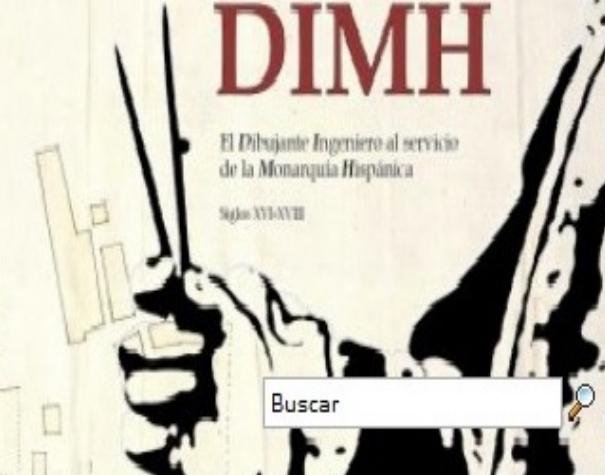


# El dibujante ingeniero

El dibujante ingeniero al servicio de la monarquía hispánica. Siglos XVI-XVIII (DIMH)



## Concept-based Organization for semi-automatic Knowledge Inference in Digital Humanities: Modelling and Visualization

Ángel Castellanos, Ana García-Serrano, Juan Cigarrán  
ETSI Informática, UNED



Natural Language Processing and  
Information Retrieval Group at UNED

nlp.uned.es

# Corpus AGS

- Digital collection of maps, plans and drawings of the *Archivo General de Simancas* (AGS).

[http://www.mcu.es/ccbae/es/consulta/resultados\\_busqueda.cmand?tipo\\_busqueda=mapas\\_planos\\_dibujos&posicion=1&id=30485](http://www.mcu.es/ccbae/es/consulta/resultados_busqueda.cmand?tipo_busqueda=mapas_planos_dibujos&posicion=1&id=30485)

- The data are provided in textual cards (7792 cards), one per each item (map, plan or draw) in the collection.

# Corpus AGS

[http://www.mcu.es/ccbae/es/consulta/resultados\\_busqueda.cmd?tipo\\_busqueda=mapas\\_planos\\_dibujos&posicion=1&id=30485](http://www.mcu.es/ccbae/es/consulta/resultados_busqueda.cmd?tipo_busqueda=mapas_planos_dibujos&posicion=1&id=30485)

Sección:	Material gráfico AGS
Número de control:	BAB20100018941
Autor:	Huet, Luis 
Título:	Plano y perfiles que manifiestan el estado en que se alla la Real obra de Fuerte-Príncipe en 30 de junio de 1779 [Material gráfico no proyectable] / [rúbrica] Luis Huet
Área de datos:	Escala [ca.1:816], 200 varas reales [=20,6 cm]
Publicación:	Habana, 30 de junio de 1779
Descripción física:	1 plano : ms., col. ; 36,5 x 47 cm
Notas:	Referencias: Mapas, planos y dibujos (Años 1503-1805). Volumen I : p. 576 Tinta y colores a la aguada ocre y encarnado. Explicación con clave alfabética Manuscrito sobre papel. AGS. Secretaría de Guerra, Legajos, 03222. Acompañía a carta y relación delas obras de don Luis al Conde de Ricla de 1 de julio y de la misma fecha del plano
Materia / geográfico:	Fortificaciones-La Habana-Dibujos  La Habana-Edificios, estructuras, etc.-Dibujos 
Género / forma:	Dibujos de arquitectura-España-S.XVIII 

Ficha: 176927

# Corpus AGS

- Spanish Project (HAR2012-31117)

**El dibujante ingeniero al servicio de la monarquía hispánica. Siglos XVI-XVIII (DIMH)**

<http://dimh.hypotheses.org/>

Main goals of the project are:

- Knowledge organization of the cards contents
  - Identification of data relationships
  - Visualization of the results
- In order to support the research of the historian researchers of the project

# Corpus DIMH

- The data provided in the textual cards (7792) have been pre-processed in order to:
  - Convert the cards from RDF:DC to XML.
  - Identify the named entities
  - Identify the nominal groups
  - Identify the lemmas
- Not supervised process



# Corpus DIMH: Enrichment using Linguistic Resources (not supervised)

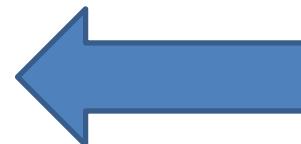
```
<Ficha id="176927">
<Tematica> La Habana-Edificios,
estructuras, etc.-Dibujos</Tematica>
<Materia>-La Habana</Materia>
<Materia>-España</Materia>
<Formato>image/jpeg</Formato>
<Idioma>spa</Idioma>
<Materia>La Habana</Materia>
<Materia>-S.XVIII</Materia>
<Notas> AGS. Secretaría de Guerra,
Legajos, 03222. Acompaña a carta y
relación delas obras de don Luis Huet al
Conde de Ricla de 1 de julio y de la
misma fecha del plano</Notas>
<Creador>Huet, Luis</Creador>
<Publicacion>1779</Publicacion>
<Notas> Manuscrito sobre
papel.</Notas>
<Materia>Cuba-La Habana</Materia>
... </Ficha>
```

```
<nes>La_Habana Fortificaciones-Dibujos_Dibujos
Dibujos Luis_Huet Huet AGS Cuba España
Real_Sociedad Real Conde_de_Ricla Fuerte-
Príncipe Luis_La_Habana </nes>
<nes_person>Luis_Huet Huet Fuerte-
Príncipe</nes_person>
<nes_organization>Real_Sociedad
Real</nes_organization>
<nes_location>La_Habana Cuba España
</nes_location>
<nes_misc>Fortificaciones-Dibujos_Dibujos
Dibujos AGS Conde_de_Ricla
Luis_La_Habana</nes_misc>
<lemas> plano y perfil que manifestar el estado en que
se allá el Real obrar de Fuerte - Príncipe en 30 de Junio
de 1779 ... </lemas>
<sintagmas> plano_perfil real fuerte_príncipe
junio_1779_material_grafico proyectable_huet
luis_habana_edificio dibujos
... </Ficha>
```

**Ficha: 176927**

# Main Goal

- End-Users interested in **hidden data** and/or **data relationships**
  - Knowledge Extraction and
  - Latent Organization Discovering
- Available techniques and technologies
  - Ontologies and LOD resources
  - Quantitative approaches



# Motivation of the Work

Probabilistic techniques as Latent Dirichlet Allocation (LDA), have become almost a standard for the **Content Organization** in the Digital Humanities (DH) (Meeks & Weingart 2012, Yin et al. 2013, ...)

It suffers from:

- the need to fix the number of topics to be detected or
- the non-trivial interpretation by the humanists.

Our approach is based on the Formal Concept Analysis (FCA) that takes advantage of a well-founded mathematical background

# FCA for Content Organization

- **Formal Concept Analysis is a priori solution because:**
  - Allows the organization of objects according to their shared attributes into a generalization-specification relationship.
  - Organize the latent structure according to the shared terms in thematic-based concepts.
  - Generate a hierarchical structure (Lattice) allowing its navigation and visualization.

But, We need to develop our own **framework**

# FCA at a Glance

**Formal Context** is a triple  $(G, M, I)$ , where

- $G$  is a set of (formal) objects
- $M$  a set of (formal) attributes
- $I \subseteq G \times M$  is the incidence relation

i.e.  $(I \subseteq G \times M)$  means that the object  $g$  has the attribute  $m$ .

**A Formal Context**  
is an incidence  
matrix that indicates  
whether or not an  
attribute is related to  
an object.

The main concept of FCA is the formal context pair  $(A, B)$  where  $A \subseteq G$  is a set of objects (extent) and  $B \subseteq M$  is the minimal set of attributes (intent) shared by all the objects in  $A$ .

# FCA at a Glance

*Formal Context* is a triple  $\mathcal{K} := (G, M, I)$ , where

- $G$  is a set of (formal) objects,
- $M$  a set of (formal) attributes and

i.e. ( $I$ )  
object  $g$  has

**A Formal Concept**  
is a set of objects  
sharing a set of  
attributes

The main concept of FCA is the **Formal Concept**, a pair  $(A, B)$ , where  $A \subseteq G$  is a set of objects (*extent*) and  $B \subseteq M$  is the maximal set of attributes (*intent*) shared by all the objects in  $A$ .

# FCA at a Glance

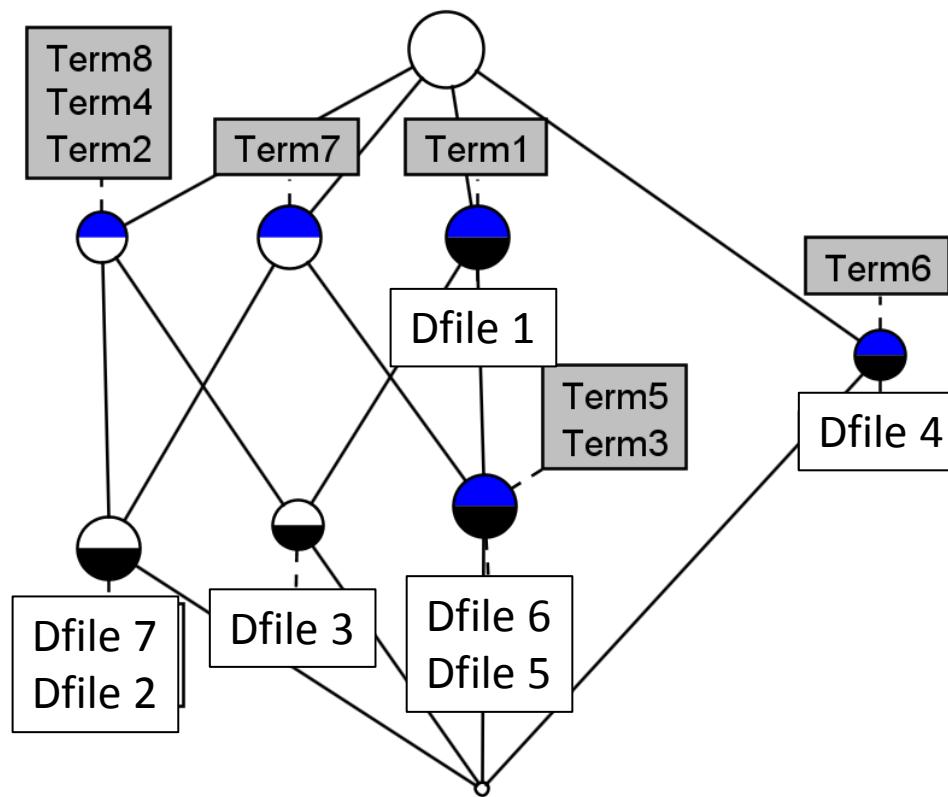
*Formal Concepts* can be formally ordered in a **subconcept-superconcept-relation** according to their extents:

where (C,  
conversely,  
specific than  
to be a *latt*

Since concepts  
displayed in

and,  
B) is more  
e proven

be



# FCA application to DIMH corpus

Main goals:

- how FCA performs for a content modelling task and
- whether the obtained model infers new knowledge from the original data.

The steps in the FCA application for content modelling are:

**1. Information Extraction:** It extracts the data in the processed AGS files (DIMH corpus)).

Selected Metadata: publication, reference, notes, named entities, topic, material and title.

# FCA application to DIMH corpus

- 1. Information Extraction**
- 2. Formal Context Generation:** The *objects* are the AGS files and the *attributes* are the selected terms related to the files.
- 3. Formal Context Reduction:** The formal context generated in the previous step includes redundant and not valuable information.

The formal context reduction takes only those features in the *formal context* that allow the identification of more relationships among the cards, avoiding information loss.

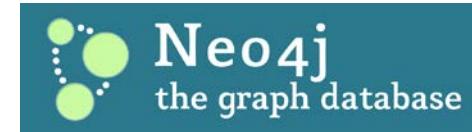
# FCA application to DIMH corpus

- 1. Information Extraction**
- 2. Formal Context Generation**
- 3. Formal Context Reduction**
- 4. FCA execution:** formal concepts and its hierarchical structure generation.

The FCA algorithm is based on a self-implemented version of the Next Neighbourhood algorithm (Carpinetto & Romano 2004).

# Applied Research: Own ToolKit

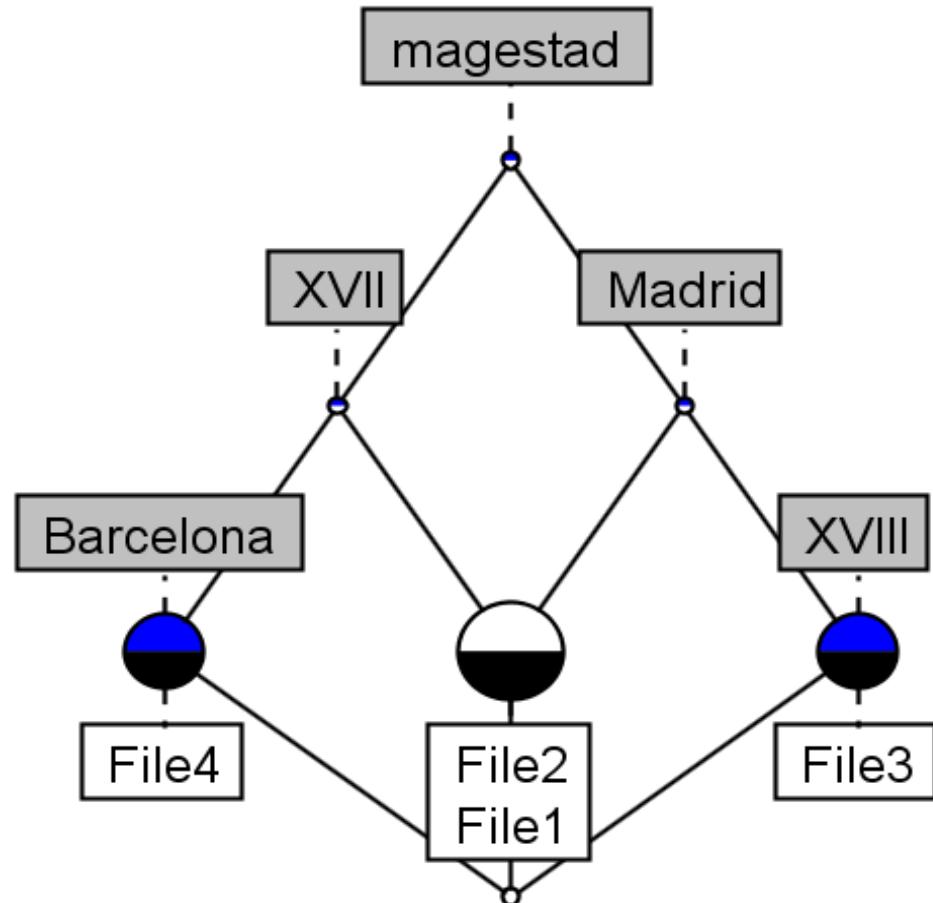
- Modelling step
  - Library for KLD Pre- Filtering
  - FCA Algorithm (Jbrainded, la4j for huge matrix)
  - Graph Representation  
(Neo4j, graph DB based on Lucene)
  - Navigation Algorithms (Code for graph navigation)
- Interface to facilitate the interaction within the Lattice  
(Alpha version using Google Web Toolkit)
- Tools for Lab experimentation (refinement, evaluation)



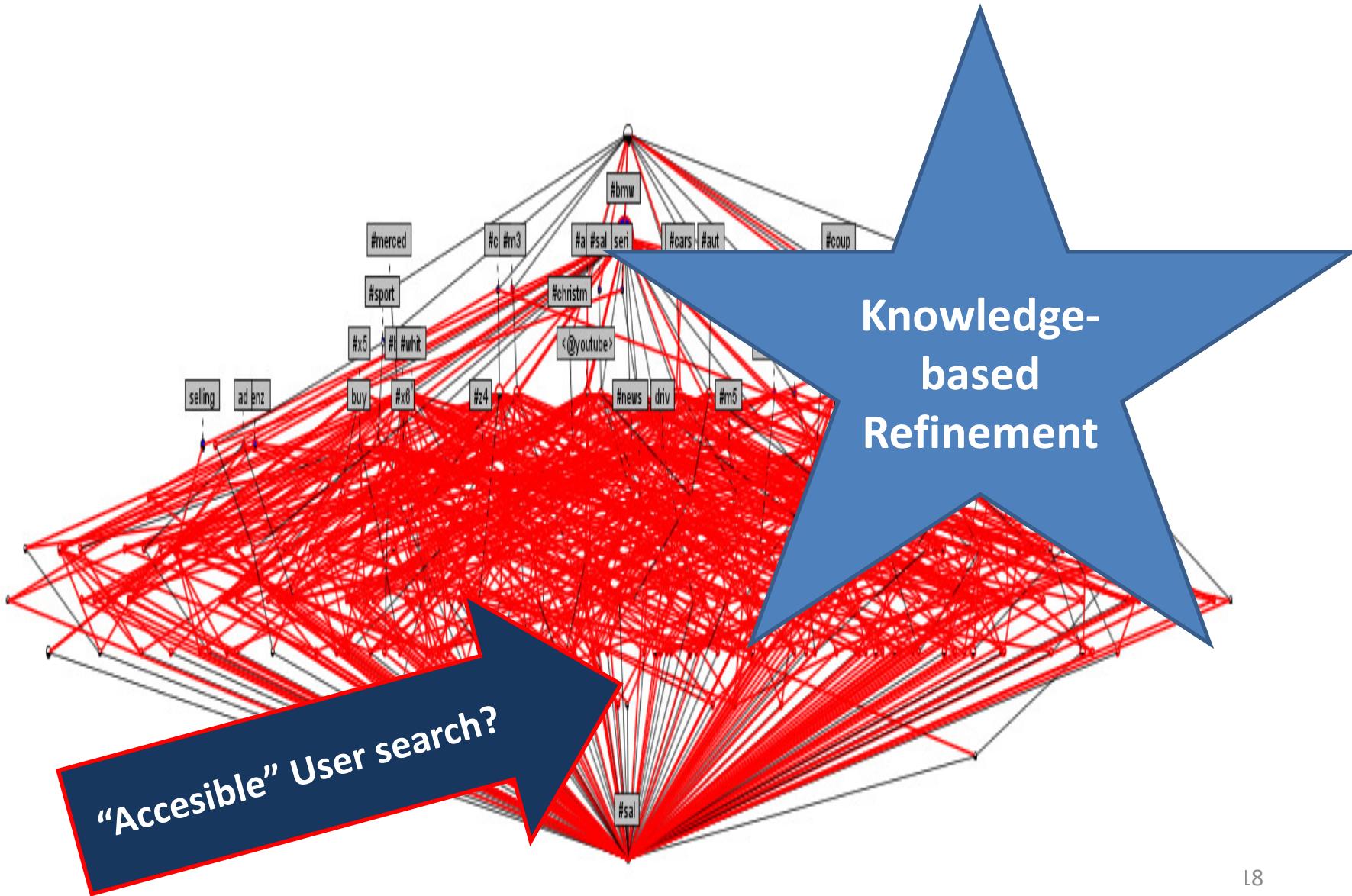
# FCA-DIMH at a Glance

FCA components:

- Objects = DIMH Cards
- Attributes = terminology + special features
- Formal Concepts = similar DIMH cards together according to their features
- Structure = Formal concepts ordered in a hierarchical structure (lattice)



# FCA-DIMH First Results



# First Refinement: Attribute Selection

To identify the “interesting” terms, it is applied the Kullback-Leibler Divergence (Kullback & Leibler 1951).

KLD analyses the probability distribution of an attribute over a document and a collection of documents, **to identify those attributes that best represent a document by differentiating it from the rest in the collection.**

	Files	Attributes	Relationships	Formal Concepts
BEFORE KLD	7,792	130	+ 32,000	+ 25,000
AFTER KLD	7,792	103	29,874	17,501

# FCA -DIMH Browser

<http://albali.lsi.uned.es/DIMHDemo-2/>

Busqueda

Mapas

[Buscar Concepto Formal](#) [Buscar Contenido](#)



[mapas] (846)
[mapas, galicia] (100)
[andalucía, mapas] (80)
[madrid, mapas] (62)
[mapas, hacienda] (37)
[mapas, barcelona] (119)
[mapas, ceuta] (43)
[mapas, valencia] (31)
[mapas, juan] (26)
[mapas, papeles] (17)
[mapas, expediente] (11)
[mapas, domingo] (10)
[mapas, capitán] (7)
[mapas, incluido] (50)
[mapas, duque] (45)
[mapas, acompaña] (29)
[mapas, santa] (14)
[mapas, azul] (12)

- [mapas, costa] (19)
- [mapas, conde] (17)
- [mapas, carlos] (9)
- [mapas, nāº] (7)
- [mapas, cádiz] (39)
- [mapas, españacastilla] (35)
- [mapas, sevilla] (32)
- [mapas, ciudad] (26)
- [mapas, granbretaña] (22)
- [mapas, paísvasco] (21)
- [mapas, castilla] (8)
- [mapas, magestad] (6)
- [mapas, marquésdelaensenada]
- [mapas, italia] (38)
- [mapas, pamplona] (34)
- [mapas, perfil] (27)
- [mapas, argelia] (26)
- [mapas, real] (23)
- [mapas, aragón] (16)**

# FCA – DIMH Contents Browser

- [mapas, aragón] (16)
  - [mapas, aragón, castillo] (5) 
  - [madrid, mapas, aragón] (4)
  - [mapas, aragón, febrero] (4)
  - [mapas, aragón, ciudad, fuerte, febrero] (2)

**Perfil del Castillo d Monzón cortado por los puntos E y F de Poniente a Oriente [Material cartográfico] Monzón (Huesca). Fortificac**

**1797**

**Idioma: spa**

España-Aragón-Huesca (Provincia)-Monzón Manuscrito sobre papel. AGS. Secretaría de Guerra, Legajos, 05866. Con oficio del marqués de Alos a Juan Manuel Alvarez Zaragoza (Años 1508-1962). Volumen II : p. 287 Perfil del Castillo d Monzón cortado por los puntos E y F de Poniente a Oriente [Material cartográfico] Monzón (Huesca). Fortificaciones Zaragoza Juan\_Manuel\_Alvarez\_Perfil\_del\_Castillo\_Monzón Huesca Aragón España Zaragoza Mapas\_Referencias Mapas Perfil\_del\_Castillo\_Alos Juan\_Mon...

**Jaca (Huesca). Fortificaciones. Planos. 1642 Planta del Castillo de Jaca [Material cartográfico]**

**1642**

**Idioma: spa**

España-Aragón-Huesca (Provincia)-Jaca Manuscrito sobre papel. Tinta AGS. Guerra y Marina, Legajos, 01454. Acompañía a carta del Marqués de Távara a S. M., Zaragoza, 26 de febrero de 1642. Volumen I : p. 605 Jaca (Huesca). Fortificaciones. Planos. 1642 Planta del Castillo de Jaca [Material cartográfico] Huesca Jaca Mapas\_Referencias Mapas Aragón España Zaragoza Zaragoza Simón\_Manuscrito Huesca Jaca Aragón España Zaragoza Mapas\_Referencias Mapas Cornacholo Marqués\_de\_Távara Planta\_del\_Castillo\_de\_Jaca Simón\_Mar...

**Planta del castillo de Canfranc con los reparos que hay que hacer en Él [Material cartográfico] Canfranc (Huesca). Fortificacio**

**1617**

**Idioma: spa**

España-Aragón-Huesca (Provincia)-Canfranc Tinta y color amarillo. Con rotulación Manuscrito sobre papel. AGS. Guerra y Marina, Legajos, 00823. Con carta de Felipe de Sotomayor a su hermano don Juan de Sotomayor (Años 1508-1962). Volumen II : p. 96 Planta del castillo de Canfranc con los reparos que hay que hacer en Él [Material cartográfico] Canfranc (Huesca). Fortificaciones. Planos. 1617 Canfranc Huesca Aragón España Mapas\_Referencias Mapas Felipe\_de\_Soria

**Jaca (Huesca). Fortificaciones. Planos. 1645 Planta del castillo de Jaca [Material cartográfico]**

**1645**

**Idioma: spa**

España-Aragón-Huesca (Provincia)-Jaca Manuscrito sobre papel. Dibujo a plumilla. Con rotulación AGS. Guerra y Marina, Legajos, 01594. Con carta de D. Luis Carrillo de Toledo a su sobrino don Francisco de Toledo (Años 1503-1805). Volumen I : p. 605 Jaca (Huesca). Fortificaciones. Planos. 1645 Planta del castillo de Jaca [Material cartográfico] Jaca Huesca Mapas\_Referencias Mapas Luis\_Carrillo\_de\_Toledo Toledo Mapas\_Referencias Mapas AGS Castillo\_de\_Jaca Luis\_Carrillo

**Berdún (Huesca). Mapas generales. 1592 Berdún (Huesca). Fortificaciones. Planos. 1592 Planta de Verdun ; Planta del castill**

**1592**

**Idioma: spa**

España-Aragón-Huesca (Provincia)-Berdún Tinta. Con rotulación Manuscrito sobre papel AGS. Guerra y Marina, Legajos, 00356, 181. Con carta de D. Alonso de Vargas al Rey (Años 1503-1805). Volumen I : p. 1000 Berdún (Huesca). Mapas generales. 1592 Berdún (Huesca). Fortificaciones. Planos. 1592 Planta de Verdun ; Planta del castillo de Berdún Huesca Mapas\_Referencias Mapas Mapas Verdun AGS Acosta Aragón España Jaca Tiburcio Vargas D.\_Alonso Hernando\_de\_Spannocchi Hernando\_de\_Spannocchi Berdún Verdun Acosta Tiburcio D.\_Alonso

# User Evaluation

- End-Users Interested in already known Topics/A priori taxonomy of concepts
  - Authors (cited)
  - Measures Units
  - Dates, Places,...



**Next solution:**

**take only the terminology defined by experts!**

# 2nd Refinement: Taxonomy

The inclusion of experts terminology (taxonomy) in the FCA execution lead to a model organizing the data according to it.

¿It is reduced the number of formal concepts and relationships?

	Files	Attributes	Relationships	Formal Concepts
BEFORE KLD	7,792	130	+ 32,000	+ 25,000
	7,792	36	13,719	1,197

# Research Support

- By means of the developed visualization, the experts are allowed to explore the data and the inferred relationships, drawing new conclusions about the contents based on their expertise.
- In addition, **associations in the data can be automatically discovered from the lattice structure.**  
Hypothesis: This Knowledge Extraction offers to the experts **the starting point for a deeper analysis** of the discovered implications.

# Automatic Association Rules

Automatic knowledge inference from the FCA lattice is carried out in two steps:

- Find the most common feature by the FP-Growth algorithm (Han et al. 2004).
- Find the association rules related to these frequent features by algorithm (Agrawal et al. 1994).

	# Features	# Rules
Named Entities	22	29
Taxonomy	102	133
All	27	121

# Results (not yet by the historians)

- First, the rules inferred from the named entities are mostly related to well known locations.
- By means of the taxonomy, more specific rules are obtained.

**plaza → paper plan**

- support: 0.01196 (87/7274)
- confidence: 0.12850
- Finally, by using all the information, it seems that new unknown information is offered by the rules.

**representation\_system material graphic → plans**

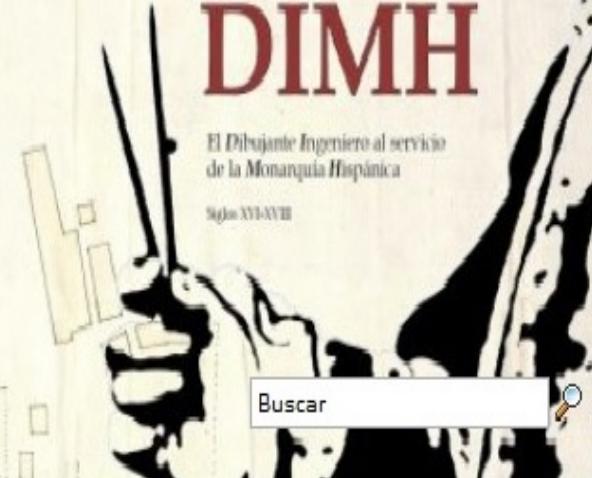
- suport: 0.7066221765913757 (5506/7792)
- confidence: 0.8111373011196229

# Next Steps

- New Visualization Metaphor
- Evaluation of the quality of the results
  - Laboratory-based Evaluation???
  - User-based Evaluation

# El dibujante ingeniero

El dibujante ingeniero al servicio de la monarquía hispánica. Siglos XVI-XVIII (DIMH)



<http://dimh.hypotheses.org/>

**Thanks!!! to Alicia Cámara, DIMH  
Project Responsible, Departamento  
de Historia del Arte, UNED**



**Ana García-Serrano,  
Ángel Castellanos,  
ETSI Informática, UNED**

