



# Visualització i Gestió de la Informació Científica

[www.dama.upc.edu](http://www.dama.upc.edu)

[www.sparsity-technologies.com](http://www.sparsity-technologies.com)



Josep L. Larriba Pey, larri@ac.upc.edu

# DAMA-UPC: research university group

Our goal



Research and Technology Transfer  
on Managing Very Large Data Volumes



Since 1999

Founded inside UPC

The beginnings



22 researchers and  
developers

The team



Funding,  
agreements and  
collaborators



Our software

\*sparsity  
technologies  
performance in action

Massive graph  
analysis

AURUM: data  
cleansing and fusion



Awards

IBM Faculty Awards (2004, 2009,  
2012)

IBM PhD Award (2004)

CINC prize for novel  
entrepreneurs (2009)

Digital prize for Universities,  
2011



Publications

> 75 research papers  
6n patents

# Visualització i Gestió de la Informació

---

1. Fonts de dades
2. Integració/qualitat de la informació
3. Consultes de valor afegit
4. Automatització de processos
5. Visualització

# 1. Fonts de dades

---

- Bibliogràfiques:
  - Obertes: PubMed, ArXiv, DBLP
  - Propietàries: Scopus, Thomson-Reuters, ACM, IEEE, etc
- Patents:
  - Obertes però pagant
- Xarxes socials:
  - APIs per accedir a la informació social
  - Relació entre persones
- Altra informació:
  - Cordis
  - Proteïnes
  - Productes químics

## 2. Integració/qualitat

---

- Qualitat de les fonts de dades
  - Errors tipogràfics
  - Falta d'informació
- Desambiguació d'autors
  - Gran problema actual
- Citacions incomplertes
  - Molt car de mantenir (Scopus >20K articles diaris)
- Preservar la privacitat de les dades donant accés a les metadades
  
- Solucions possibles
  - Us de tecnologia de grafs i desambiguació
  - Metaintegradors de dades que millorin les dades a partir de la interacció dels usuaris



### 3. Consultes de valor afegit

---

Amb l'ús de paraules clau, resums o metadades en grl.

- Reputació
- Cerca d'obres més valuoses (citacions, reputació dels autors)
- Evolució històrica d'àrees científiques a partir de paraules clau
  
- Solucions possibles:
  - Us de tecnologia de gestió de grafs
  - Consultes avançades

# 4. Automatització de processos

- Cerca i gestió de revisors

Decision Trees For Error-tolerant Graph Database Filtering

graph graph matching database graphs candidate patterns database retrieval retrieval subgraph isomorphism error tolerant matching

Type a search expression like database performance or "conceptual schema"

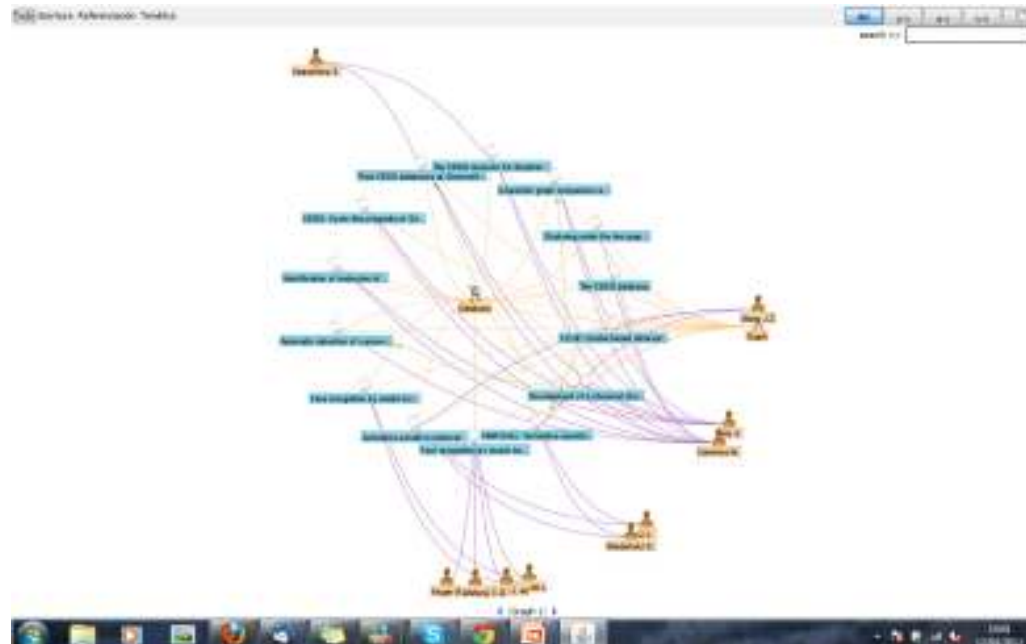
"graph" "graph matching" Search

Show: All Country: Please Select a Country

Candidates	Reviewers	Authors
<b>Hancock S.T.</b> Number of documents: 58 Related documents: 41 Weight: 100.0	Show candidates Assign review	Show candidates Assign review
<b>Escotano E.</b> Number of documents: 27 Related documents: 16 Weight: 24.35	Show candidates Assign review	Show candidates Assign review
<b>Iccano M.A.</b> Number of documents: 70 Related documents: 12 Weight: 16.25	Show candidates Assign review	Show candidates Assign review

## 5. Visualització

- Com fer obvis els problemes causats per les dades de forma visual
- Com mostrar diferents fonts de dades
  - Aportant valor afegit
  - Mostrant que les dades són de fonts diferents
  - Integrant sense molestar als propietaris

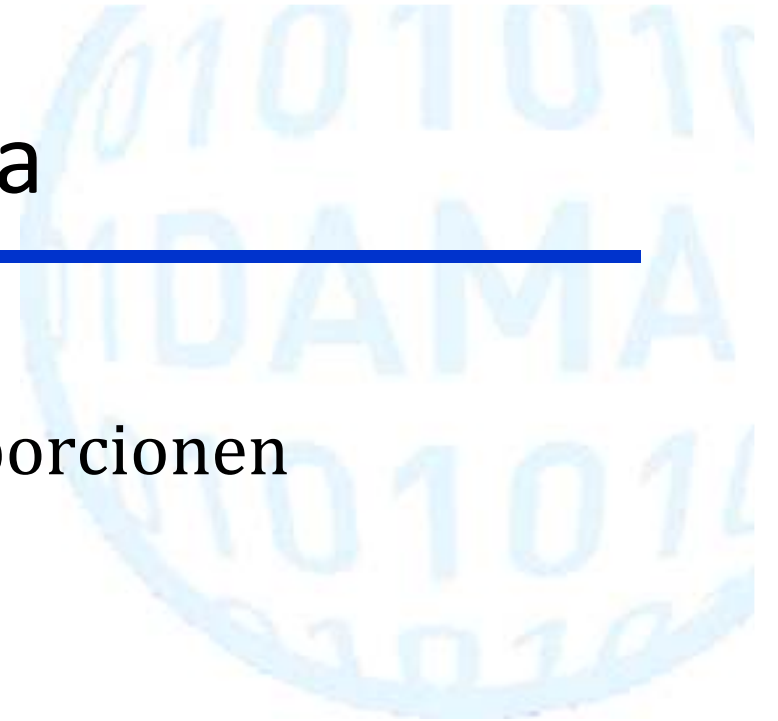




# Solucions al problema

---

- DAMA-UPC i Sparsity proporcionen tecnologia:
  - Gestió de grafs, Dex
  - Reputació
  - Cerca de revisors
  - Suport a l'elaboració de projectes, Sciencea



# Altres preguntes

---

- Donat que les dades amb una mínima qualitat són cares (Scopus, Thomson...), quina qualitat es pot aconseguir amb dades públiques (DBLP, PubMed, etc)?
- Podem aconseguir mètodes alternatius per a donar qualitat a les dades públiques que siguin ràpids, barats (sense intervenció humana) i ens donin fiabilitat?
- Què volem saber de la informació científica? Com volem que es visualitzi, en forma de llistat o en forma gràfica?

# Altres preguntes (cont)

---

- Integrar informació de forma ràpida és possible, i.e. Scopus rep més de 20K articles per dia i els ha d'integrar a la seva base de dades de milions d'articles?
- Que ens aportaria tenir bases de dades integrades amb fons bibliogràfiques, de patents, d'institucions, d'experiments científics, etc?
- Com podriem representar informació integrada de forma visual?

# Any questions?

---



Contact e-mail: [larri@ac.upc.edu](mailto:larri@ac.upc.edu)

DAMA Group Web Site: <http://www.dama.upc.edu>