataRSS = Atom + RDFa
  - Industry standard vocabularies
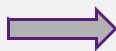
# What is SearchMonkey?

an open platform for using structured data to build more useful and relevant search results

## Before

Topics for **Getting Pregnant** - BabyCenter
Find out how to boost your chances of **getting pregnant**, what you can do if you're having a problem conceiving, and more. **...** do before you try to **get pregnant ...**
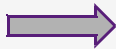www.**babycenter.com/getting-pregnant** - 91k - Cached

WebMD **Allergies** Health Center - Find allergy information and latest **...**
Information and articles on the diagnosis, symptoms, treatment, and the prevention of **allergies**.
www.**webmd.com**/diseases_and_conditions/**allergies**.htm - 132k - Cached

**Italy Travel Guide** and **Travel** Information - Lonely Planet
Lonely Planet **Italy** includes information on events, attractions, activities, and transportation for the independent traveler.
www.**lonelyplanet.com**/worldguide/destinations/europe/**italy** - 57k - Cached

## After

Topics for **Getting Pregnant** - BabyCenter
Tools | BabyCenter's getting **pregnant** tools and information
Before you Try | can boost your chances of conception by helping you
Active Trying | chart your cycle, read your cervical mucus, and
Having Trouble | pinpoint ovulation.
www.**babycenter.com/getting-pregnant** - 76k - Cached

WebMD **Allergies** Health Center - Find allergy information and latest .
Overview | **Allergies** are an abnormal response of the immune system.
Symptoms | People who have allergies have an immune system that reacts
Diagnosis | to a usually harmless substance in the environment. This
Treatment | substance (pollen, mold, dander, etc.) is called an allergen.
www.**webmd.com/allergies** - 111k - Cached

**Italy Travel Guide**: Overview - Lonely Planet WorldGuide
Overview | When to go: **Italy** is at its best in spring (April-May) and
Money & Costs | autumn (October-November). During these seasons, the
Itineraries | scenery is beautiful, the temperatures are pleasant and
Sights | are relatively few crowds.
www.**lonelyplanet.com**/worldguide/destinations/europe/**italy** - 46k - Cached

# Enhanced Result

Yahoo!  My Yahoo!  Mail    Welcome, **Guest** [Sign In]  Help

Web | Images | Video | Local | Shopping | more ▾
art of pizza chicago    [Search]   Options ▾

YAHOO!

1-10 of about 9,410,000 for **art of pizza chicago** (About this page) - 0.23 sec.

SPONSOR RESULTS

**deep links**

**image**

The **Art of Pizza** - Lakeview - **Chicago**, IL 60657
Reviews | Ratings: ★★★★★ (173)
Photos | Address: 3033 N Ashland Ave, Chicago, IL
Send to a friend | Phone: (773) 327-5600
Send to Phone | Price Range: $
www.acmereviews.com/biz/the-**art-of-pizza-chicago** - 145k - Cached

**name/value pairs or abstract**

UGC: Exploiting Flickr Tags

## Tag Mining

- Objective:
  - Deploy collective knowledge that exists within Social Media services (Flickr and Delicious)

- Approach:
  1. Use tag co-occurrence statistics for media annotation and retrieval
  2. Semantic analysis of large tag-spaces

## Tag Mining - Collective Knowledge



- Many users annotate photos of "La Sagrada Familia":
  - Sagrada Familia, Barcelona
  - Sagrada Familia, Gaudi, architecture, church
  - church, Sagrada Familia
  - Sagrada Familia, Barcelona, Spain

- Derived collective knowledge:
  - Barcelona, Gaudi, church, architecture

## Tag Mining - Semantics

- Assign tag semantics using WordNet broad categories



- – Paris :: location
- – Eiffel Tower :: artifact
- – Coverage: 52% of tag volume

## Tag Mining - Semantics

- Extend this mapping using patterns found in Wikipedia
  - – Upperbound for coverage: 78.6% of the tag volume
  - – Based on SVM approach
    - Features: Wikipedia templates and categories
    - Training data: Wikipedia entries found in WordNet
  - – Extended coverage: 68% of the tag volume
  - – Mapping from Wikipedia pages to tags
    - Reduces ambiguity in the classification

# Understanding tags

**London Eye**

London Eye and Golden Jubilee Bridge seen from Westminister Bridge.

**Tag list**

london eye, thames,

**Suggested tags**

- ☑ london
- ☑ england
- ☑ uk
- ☑ river
- ☐ eye
- ☑ south bank
- ☐ big ben
- ☐ night
- ☑ bridge
- ☐ 2006

[ Update annotation ]

# Semantic Photo Search

Query → **flickr** → unsorted photos

tags

tag graph

tag:type

tag:type

tag:type

**WORDNET**

**flickr**

lights

architecture

sun set

# Media Search - Demo

# Understanding Text

## Document Understanding Cartoon

our work!

Complexity of Document Understanding

grep

search engines

semantic web?

domain expert

Q & A

# Extending metadata

**Pablo Picasso** was born in <u>Málaga</u>, <u>Spain</u>.
       **PER**                          **LOC**       **LOC**

E:PERSON                     GPE:CITY   GPE:COUNTRY

artist:name               artist:placeofbirth   artist:placeofbirth

> If most artists are persons, than let's assume all artists are persons.
> If most places of birth are locations, then let's assume all are.



- 78 -

# Entity Containment Graph

# Example: Picasso

# Synthetic Document



Query

Article
Article
} Syntactic matches: extract snippets

Sentence
Sentence
Sentence
Sentence

Article
Article
Article
Article
} Other: group by relevant categories

# Synthetic Document

Article

**Article**

**Climbing**

Climbing is the activity of using one's hands and feet to move up the surface of a steep object. It is pursued both recreationally, either to get to a destination otherwise inaccessible or for its own enjoyment, and also professionally, as part of activities such as maintenance of a structure, or military operations.
Rock climbing, the scaling of steep rocky surfaces, is perhaps the most familiar sort of climbing; other types of climbing include ice climbing, tree climbing, buildering (climbing on the outside of buildings), and pole climbing.
Mountaineering, the general activity of ascending mountains, often requires the use of climbing techniques.
Climbing may be divided into two broad categories : aid climbing and free climbing .
Climbing communities in many countries , as well as individual regions , have developed their own climbing rating systems .
Hiking , Bouldering , Roped free climbing , and Aid climbing all share these factors to one degree or another .
more about Climbing ...

**Category**

Category
Article
Article
Article
Article

**Climbing areas**

Rock climbing in the Peak District: Generally the climbing style is free climbing ( as opposed to aid climbing ) and the rock is either gritstone or limestone . ... There is a long-standing practice of climbing routes in the traditional climbing style .
Skaha Bluffs: The area is mostly a sport climbing area , though many traditional climbing opportunities also exist . ... Skaha Bluffs is slowly gaining recognition as a destination climbing area due to its mild weather , easy access and high number of sport climbing routes .
Fair Head: The Dal Riada Climbing Club maintains a climbing hut in the area . ... Categories : Northern Ireland geography stubs | Climbing stubs | Climbing areas | Headlands of County Antrim
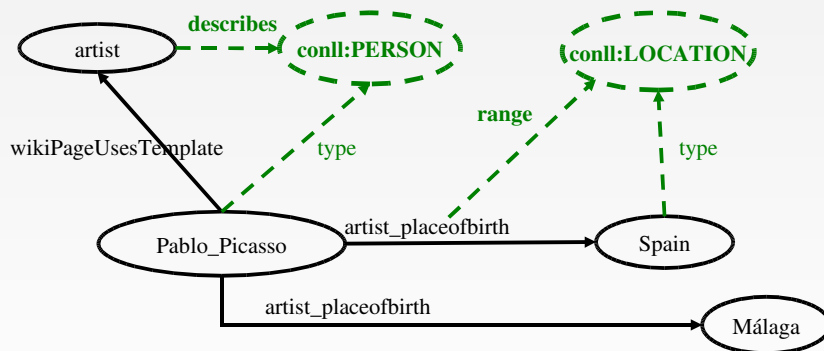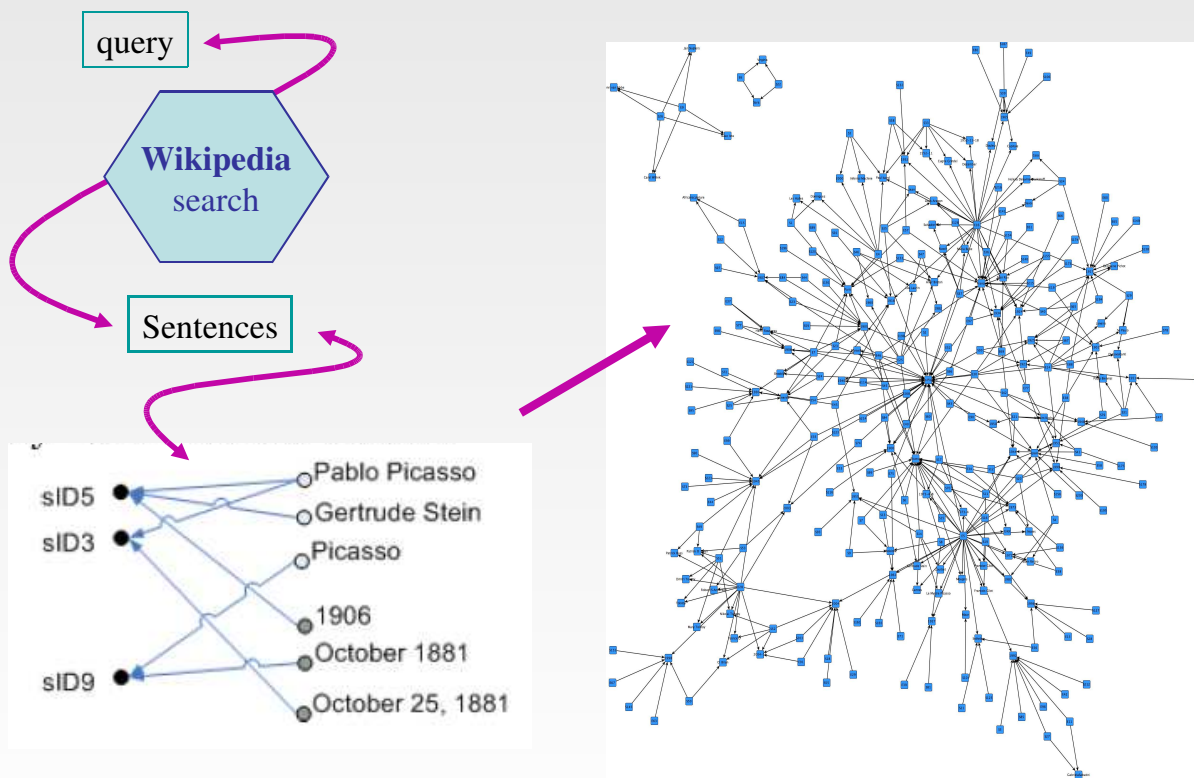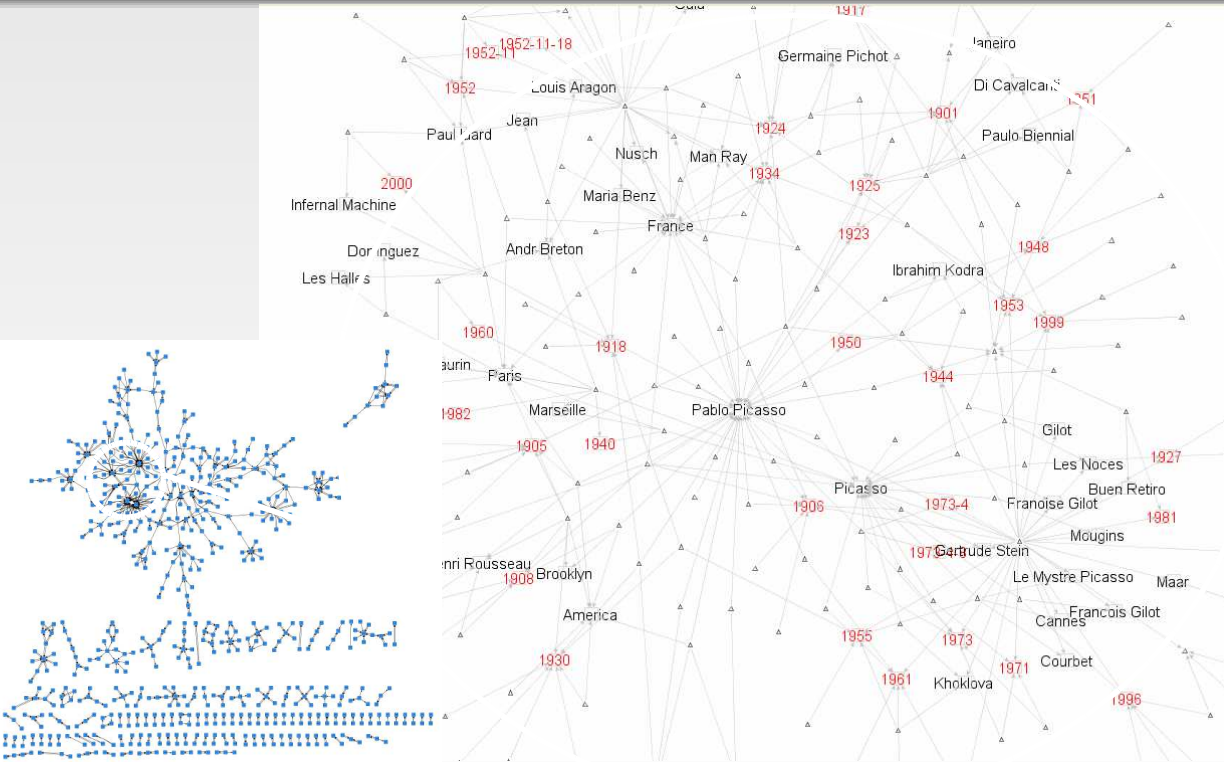Mount Nemo Conservation Area: Categories : Climbing areas | Halton Region , Ontario | Parks in Halton | Climbing stubs | Golden Horseshoe geography stubs
more about Climbing areas ...

**Category**

Category
Article
Article
Article

**Climbing techniques**

Head point: In climbing , head pointing a trad climbing route essentially means `` leading the climb after top-roped practice '' . ... Categories : Orphaned articles from September 2006 | Climbing techniques | Climbing stubs
Climbing command: A climbing command is a short standard phrase used in climbing to ensure the smooth operation of the climbing system .
Autoblock: It involves using a friction hitch around the climbing rope , and may be combined with other climbing equipment .
more about Climbing techniques ...

# Related Entities Search



Sentences      Entity graph      Top entities
Possibly filtered by type

Query

1 relevant
sentence per entity

emergency landing

Overview | Names | Places | Concepts | Events | Photos | Queries | News | Answers | Sites | All related

Locations related to 'emergency landing'

"emergency landing"

**Town**
Okinawa Baghdad Goose Bay Los Angeles

**Country**
Cuba

**Island**
Hainan Island

- (From W Narita_International_Airport) "The airliner was able to make an emergency landing in **Okinawa** . "
- (From W Philippine_Airlines_Flight_434) "The Boeing 747-283B , tail number EI-BWF , made an emergency landing in Naha Airport , **Okinawa** , one hour after the bomb exploded . "
- (From W Oplan_Bojinka) "The Boeing 747-200 safely made an emergency landing in Naha , **Okinawa** . "

## Syntactic/Semantic Tagging

- Goals: Identify multiword expressions and entities, support generalization, coarse disambiguation

- Tagger: Average perceptron HMM (Collins, 2002) general purpose tagger: efficient (millions of features, hundreds of classes), fast (thousands of sentences/sec)

- State of the art: 3rd NIST Automatic Content Extraction Evaluation (ACE) 2007 (Surdeanu & Ciaramita, 2007)

- Tasks: PoS Tagging, supersense tagging (Ciaramita & Altun, 2006), named entity detection (CoNLL, BBN-WSJ, ACE, etc.)

- Research problems: robustness on Web data (domain adaptation), learning/evaluating from user-generated data (Mika et al., forthcoming)

## Parsing/SRL

- Goal: extract structured information at sentence level (beyond the bag of words/document-centric models)

- Dependency parsing:

  - Parser: Fast discriminative multilingual Shift/Reduce parsing (hundreds of sentences/sec) (Attardi, 2006); 2nd in Adaptation Task of CoNLL 2007 (Attardi et al., 2007)

- Semantic role labeling:

  - Joint parsing and SRL: 3rd best system at CoNLL 2008 (Forthcoming)

- Research problems: how can structured linguistic representations be used to improve search/ranking problems? (Surdeanu, Ciaramita & Zaragoza, ACL 2008)

# Web Usage:
## Extracting Semantics from Queries

## Relating Queries (Baeza-Yates, 2007)



**common session**

q1 ⟷ q2          q3          q4          queries

**common words**

**common clicks**

clicks

pages

**links**

w          w

**common terms**

| Graph | Strength | Sparsity | Noise |
|-------|----------|----------|-------|
| Word | Medium | High | Polysemy |
| Session | Medium | High | Physical sessions |
| Click | High | Medium | **Multitopic pages Click spam** |
| Link | Weak | Medium | Link spam |
| Term | Medium | Low | Term spam |

## Click Graph

# Node Degree Distribution

# Connected Components

## Implicit Folksonomy?



## Set Relations and Graph Mining

- Identical sets: **synonyms**

- Subsets: **specificity**     **Baeza-Yates & Tiberi**
  **ACM KDD 2007**
  – directed edges

- Non empty intersections (with threshold)
  – degree of relation

- Dual graph: URLs related by queries
  –High degree: multi-topical URLs

## Evaluation: ODP Similarity

- A simple measure of similarity among queries using ODP categories

  - Define the similarity between two categories as the length of the longest shared path over the length of the longest path

  - Let $c_1,.., c_k$ and $c'_1,.., c'_k$ be the top $k$ categories for two queries. Define the similarity (@$k$) between the two queries as $max\{ sim(c_i,c'_j) \mid i,j=1,..,K \}$
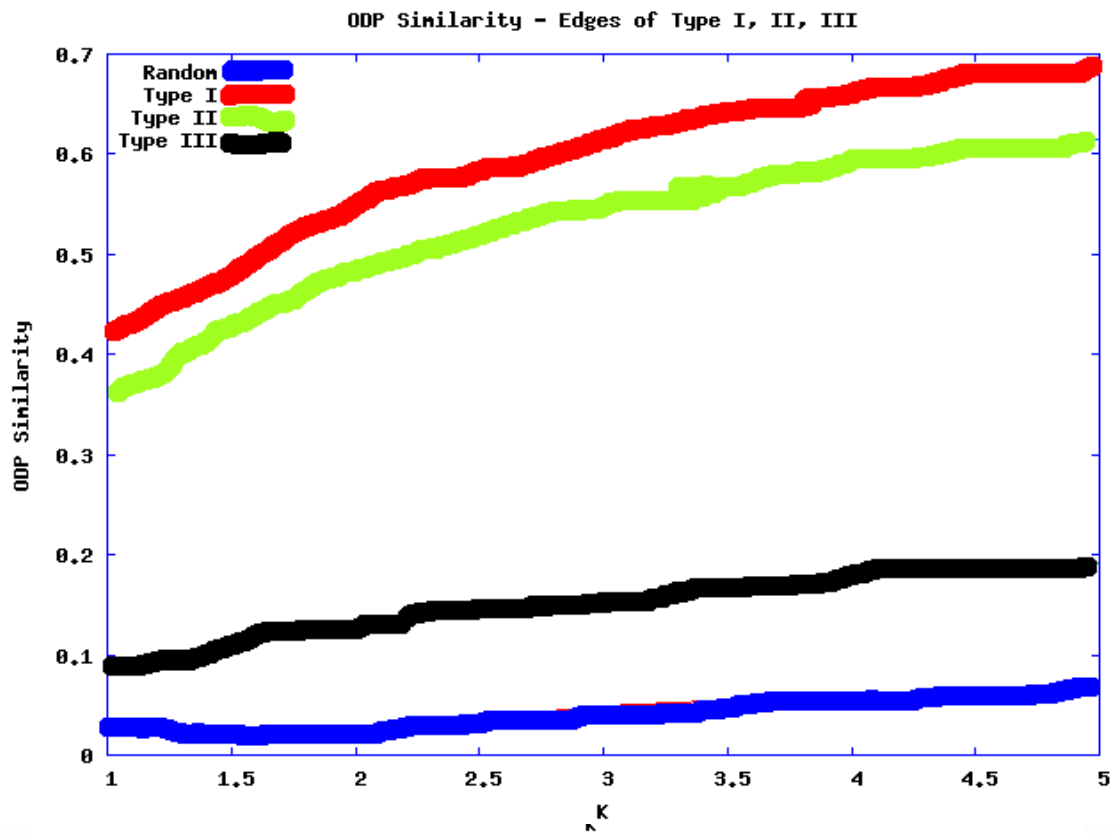
## ODP Similarity

- Suppose you submit the queries "*Spain*" and "*Barcelona*" to ODP.

- The first category matches you get are:

  – Regional/ Europe/ Spain

  – Regional/ Europe/ Spain/ Autonomous Communities/ Catalonia/ Barcelona

- Similarity @1 is 1/2 because the longest shared path is "Regional/ Europe/ Spain" and the length of the longest is 6

## Experimental Evaluation

- We evaluated a  sample of 1,000 thousand edges for each kind of relation

- We also evaluated a sample of 1,000 random pairs of not adjacent queries (baseline)

- We studied the similarity as a function of $k$ (the number of categories used)

ODP Similarity – Edges of Type I, II, III

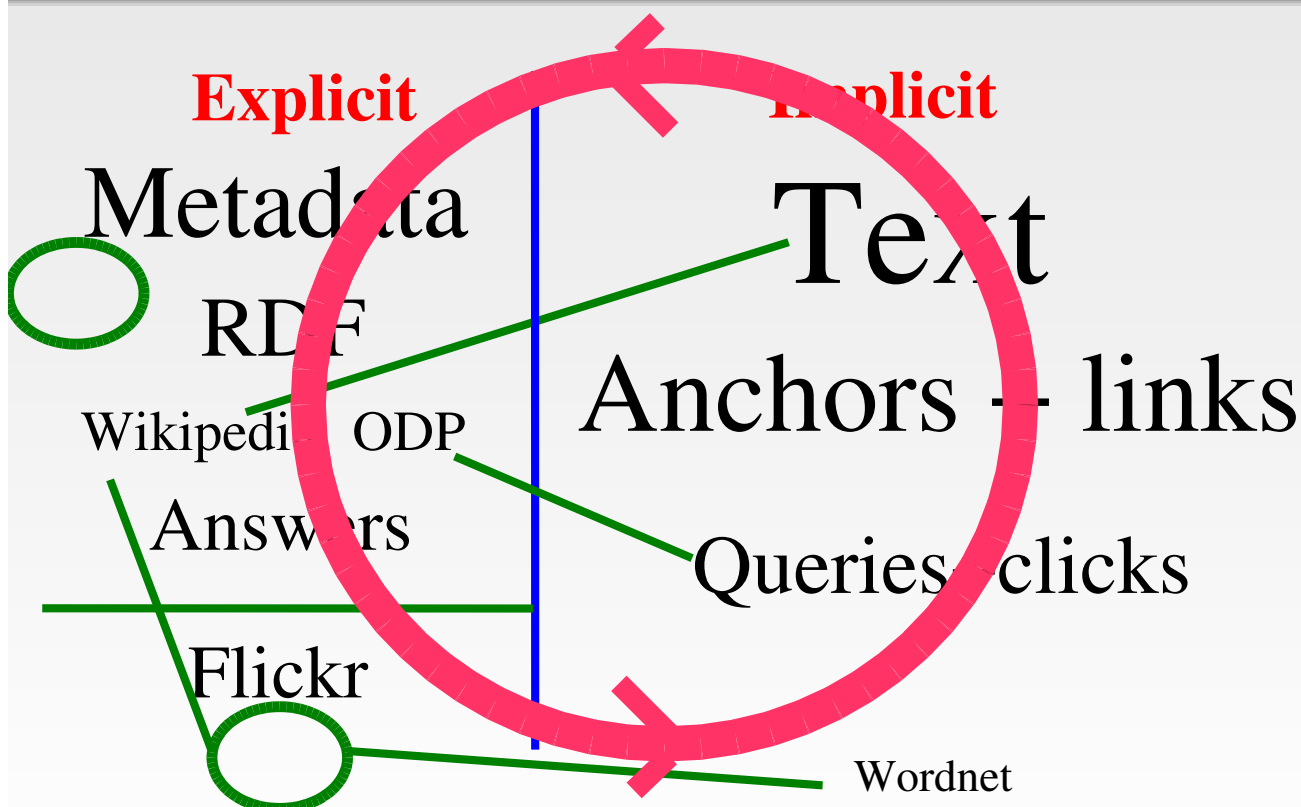# Final Remarks

- Explicit vs. implicit social networks
  - Any fundamental similarities?

- How to evaluate with (small) partial knowledge?
  - Data volume amplifies the problem

- User aggregation vs. personalization
  - Optimize common tasks
  - Move away from privacy issues

## The Virtuous Cycle

**Explicit**  **Implicit**

Metadata  Text

RDF

Wikipedia ODP  Anchors + links
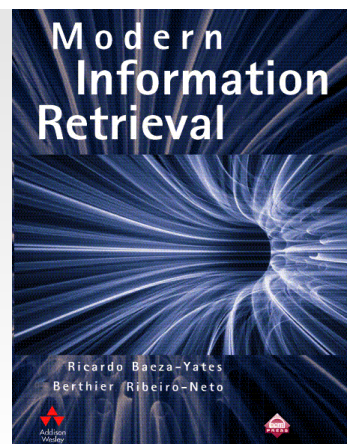
Answers

Queries clicks

Flickr

Wordnet

# The Future: Web 3.0?

- We are at Web 2.0 beta

- People wants to get tasks done
  - Where I do go for a original holiday with 1,000 US$?

- Take in account the context of the task

Start ──── [ **Yahoo! Experience** ] ──▶ Finish

**Second edition coming soon**

Modern Information Retrieval

Ricardo Baeza-Yates
Berthier Ribeiro-Neto

## No Questions?
**Contact:** rbaeza@acm.org

**Thanks to** Massi Ciaramita, Peter Mika, Borkur Sigurbjornsson, Mihai Surdeanu, Andrew Tomkins, Roelof van Zwol, Hugo Zaragoza