# WeKnowIt

## Making the Collective Intelligence of Social Media Searchable

Yiannis Kompatsiaris

CERTH – ITI

CIVR 2010, July 2010

http://www.weknowit.eu

# Contents

- Introduction
- Social Media - Collective Intelligence
- WeKnowIt approach
- Community detection in Social Media
- Demos
- Conclusions - Issues

# Web 2.0 content

### flickr
- 3,190 uploads in the last minute
- 3.2 million things geotagged this month
- 4,754,012,299 photos (2 July 2010)

### YouTube
- 24h of video content uploaded every minute
- 2 billion movies watched every day

### facebook
- More than 400 million active users
- More than 200 million users log on at least once each day
- 2.5 billion photos uploaded each month

**flickr** from YAHOO!

Home    You -    Organize & Create -    Contacts -    Groups

## Winner



The winner of the WeKnowIt Grand Travel Challenge

weknowit

# Tags, content everywhere
## Upload, tag, share, search

# Can we do more things?

Search     Photos | Groups | People

SEARCH   Full Text | Tags Only
Advanced Search

Sort: **Relevant** | Recent | Int

**Tag Clusters**

Photos with tags like **nyc, newyork and manhattan**

Photos with tags like **fruit, red and green**

Photos with tags like **ipod, iphone and music**

Small | Medium | Detail | Slideshow

From sonnyhung

From Wan

From Taxi Lady...

From ( karen )

From HAZEL- 名 b-女港

From jonbradbury

**By combining information from many photos - tags, it seems that we can extract**

**Stable patterns**

**in tagging systems over time**

From amy johanns

From nnvica

From fernando780

From jordanmerric...

From humedni

weknowit

Low-level

Time

Geo

User
Profile

Groups

**flickr**

**Deutsches Eck from Ehrenbreitstein
Fortress, Koblenz, Germany**

by **schaengel**

💬 121 comments ⭐ 69 faves

Tagged with **koblenz, ehrenbreitstein** ...
Taken on **November 15, 2009**, uploaded
**November 17, 2009**

📷 See **more of schaengel photos**, or visit
his **profile**.

When you're high up on the hill above Koblenz at Ehrenbreitstein Fortress you can get a
great panoramic view of the city and the surrounding area.

Comms

Favs

Tags

Caption

Social
network

**weknowit**

# ... and more: Travel trends using flickr



Trace Flickr users from a chronologically ordered set of geographically referenced photos

*Who are the Italians and who are the Americans?*

*MIT SENSEABLE CITY LAB, "The World's eyes"*

# Defining Collective Intelligence

Collective Intelligence is the Intelligence which emerges from the collaboration, competition and coordination among individuals.



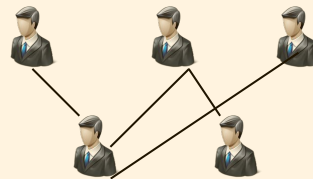...an Intelligence greater than the sum of the individuals' intelligence

# WeKnowIt and CI

## Collective Intelligence

### Social Intelligence



Social Networks

### Media Intelligence



User-generated content, social context

### Mass Intelligence



Blogs, forums, ratings, voting

### Organizational Intelligence

Knowledge Management

### Personal Intelligence

Upload and Access

weknowit

# Personal Intelligence

>> Login, Upload

>> Tag recommendation,

>> Spam detection

## Organisational Intelligence

>> Log Merger

## Access

# Emergency Response

## Media Intelligence

**Picture arrives at emergency response**

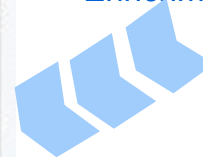>> Automatic localisation of photo

## Mass Intelligence

**Many contributors**

>> Clustering
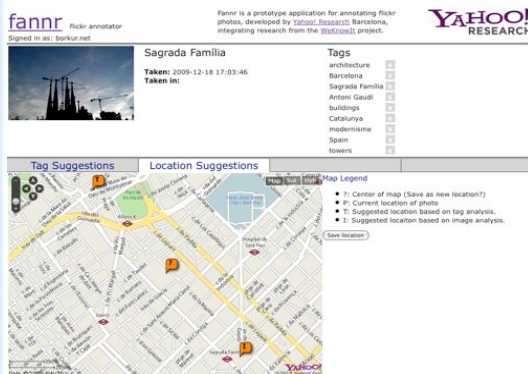
>> Tag Harmonization

>> Enrichment from add. sources

## Social Intelligence

>> ER Alert Service

# Travel prototype

## Mass Intelligence

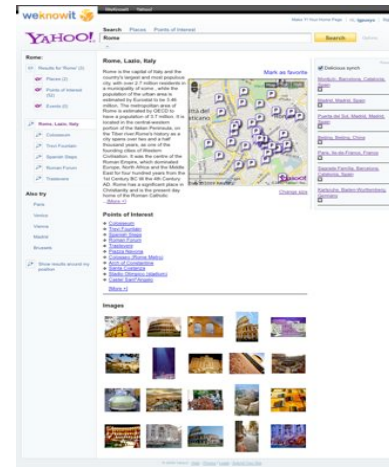**>> Media Collection:** flickr, query logs,

**>>** Automatic generation of ranked facet lists of POIs

## Media Intelligence

**>>** Hybrid Clustering

**>>** Image Localisation

**>>** Tag suggestions

## Travel Prepa-ration

## Personal Intelligence

**Profile of contributor**
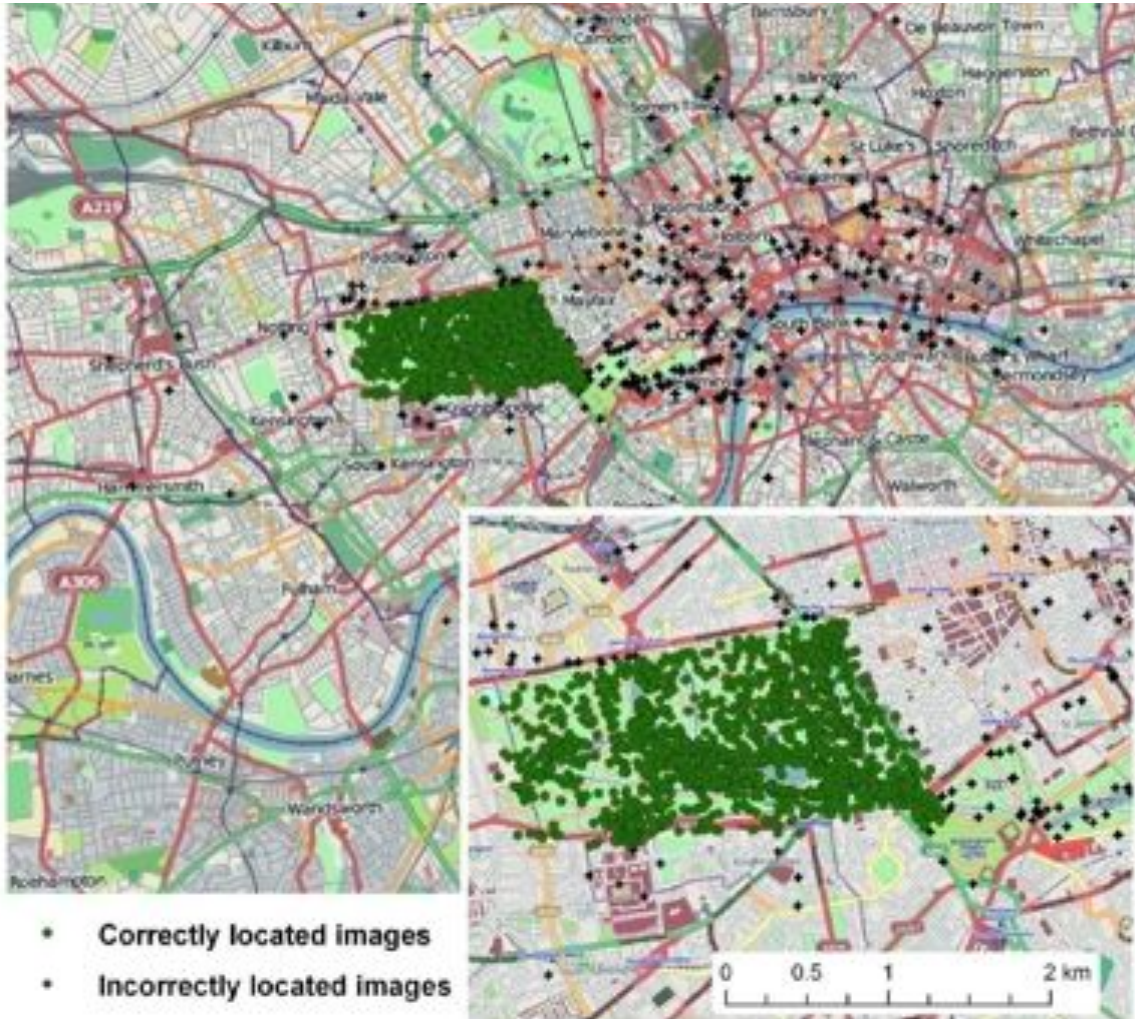**>>** Recommendations

## Social Intelligence

**Can your social network be of help?**

**>> Friends position, alert**

## Post Travel

## Mobile Guidance

# weknowit

# Relevant Activities (*ER*)



- Correctly located images
- Incorrectly located images

Automatically describe city cores

Distinction between administrative and vernacular uses of place names

Potential for confusion in the dispatch of emergency services

*Livia Hollenstein and Ross S. Purves, "Exploring place through user-generated content: using Flickr to describe city cores", JOURNAL OF SPATIAL INFORMATION SCIENCE*

# Relevant activities



MIT Center for Collective Intelligence

http://cci.mit.edu/index.html

The Climate Collaboratorium

Collective prediction  ...*accurate predictions about future events such as product sales, political events, and outcomes of medical treatments*....

Collective intelligence in healthcare

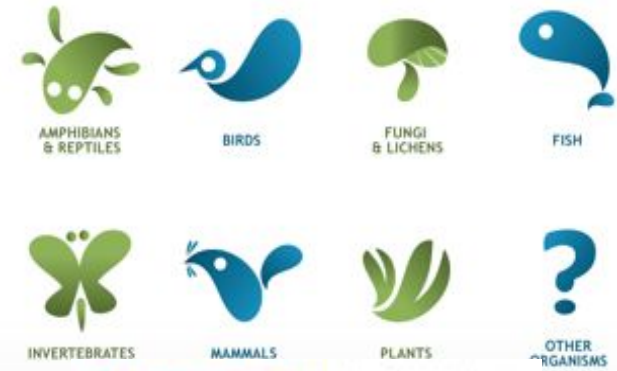Studying collective intelligence in today's organizations

# Relevant activities

collective intellect

Live Traffic Information

NOKIA

FixMyStreet

Mobnotes beta™
The World, Around You!

Citizens Connect

Graffiti
Public Alley 421, Boston

Map    List

iSpot
your place to share nature

AMPHIBIANS & REPTILES    BIRDS    FUNGI & LICHENS    FISH
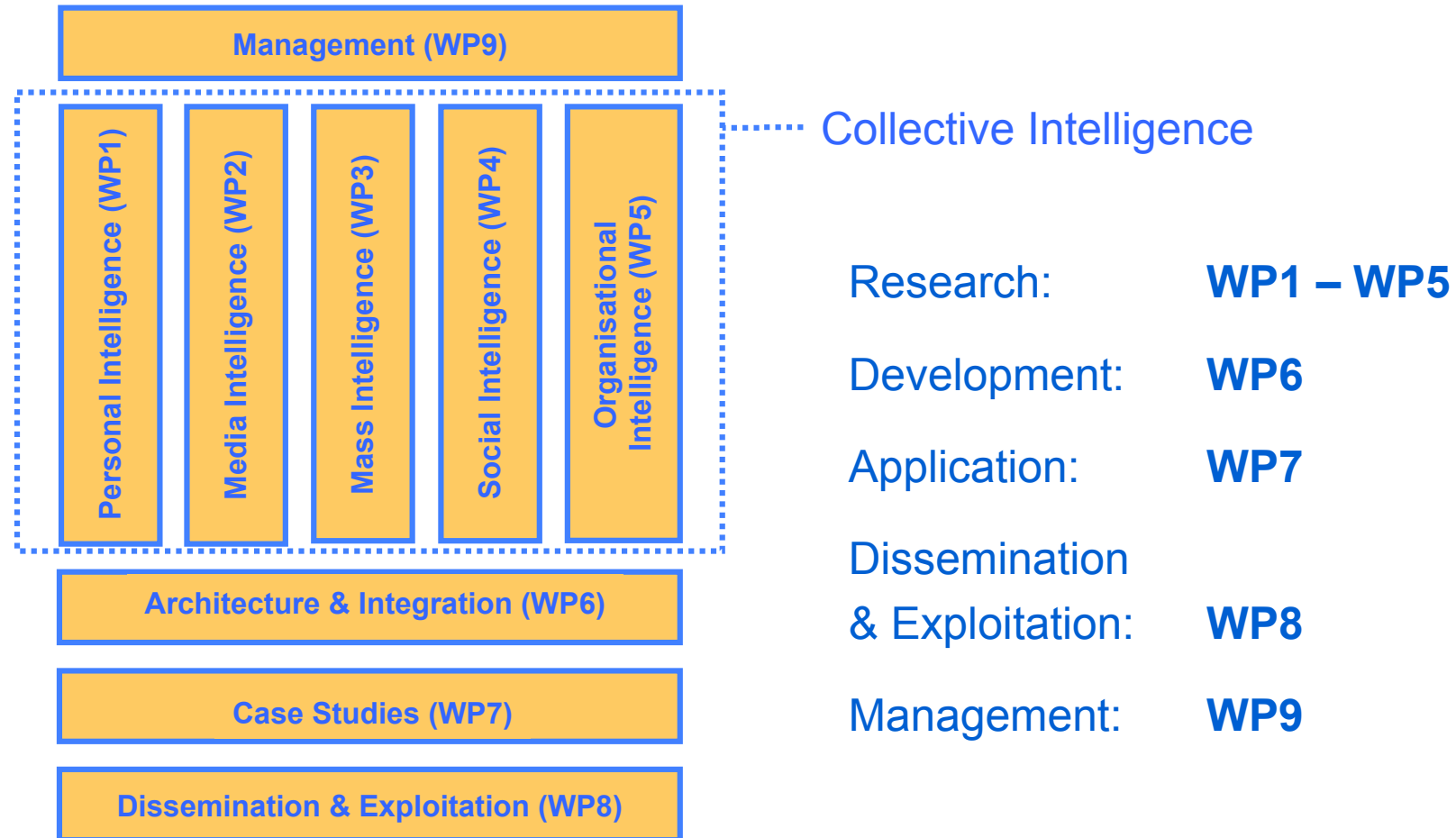
INVERTEBRATES    MAMMALS    PLANTS    OTHER ORGANISMS

DOPPLR

Dopplr helps you share your ... travel .... and exchange tips ... presents this **collective intelligence** - the travel patterns and advice ... as the Social Atlas.

http://traffic.berkeley.edu        Boston Citizens Connect

weknowit

# Relevant activities

- Most applications are still harnessing Collective Intelligence
  - Upload applications
- Emphasis is on visualization of results
- Few focus on analysis
- No fusion of modalities – sources
- Unlimited applications
    - Draught detection, through estimation of green levels in flickr photos for fire protection in Spain (MIT Senseable lab)
    - Hollywood stock Exchange – HP Labs

# Project work Overview

# Content in WeKnowIt

offline → model creation, training

**Non-Web 2.0 training data**

Standard annotated corpora used for training.
- **Single-modality:** text (Brown corpus), speech (TIMIT database), image (Corel database)
- **Single-source:** prepared by a single person/organization
- **Consistent quality:** absence of spam, malicious or erroneous data
- **Small-moderate volume:** Manually produced

**Massive Web 2.0**

Massive user generated content and feedback from Web 2.0 applications
- **Multi-modality:** e.g. image + tags, image + geo-location + time
- **Multi-source:** may be generated by different applications, user communities, e.g. delicious, StumbleUpon and reddit are all social bookmarking sites
- **Inconsistent quality:** noise, spam, ambiguity
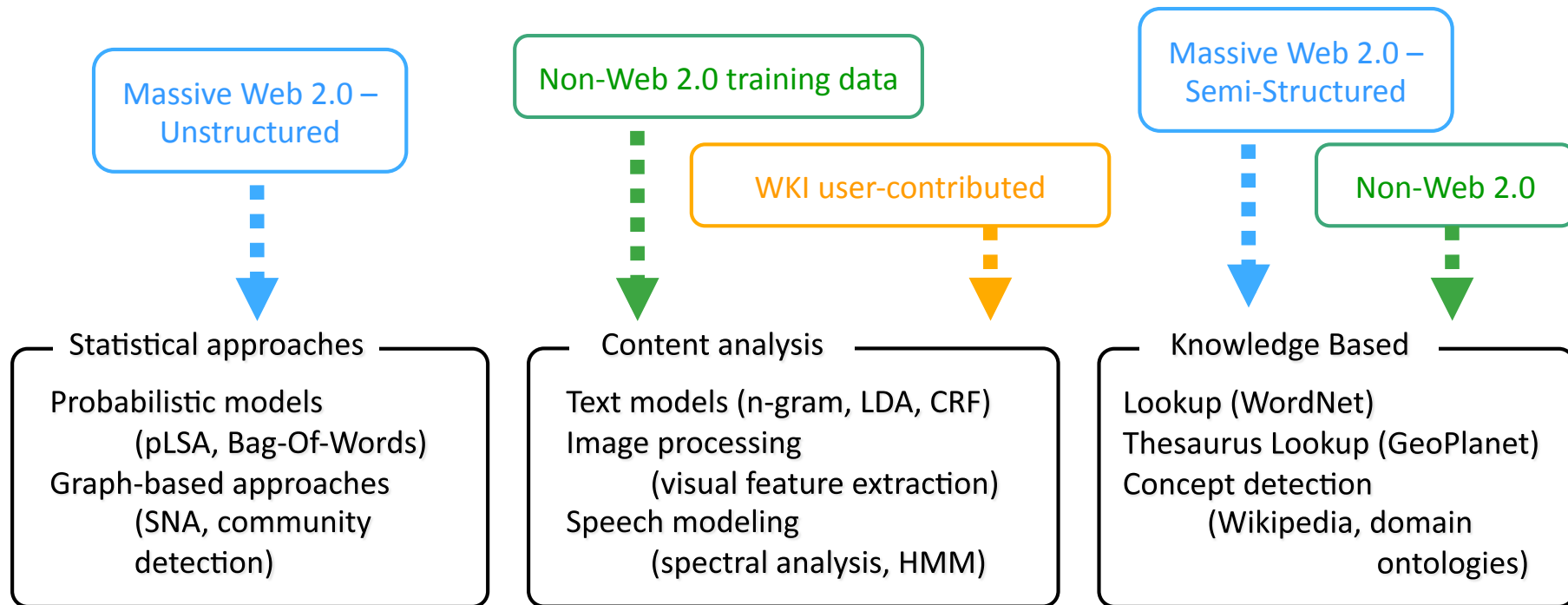- **Huge volume:** Massively produced and disseminated

online → user profiling, method invocation

**WKI user-contributed**

Online content and user actions by WeKnowIt users. It is mainly used for triggering WeKnowIt services and for providing context to them, e.g. user profile, input content to be used as example for querying, etc.

# Technical approach

Variety of approaches depending on content-metadata input.

Massive Web 2.0 – Unstructured

Non-Web 2.0 training data

WKI user-contributed

Massive Web 2.0 – Semi-Structured

Non-Web 2.0

**Statistical approaches**

Probabilistic models
(pLSA, Bag-Of-Words)
Graph-based approaches
(SNA, community
detection)

**Content analysis**

Text models (n-gram, LDA, CRF)
Image processing
(visual feature extraction)
Speech modeling
(spectral analysis, HMM)

**Knowledge Based**

Lookup (WordNet)
Thesaurus Lookup (GeoPlanet)
Concept detection
(Wikipedia, domain
ontologies)

Massive → Collective Intelligence

Multi-Modal (Fusion) →Combined CI
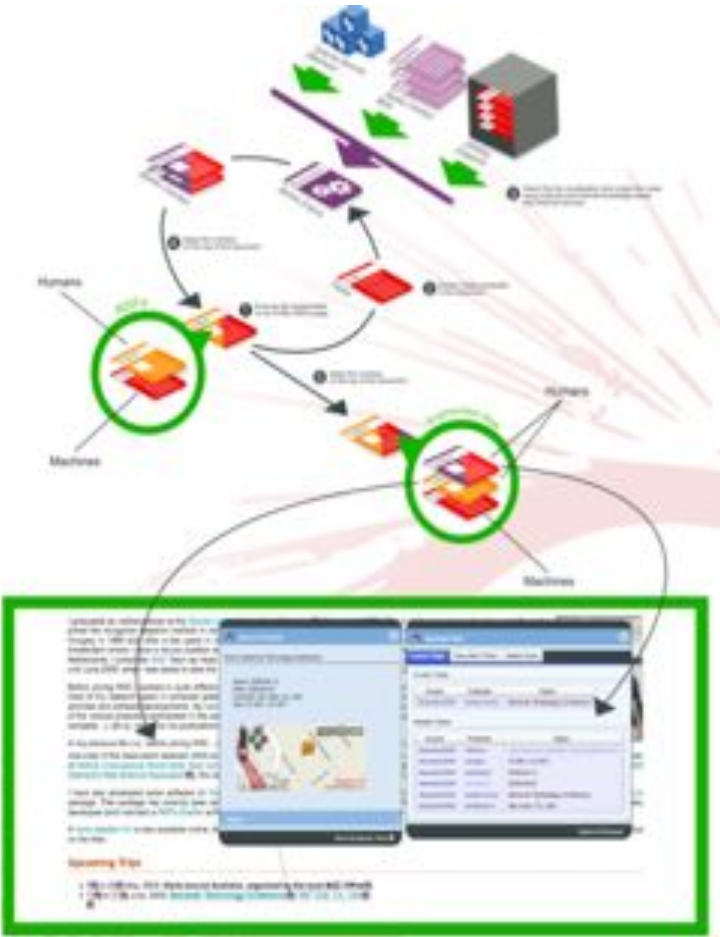
weknowit

# WP1: Personal Intelligence



**Access - Touchpad**

**CURIO Interaction Ontology**

**Verosity**

**(WP1/WP5)**

**Sparks: Semantically Aware Visualization Framework**

**Context – Attention Streams**

weknowit

# WP2: Media Intelligence

**Visual analysis – Localization - Clustering**

**Ranked entities**

**Speech recognition**

**Text Analysis**

**Combined CI Services**

weknowit

co-funded by the European Union

# WP3: Mass Intelligence



visual

tag

visual + tag

| CF | Normalized CF | Structural Similarity |
|---|---|---|
| **Heineken Music Hall** (neighbors: 172) | | |
| Le Zenith (11) | Plaza d. Toros de Valencia (0.25) | Le Zenith (0.6201) |
| Forest National (7) | Le Zenith (0.1078) | Halle Tony Garnier (0.5244) |
| Werchter (5) | Belgrade Fair-Hall (0.0909) | Principe Felipe Arena (0.4878) |
| Lotto Arena (5) | Fonix Hall (0.0667) | Rockhal (0.4672) |
| Oslo Spektrum (4) | Samsung Arena (0.0625) | Werchter (0.4590) |
| Cardiff International Arena (4) | Paradiso (Amsterdam) (0.0526) | Coliseu dos Recreios (0.4590) |
| Skandinavium (4) | Bang Your Head! (0.0455) | Coliseu de Porto (0.4497) |
| Melkweg (4) | Hala Rondo (0.0455) | Killesbergpark (0.4462) |
| **Degree filtered:** France, Paris, Switzerland, Brussels, Madrid, Czech Republic, Lisbon, etc. | **Degree filtered:** France, Switzerland, Paris, Brussels | **Degree filtered:** - |
| **Madame Tussauds** (neighbors: 346) | | |
| Alton Towers (10) | Rock Circus (0.2857) | Historic House Trust (0.3413) |
| Thorpe Park (9) | National Wax Museum (Ireland) (0.1429) | Hudson River Maritime Museum (0.3409) |
| London Eye (8) | The Amsterdam Dungeon (0.1) | Mabee House (0.3366) |
| Chessington World of Adventures (7) | Fort Decker (0.0833) | Johnson Hall State Historic Site (0.3328) |
| Baker Street (5) | Glaspalast (Munich) (0.0714) | Empire State Railway Museum (0.3298) |
| Legoland (4) | Johnson Hall State Historic Site (0.0714) | |
| Natural History Museum (4) | | |
| **Degree filtered:** New York, Victoria and Albert Museum, Buckingham Palace, Westminster Abbey, London | **Degree filtered:** - | **Degree filtered:** - |

| | | | | | | |
|---|---|---|---|---|---|---|
| microsoft | 0.139 | girls | 0.110 | | deportivo | 0.077 |
| teched | 0.136 | show | 0.108 | | sports | 0.076 |
| autumn | 0.132 | sexy | 0.108 | 3gsm | 0.408 | camp |
| student | 0.132 | ficeb | 0.108 | fira | | soccer |
| partners | 0.132 | erotic | 0.107 | fair | | lacoru | 0.055 |
| msp | 0.132 | model | 0.107 | trade | 0.058 | futbol | 0.054 |
| nov | 0.125 | festival | 0.107 | montjuic | 0.051 | liga | 0.052 |
| water | 0.031 | ficeb09 | 0.107 | fountains | 0.037 | deportes | 0.041 |
| bycic | 0.017 | tattoo | 0.105 | amazing | 0.029 | barca | 0.038 |
| food | 0.010 | fishnet | 0.019 | show | 0.027 | depo | 0.036 |
| | | | | night | 0.012 | camprou | 0.036 |
| | | | | wate | 0.011 | | |

graffiti 0.199
streetart 0.199
art 0.199
street 0.199
mtn 0.012
happy 0.012
blameless 0.011
montana 0.009
parr 0.007
spraypaint 0.003

Microsoft students program

teched-09-area

Erotic fairs

GSM conference

Football

weknowit

# WP4: Social Intelligence

**Community administration platform**

**Community analysis tool**

**Community browser**

weknowit

co-funded by the European Union

# WP5: Organisational Intelligence

**Event Model-F**

**Multimedia Metadata Ontology (M3O)**

**dgFOAF Membership Evaluation**

**Log Merger**

# WP6: Architecture & Integration



**Architecture Layers**

**WKI Data Storage**

## 30 services integrated

weknowit

# WeKnowIt Community Detection

# Challenges in Social Media network mining

No prior assumptions about structure:

- Complex & evolving structure
- No possibility for knowing structural features (e.g. number of clusters on a graph) in advance

→ Unsupervised

Scale

- Tens of millions of active users frequently contributing loads of content links + metadata (tags, comments, ratings)
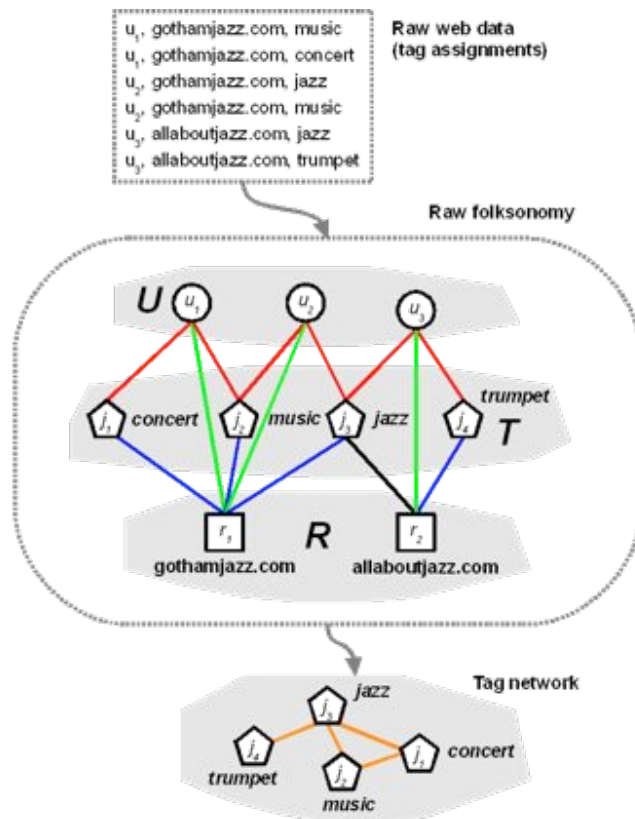
→ Efficient - scalable

Quality

- Spam is very common. Only a portion of user contributions is worth further analysis.
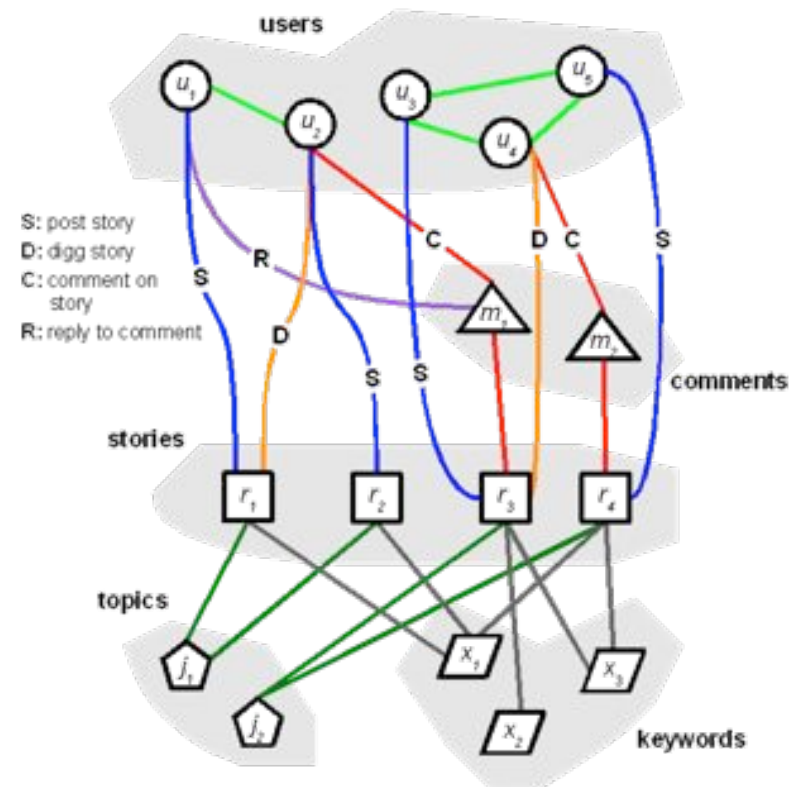
→ Noise resilient

# Examples of Social Media networks

### Folksonomy (Delicious)

Mika, P. (2005) Ontologies Are Us: A Unified Model of Social Networks and Semantics.  Proceedings of the 4th International Semantic Web Conference (ISWC 2005), Springer Berlin / Heidelberg, pp. 522-536

### MetaGraph (Digg)

Lin, Y., Sun, J., Castro, P., Konuru, R., Sundaram, H., and Kelliher, A. (2009) MetaFac: community discovery via relational hypergraph factorization. Proceedings of KDD '09, ACM, pp. 527-536

weknowit

# What is a community in a network?

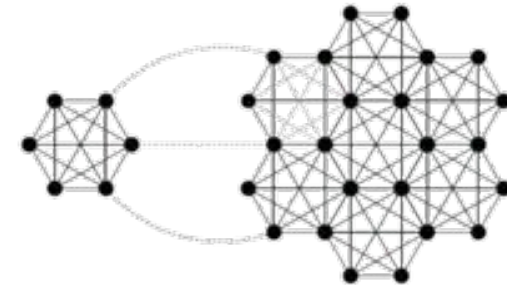Group of vertices that are more densely connected to each other than to the rest of the network.

Multiple definitions to quantify communities:

Fortunato S. (2010) Community detection in graphs. Physics Reports486: 75-174

Global: N-cut, conductance, modularity
Local: Local modularity, $(\mu,\varepsilon)$-cores
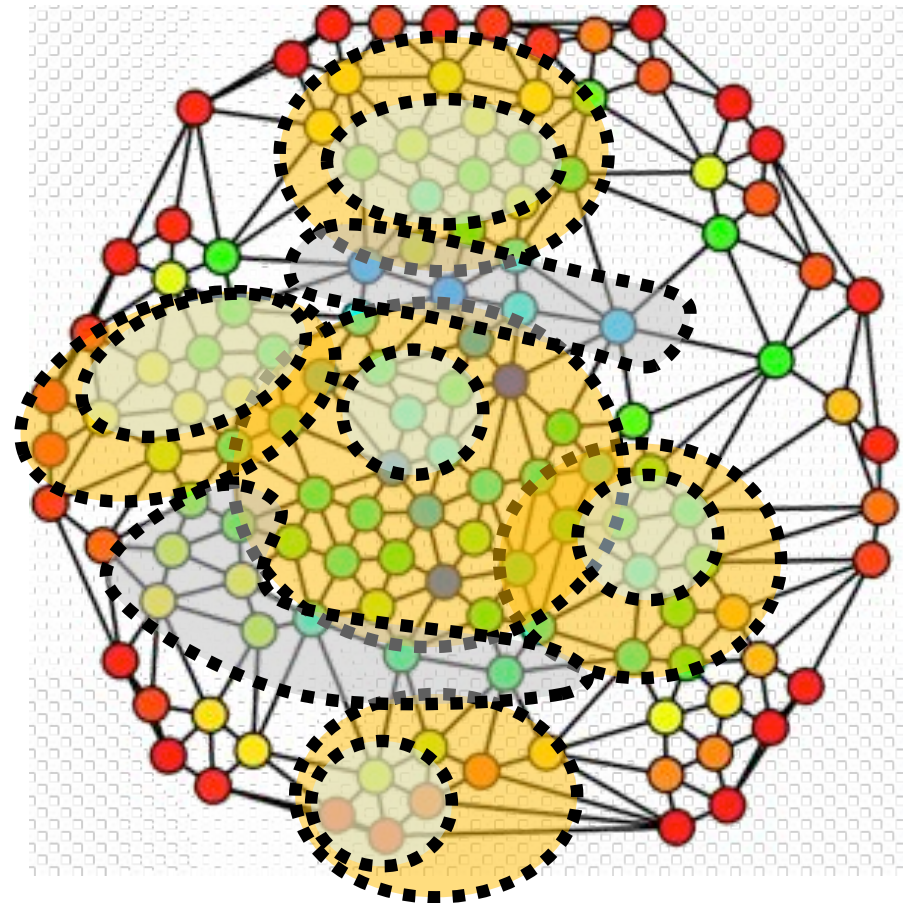Ad hoc: Label propagation, dynamic synchronization

Related to clustering, but: (a) not necessary to know number of communities, (b) computationally more efficient

In WeKnowIt, we focus on local definitions, because of the properties of Social Media networks: efficiency-scalability and noise resilience.

# Approach illustration

Two-step process:

- 1$^{st}$ step:

    $(\mu, \varepsilon)$ – core detection

- 2$^{nd}$ step:
    Local expansion

- 3$^{rd}$ step:
    Characterization of
    remaining vertices as *hubs*
    or *outliers*

# Hybrid Photo Clustering

Goal:

- Group large photo collections into clusters based on how much they are related to each other
- Assist browsing and navigation by means of a map-based application
- Detect landmark and event clusters.

Combine both visual features *and* tags

- Two kinds of similarity (visual and tag networks) are complementary to each other
- Many times one photo has missing tags or is hard to interpret visually
- Graph-based approach - superimpose visual and tag graphs
- Use photo cluster features for classification to landmarks/events

Results

- Higher quality clusters by use of both visual and tag similarity instead of only each one of them.
- Clusters can be used for landmark and event detection.
- Integrated in CSG prototype and ClustTour stand-alone demo.
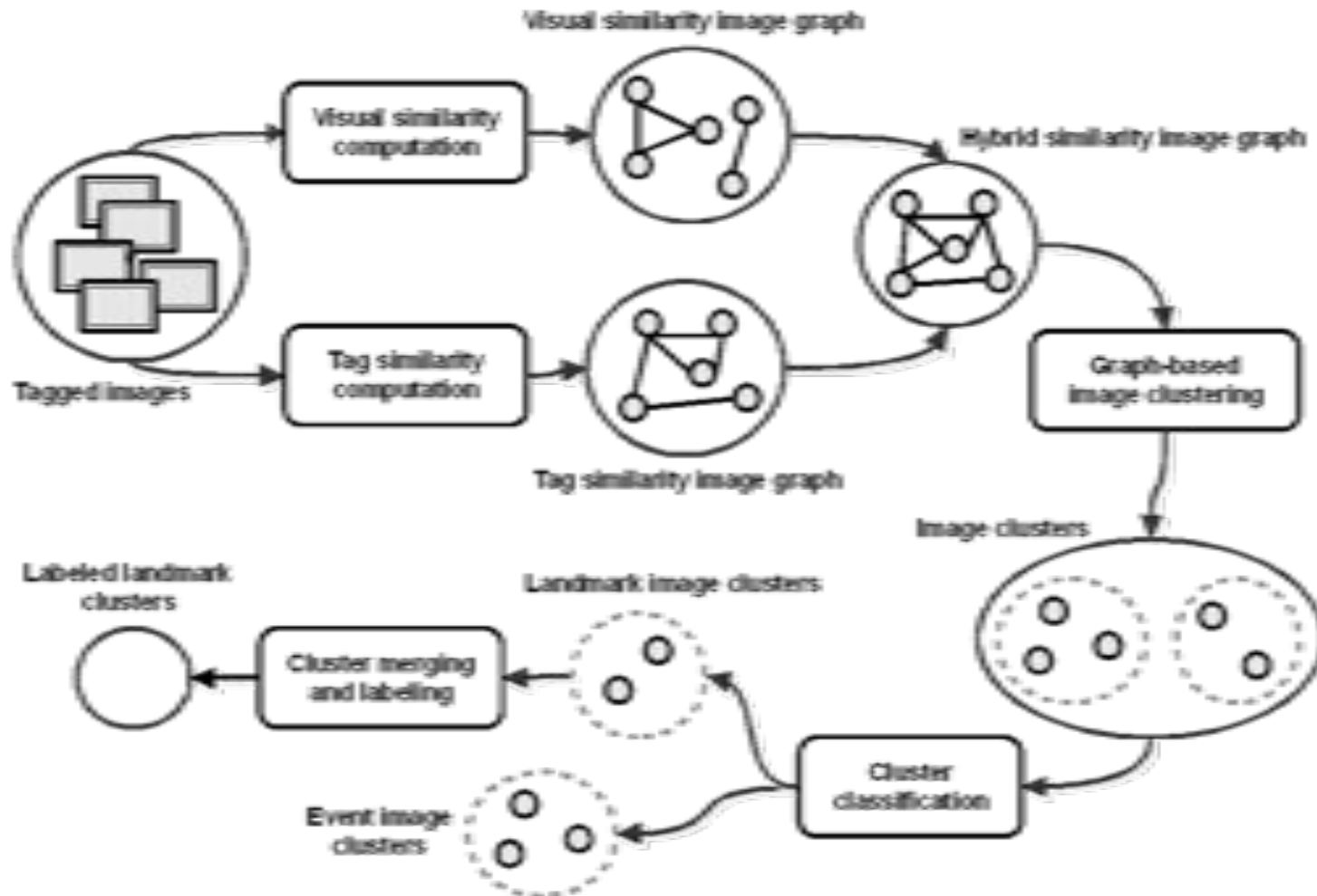
weknowit

# Overview of approach



Visual similarity image graph

Tagged images

Visual similarity computation

Tag similarity computation

Tag similarity image graph

Hybrid similarity image graph

Graph-based image clustering

Image clusters

Labeled landmark clusters

Cluster merging and labeling

Landmark image clusters

Cluster classification

Event image clusters

# Photo clustering results (1)

User study (involving 20 people)

- Users where shown photo clusters and they were asked to judge how relevant the photos of each cluster were related to each other

- Each cluster was produced by different notion of similarity (tag-only, visual-only, hybrid). Obviously, users were not aware of this information

- Hybrid clusters were found to be of superior quality (highest F-measure)

| Algorithm | Precision | Recall | F-measure | $\kappa$-statistic |
|-----------|-----------|--------|-----------|-----------|
| SCAN-VIS | 0.980 | 0.178 | 0.301 | **0.925** |
| SCAN-TAG | 0.910 | 0.197 | 0.323 | 0.688 |
| SCAN-HYB | 0.898 | **0.246** | **0.387** | 0.637 |
| EXP-VIS | **0.985** | 0.178 | 0.301 | 0.895 |
| EXP-TAG | 0.929 | 0.201 | 0.331 | 0.709 |

**weknowit**

# Photo clustering results (2)

Geographic localization of results was also found to be very high.
Most clusters correspond to landmarks or events.



LANDMARKS

EVENTS

co-funded by the European Union

# Sample results: [Visual] vs. [Tag] vs. [Visual + Tag]



VISUAL

TAG

HYBRID

# ClustTour demo: City exploration by means of photo clusters

# Travel demo

co-funded by the European Union

# WKI Grand Travel Challenge
## Barcelona, January 21st 2010

# WKI Grand Travel Challenge
## Barcelona, January 21st 2010

# WeKnowIt Grand Travel Challenge (best of)

**Group Pool**  Discussion  18 Members  Map  Invite Friends

▶ Add something?



From neilreson

From CostieC

From aizenah

From brother.logic

From Pavel.Smrz

From Pavel.Smrz

From Spikos

From dipkris

From Akis_Pag

From SonnyKA

From WeKnowIt Grand...

From pbibos

From WeKnowIt Grand...

From WeKnowIt Grand...

From WeKnowIt Grand...

From WeKnowIt Grand...

From WeKnowIt Grand...

From WeKnowIt Grand...

From YiannisKo

From lopuevo
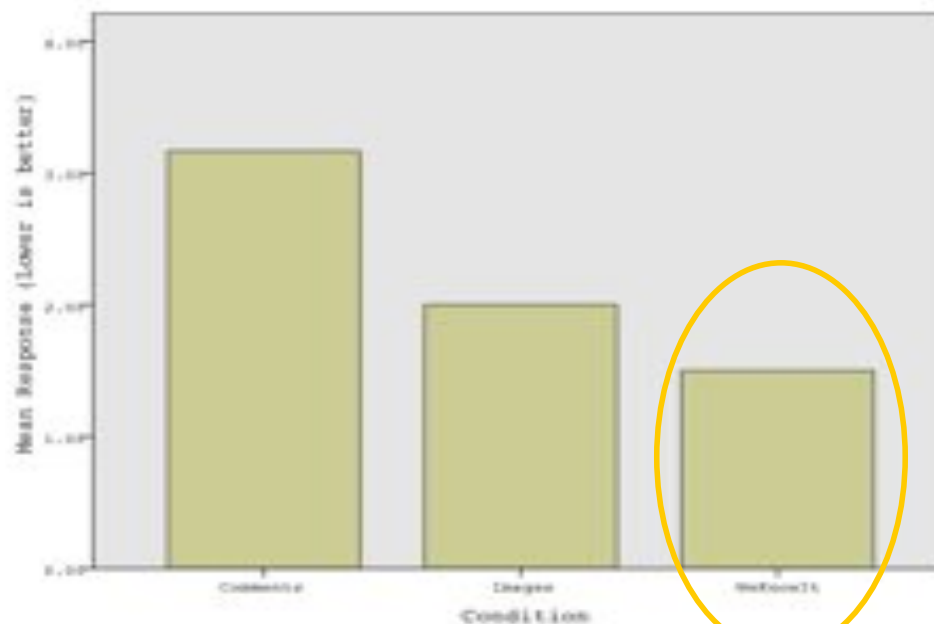
From borkur.net

# Evaluation results

- 5 different evaluation activities
  - 1 ER specific
  - 2 CSG mobile
  - 2 CGS desktop
- 59 users
- 81 system tests
  - ER personell
  - Citizens
  - Users
  - TID personell
  - WKI consortium

*"I found it easy to process all the information available to me"*



Experts

*"I felt I understood what was going on quickly"*

# Content - Emergency Response

- Text
  - Sheffield floods forum messages and posts (thousands)
  - Police & fire service logs
  - ABC news tagged articles (~7000 metadata files)
- Image
  - Flickr images + metadata
    - 136 related to June 2007 Sheffield floods
    - ~27K geo-coded photos around the area of Sheffield
      - 1400 ER images (after WKI clustering)
- Speech
  - 1000+ emergency phone calls on Sheffield flooding event
  - 1000+ voice-tagging events by at least 10 users
  - Fused text+speech dataset

**weknowit**

# Content - Consumer Social Group

- Text
  - Flickr – metadata from geotagged London images (4300+ files)
  - Wikipedia processing
  - GeoPlanet processing
- Image
  - VIRaL (1.2 million geotagged images - 22 European cities)
  - Barcelona meeting dataset - geotagged and tagged (647 images, 1669 tags)
  - 1000 restaurant images
- Social networks
  - Barcelona meeting network of contacts (14 users)

# Research Fields and Issues

- Statistical analysis, machine learning, data mining, pattern recognition, social network analysis
- Clustering
- Graph theory
- Image, text, video analysis
- Information extraction
- Fusion techniques
- Trust, security, privacy
- Performance, scalability
  - speed, storage, power, grids, clouds

# Conclusions

- Collective Intelligence can be extracted by social media

- New applications and services can be developed

- Fusion of multimodal – multisource info remains a challenge

- Scalability, quality, coverage are important issues

**weknowit**