

INDICATE

International Network for a Digital Cultural Heritage e-Infrastructure



**Introduction to geocoded cultural content (GCC):  
framework, use cases, geoparsing**

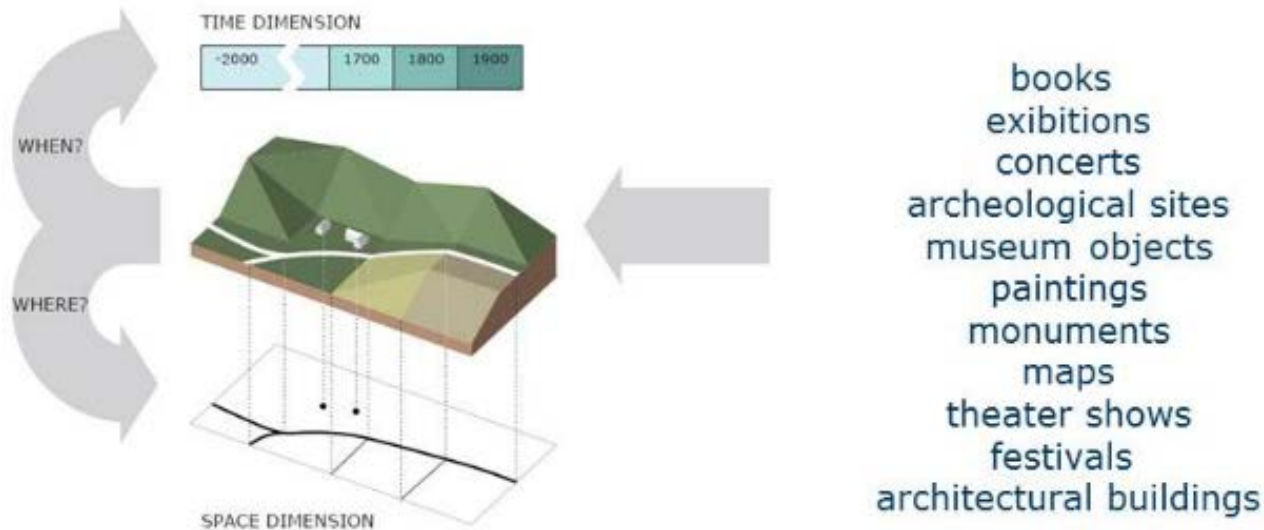
Franc J. Zakrajsek

scientific coordinator of Indicate for Slovenia, Slovenia

International workshop: Geocoded cultural content  
Ljubljana, 7<sup>th</sup> February 2012

## What is GCC ?

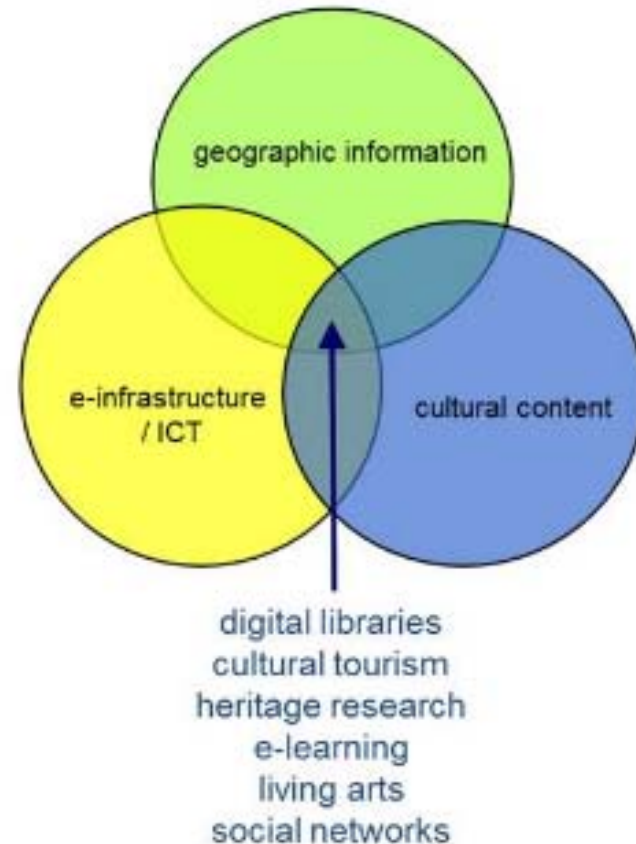
### cultural objects and events



Geographic location is one of the most important aspect of information for every cultural heritage item. A formalized location attribute (e.g. geocode or geographical coordinates) will significantly enhance the power of searching, visualization, analysis of the content.

## What we are doing ?

**The case study Geocoded cultural content:** first part reviews the current approaches and new R&D on geocoding of cultural content in digital libraries, cultural tourism, heritage, e-learning, living arts and other cultural areas. Main area of the research will identify the possibilities and benefits of using e- infrastructure. The focus will be primarily on cloud and grid computing and data infrastructures when dealing with geocoded digital cultural content. The last part of the research provides and summarizes the testing of geoparsing and geotagging e-services in digital culture and recommendations for content providers.





## Concepts and framework of GCC



**context of use**



**spatial accuracy**



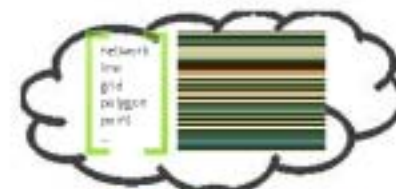
**linking open data**



**digital object types**



**standards**



**geo feature types**



**geocoding**



**devices**



**e-infrastructure**

| Category                 | Number of items | URI                    | URI class              |
|--------------------------|-----------------|------------------------|------------------------|
| Article                  | 24              | http://www.nytimes.com | http://www.nytimes.com |
| Geography                | 11              | http://www.nytimes.com | http://www.nytimes.com |
| Organization             | 24              | http://www.nytimes.com | http://www.nytimes.com |
| Publication              | 27              | http://www.nytimes.com | http://www.nytimes.com |
| City/State               | 11              | http://www.nytimes.com | http://www.nytimes.com |
| Web content              | 11              | http://www.nytimes.com | http://www.nytimes.com |
| Non-geographical content | 24              | http://www.nytimes.com | http://www.nytimes.com |

community  
New York Times  
data sets  
RDF mapping  
linked geo-data  
...

The Linking Open Data cloud diagram

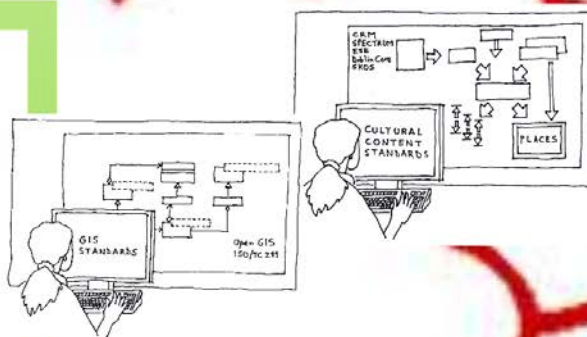


Geocoded articles: Adding the geographical coordinates to published articles; example: Linked Open Data of New York Times, connected with Geonames ontology.



# linking open data

CoreData  
OGC  
CRM  
Spectrum  
ISO/TC 211  
...



**standards**

# Archeological and architectural GCC

## Architectural / archeological heritage

### Definition

Architectural and archaeological heritage refers to a place, locality, natural landscape, settlement area, architectural complex, archaeological site, or standing structure from inventories, management, restoration, ...

### Examples

National Heritage List for England  
National Heritage Register Netherlands  
National Register of Denmark  
Heritage Register Bayern - Nürnberg

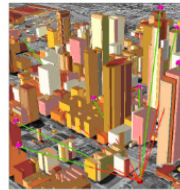
### Geographical information

Geographic information in archaeological / architectural sectors is used when capturing data, management the repositories, and processing and displaying data on the maps. The level of detail goes to individual site or object.

### e-Infrastructure

Appropriate tasks for grid computing:

- Risk scenarious simulations
- 3D visualisation
- Spatial statistics
- Spatial analyses



### Potentials of use grid computing for caching

| Steps  | Software and hardware  | Estimated time                             |
|--|--|--|
| Experiment<br>Caching area: 152 km2 scale 1:76 (approx. 2D 1:1000)<br>Tiles: 512x512 pixels (finally 104.000 tiles (3/4 tilov )) | ArcGIS Server<br>2x E5450 3GHz (8 threads)<br>32GB memory  | Caching time: 77 minut.                    |
| Generalization for the world<br>mainland: 148.429.000 km2, million times larger area than in experiment                          | ArcGIS Server<br>2x E5450 3GHz (8 threads)<br>32GB memory  | Estimated 77 millions minutes or 146 years |
| Google   | <ul style="list-style-type: none"> <li>• In 2002; upwards of 15,000 servers</li> <li>• A 2005 estimate by Paul Strassmann has 200,000 servers claimed this number to be upwards of 450,000 in 2006 900.000 (2011)</li> </ul> |  |

## Registers of immovable CH

### National Heritage List for England



Online database enables searching for listed buildings, scheduled monuments, protected wreck sites, registered parks and gardens, registered battlefields in classic way and on a map.

### National Heritage Register Netherlands



Online heritage register of national monuments (over 60.000) is split per province is searchable in classic way and on a map.

### National Register of Sites and Monuments Denmark

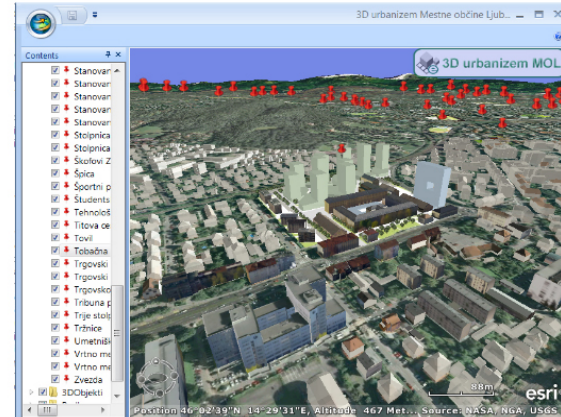


Online register of all known sites, monuments and archaeological finds (165,000 sites including shipwrecks and submarine Stone Age settlements)

### German Heritage Register Bayern - Nürnberg



Nürnberg displays part of the Bavarian monument list. Monument searching and viewing on a map is part of a city plan internet application.





## Registers of immovable CH

### National Heritage List for England



Online database enables searching for listed buildings, scheduled monuments, protected wreck sites, registered parks and gardens, registered battlefields in classic way and on a map.



### National Heritage Register Netherlands



Online heritage register of national monuments (over 60.000) is split per province is searchable in classic way and on a map.



### National Register of Sites and Monuments Denmark



Online register of all known sites, monuments and archaeological finds (165,000 sites including shipwrecks and submarine Stone Age settlements)



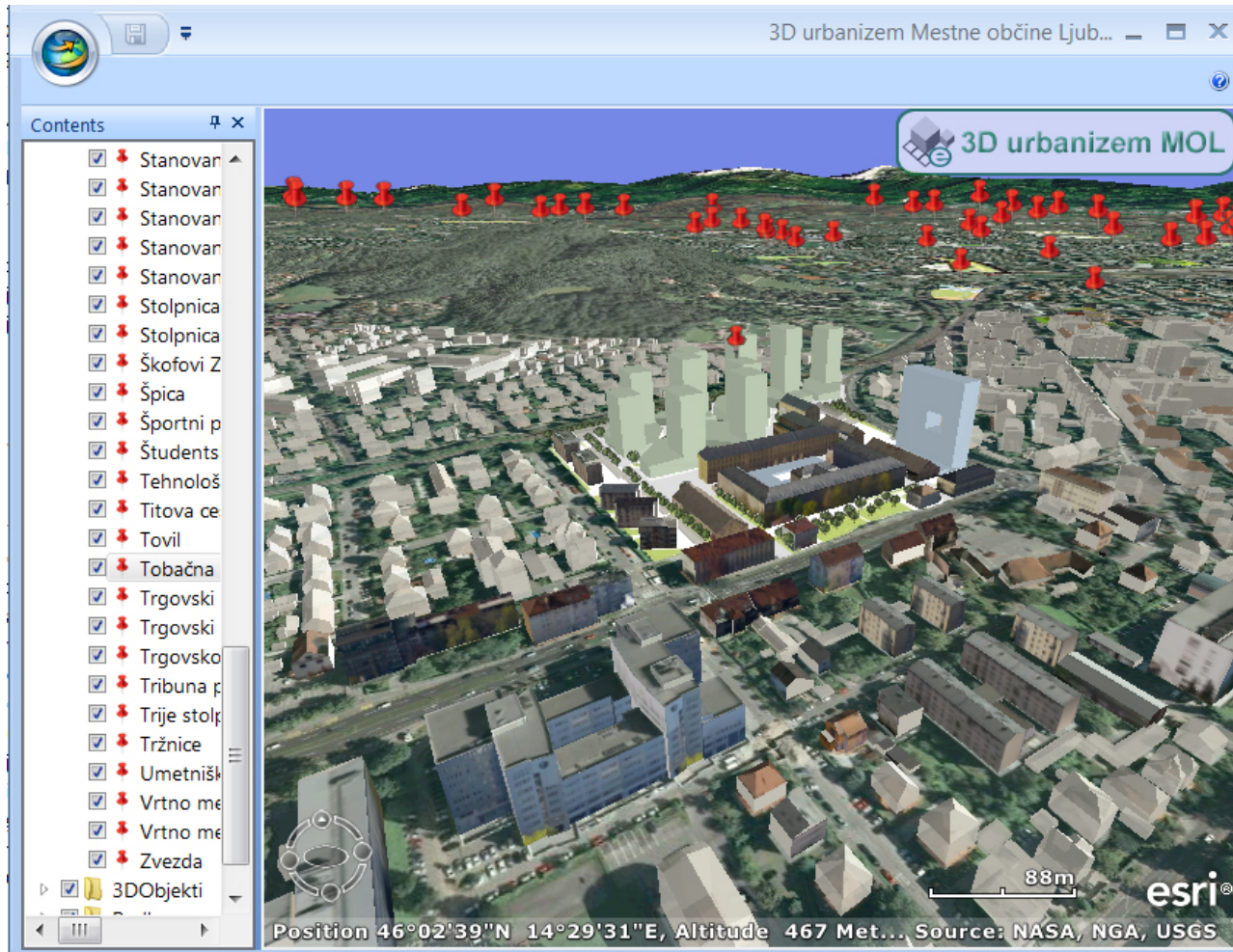
### German Heritage Register Bayern - Nürnberg



Nürnberg displays part of the Bavarian monument list. Monument searching and viewing on a map is part of a city plan internet application.







### Potentials of use grid computing for caching

| Steps  | Software and hardware  | Estimated time                                    |
|--|--|---|
| <p>Experiment</p> <p>Caching area: 152 km<sup>2</sup> scale 1:76 (approx. 2D 1:1000)</p> <p>Tiles: 512x512 pixels (finally 104.000 tiles (3/4 tilov ))</p> | <p>ArcGIS Server</p> <p>2x E5450 3GHz (8 threads)</p> <p>32GB memory</p>   | <p>Caching time: 77 minut.</p>                    |
| <p>Generalization for the world World mainland: 148.429.000 km<sup>2</sup>, million times larger area than in experiment</p>                               | <p>ArcGIS Server</p> <p>2x E5450 3GHz (8 threads)</p> <p>32GB memory</p>   | <p>Estimated 77 millions minutes or 146 years</p> |
| <p>Google</p>  | <ul style="list-style-type: none"> <li>• In 2002; upwards of 15,000 servers</li> <li>• A 2005 estimate by Paul Strassmann has 200,000 servers claimed this number to be upwards of 450,000 in 2006 900.000 (2011)</li> </ul> |   |

### Use cases of digital libraries

#### Definition

Digital library is a collection of digital content from libraries, archives, museums and other cultural institutions. It contains internal collections in e.g. museum or in certain branch e.g. movable heritage and resides at national level, european level and world level.

#### Geographical information

Support for visualization, processing geographic information about digital cultural objects. It includes also geocoded historical maps.

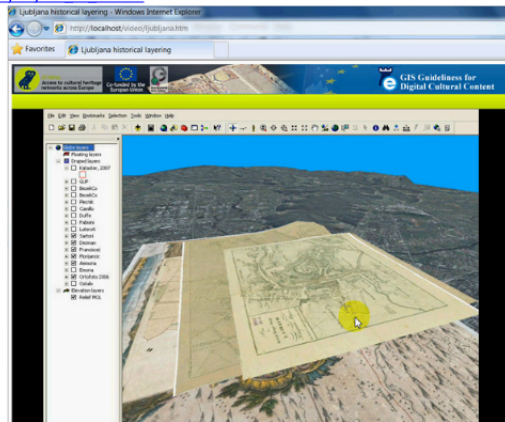
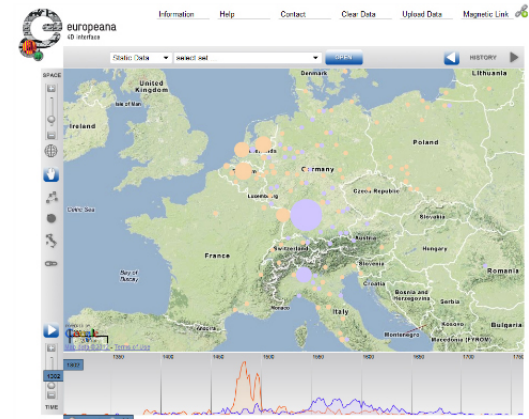
#### Examples

- Europeana globe: <http://www.europeanaglobe.eu/>
- Europeana: <http://www.europeana.eu/portal/>
- arXiv: <http://www.arXiv.org>
- American Memory: <http://memory.loc.gov/ammem/index.html>
- dLib: <http://www.dlib.si/>
- Europeana 4D: <http://wp1187670.wp212.webpack.hosteurope.de/e4d/>
- Judaica: <http://www.judaica-europeana.eu/map>
- 3d historical maps: [http://indicate.situla.org/indicate/Ljubljana\\_M\\_1.wmv](http://indicate.situla.org/indicate/Ljubljana_M_1.wmv)

#### e-Infrastructure

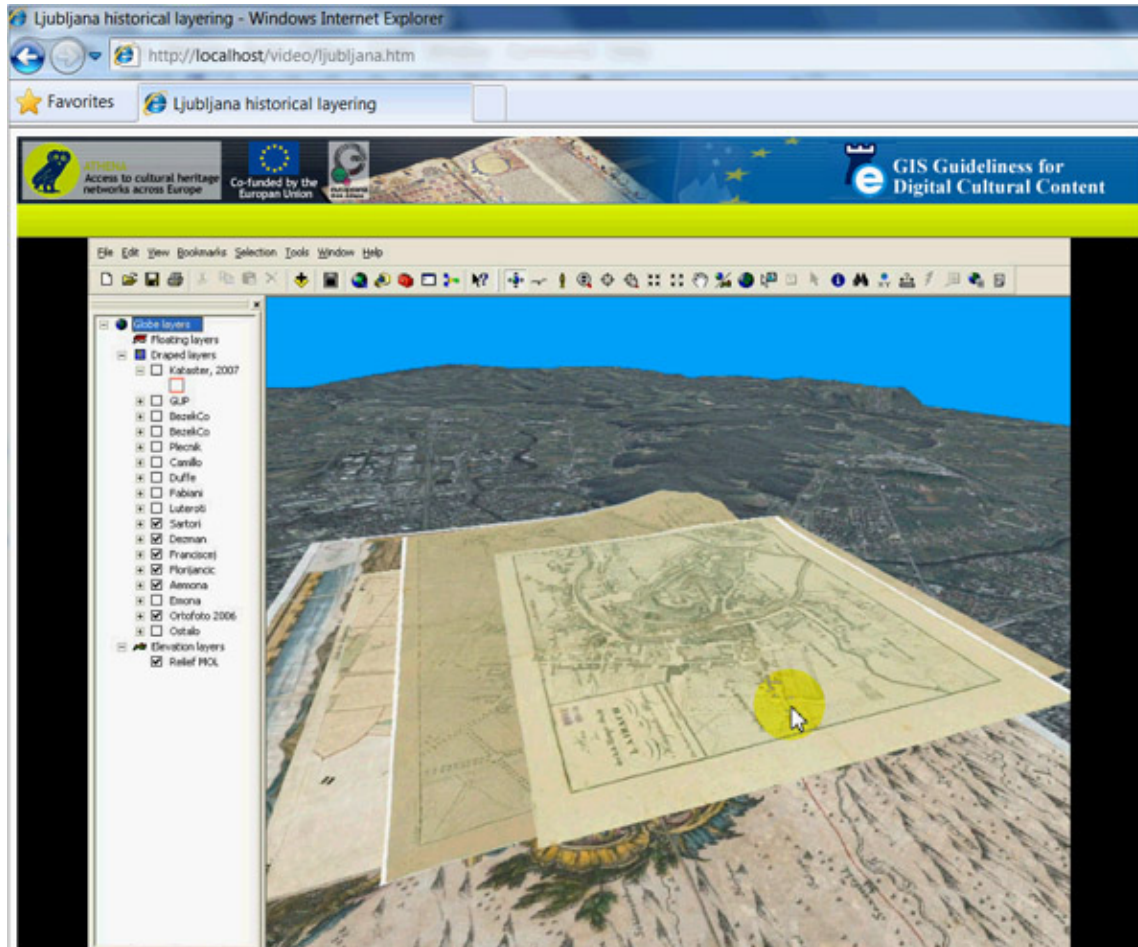
Appropriate for grid computing:

- Digital object of content providers
- GIS, ICT
- Business models ("from the street", established e.g. ECC, NREN, ...)



The screenshot shows the Europeana website interface. At the top, there is a navigation bar with 'Europeana' and 'Juxta Europa'. Below the navigation bar is a search bar and a list of search results. The search results are displayed in a grid format, with each result showing a thumbnail image and a brief description. The search results are filtered by 'Mapsearch' and 'Geograph search'. The search results include a list of countries and regions, and a search result for 'Mapsearch' with a thumbnail image and a brief description.





## **e-Infrastructure**

Appropriate for grid computing:

- Digital object of content providers
- GIS, ICT
- Business models (“from the street”, established e.g. ECC, NREN, ...)

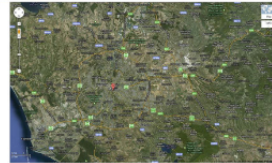
# Geoparsing

## Geoparsing ?

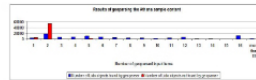


Geoparsing is the process of assigning geographic coordinates to textual words and phrases (e.g. "The author has been born in Rome"). Geoparsing is capable of handling ambiguous references in unstructured content. Geoparsed features can then be mapped and entered into a geographic information system. A geoparser is a piece of software or a (web) service that helps in this process.

## How it works ? /3



## Testing with Athena data /2



Results are for the sample of the Athena content (3,67%). The first input for the geoparser is whole LIDO object and the second input all *places* tags included in the LIDO object. The geoparser found at least one coordinate in 62,37% of LIDO objects and did not found any coordinates in 39,63% of LIDO objects. The exactness of the coordinates found in Geonames gazetteer is not the subject of the analysis.

## State of the art



## Purpose of testing

- Could we find out geographical coordinates from the textual metadata of the certain digital content ?
- What strategies and geoparsing services could we use for geoparsing ?
- What percentage of the content could be geotagged in this way, at the best?
- For what purposes / services could we use the geotagged geographical coordinates (spatial accuracy) ?
- To plan the real production of geoparsing.

## How it works ? /1



## Methods of testing



- The Europeana Geoparser 1.0 data is used
- Input for testing: 4.062.019 LIDO objects in XML format
- Analysis the geoparsing results
- Verification the results on the map

## Spatial accuracy



## How it works ? /2



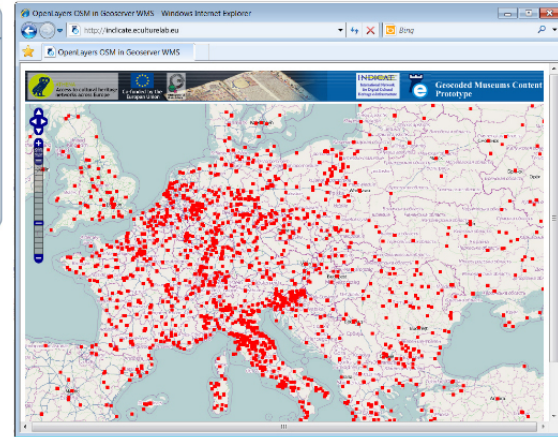
## Testing with Athena data /1



The majority of LIDO objects has at least one *places* tag: 79,53% and 14,83% of the objects do not have *places* tag. The analysis did not include the analysis of the *places* tag itself. If the places exist or their syntax was not the subject of the analysis.

## Conclusions

- Use geoparsing for upper level of LOD (Level Of Detail)
- If there are small towns or villages they were seldom found, inclusion of national register of geographic names is strongly suggested
- Use geoparsing for validation when existing coordinates are correct
- For correct locations use exact coordinates of museum or other cultural memory institution instead of geoparsing
- Use geotagging instead of geoparsing where possible
- Assigning the geographic coordinates as part of documentation process where possible



## e-infrastructure

- Appropriate for grid computing
- Natural language processing (NLP)
- Use of local Gazetteers and other data sources



## Geoparsing ?



**Geoparsing** is the process of assigning geographic coordinates to textual words and phrases (e.g. “The author has been born in Rome”). Geoparsing is capable of handling ambiguous references in unstructured content. Geoparsed features can then be mapped and entered into a geographic information system. A **geoparser** is a piece of software or a (web) service that helps in this process.

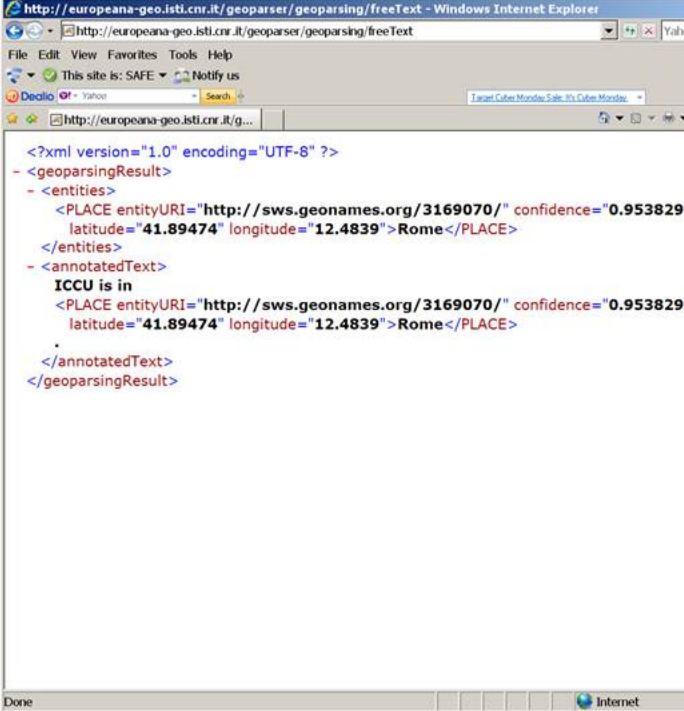
# How it works ? /1

The screenshot shows the Europeana Geoparsing Service v1.0 beta web interface. The browser window title is "Europeana Geoparsing Service - v1.0 beta - Mozilla Firefox". The address bar shows the URL "europeana-geo.isti.cnr.it/geoparser/geoparsing". The page features the Europeana logo and the text "europeana connect". Below the logo, there are several circular images. The main content area is divided into three sections:

- Europeana Geoparsing Service - v1.0 beta**  
Unstructured text and semi-structured text (metadata records) may contain mentions to places and historical periods that are not directly usable by software applications. Geoparsing consists in automatically extracting structured information about places and historical periods from these textual resources. The Geoparser is a web service where users can provide textual sentences or metadata records, and it will reply with an XML document containing the geoparsing results.
- Geoparse free text**  
Example: History of Architecture in Europe: the cases of Lisbon, Madrid and Paris of the 19th century.  
Display result in: XML
- Geoparse ESE xml metadata record**  
Note: only descriptive metadata elements are geoparsed (ie title, description, coverage, etc.).  
Example A Example B Example C Example D  

```
<record id="http://purl.org/dc/terms/1.1/"
xmlns:dc="http://purl.org/dc/terms/"
xmlns:europeana="http://www.europeana.eu" >
<dc:title>Storia della architettura e dell'arte in Europa
<dc:coverage>2008-01-11</dc:coverage>
```

## How it works ? /2



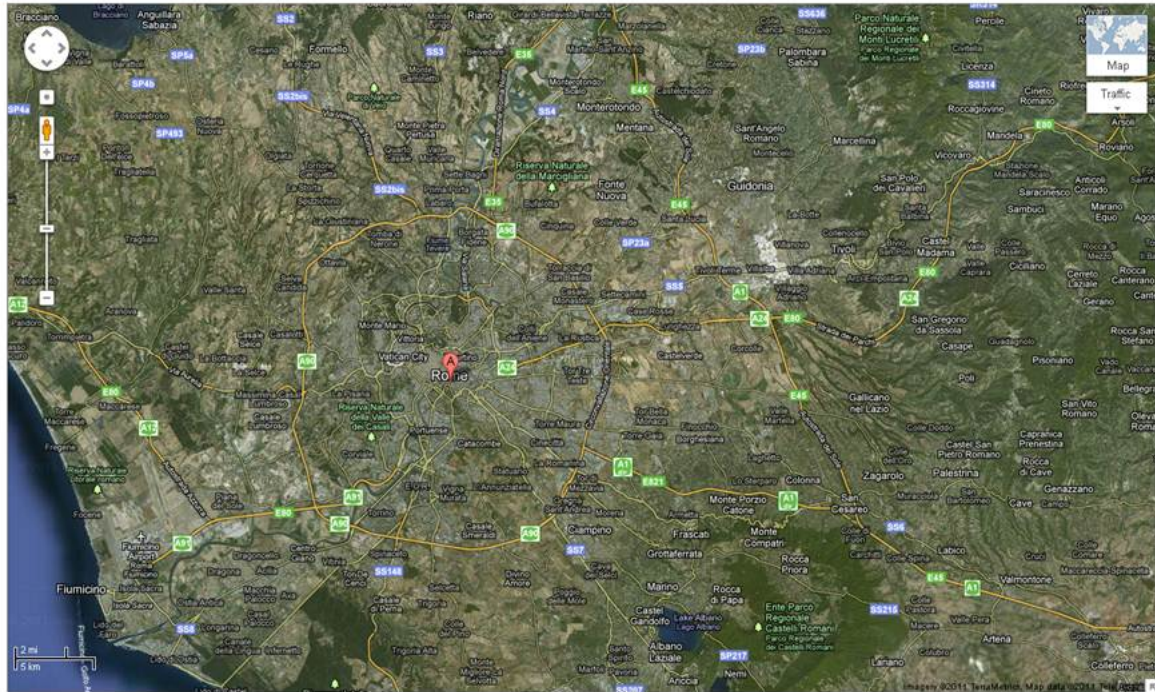
The screenshot shows a Windows Internet Explorer browser window with the address bar displaying `http://europeana-geo.isti.cnr.it/geoparser/geoparsing/freeText`. The page content is XML output from a geoparsing service. The XML structure is as follows:

```
<?xml version="1.0" encoding="UTF-8" ?>
- <geoparsingResult>
- <entities>
  <PLACE entityURI="http://sws.geonames.org/3169070/" confidence="0.953829"
    latitude="41.89474" longitude="12.4839">Rome</PLACE>
</entities>
- <annotatedText>
  ICCU is in
  <PLACE entityURI="http://sws.geonames.org/3169070/" confidence="0.953829"
    latitude="41.89474" longitude="12.4839">Rome</PLACE>
  .
</annotatedText>
</geoparsingResult>
```

The status bar at the bottom of the browser window shows "Done" and "Internet".



## How it works ? /3



## Purpose of testing

- Could we find out geographical coordinates from the textual metadata of the certain digital content ?
- What strategies and geoparsing services could we use for geoparsing ?
- What percentage of the content could be geocoded in this way, at the best?
- For what purpose / services could we use the geoparsed geographical coordinates (spatial accuracy) ?
- To plan the real production of geoparsing

## Methods of testing

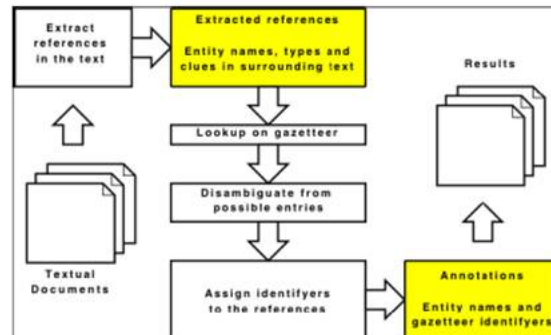
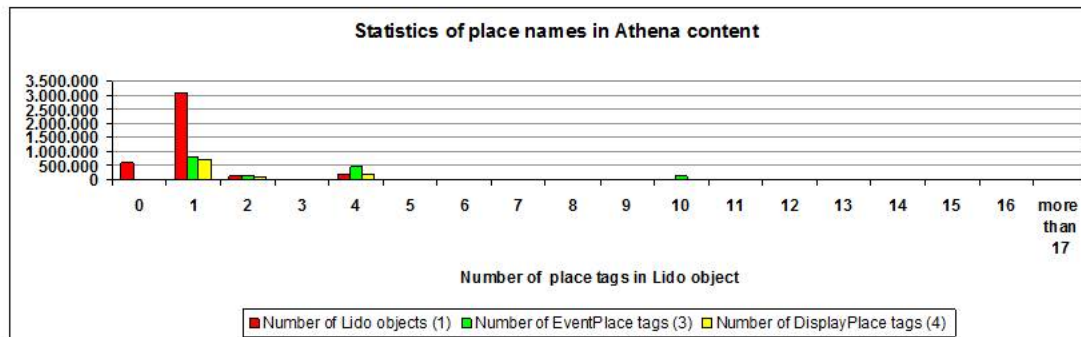


Figure 1. Typical approach for geo-parsing text

- The Europeana Geoparser v 1.0 Beta is used
- Input for testing: 4.082.619 LIDO objects in XML format
- Analysis the geoparsing results
- Verification the results on the map

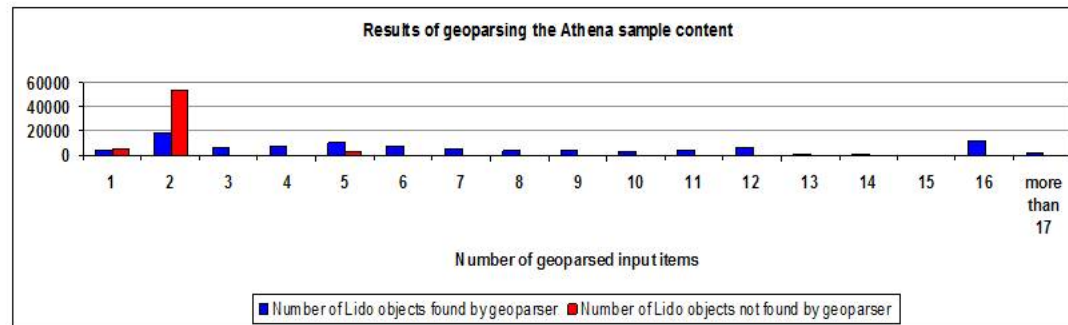


## Testing with Athena data /1



The majority of Lido objects has at least one »Place« tag: 75,53% and 14,83% of the objects do not have »Place« tag. The analysis did not include the analysis of the »Place« tag itself. If the places exist or their syntax was not the subject of the analysis.

## Testing with Athena data /2



Results are for the sample of the Athena content (3,87%) The first input for geoparser is whole Lido object and the second input all »place« tags included in the Lido object. The geoparser found at least one coordinates in 60,37% of Lido objects and did not found any coordinates in 39,63% of Lido objects. The exactness of the coordinates found in Geonames gazetter is not the subject of the analysis.

## Spatial accuracy



**The same geographic name for different places**

e.g. "Paris" addresses 93 places



**The different name for the same place**

e.g. Istanbul (Nova Roma, Constantinople, Tsargrad, ...)

OpenLayers OSM in Geoserver WMS - Windows Internet Explorer

http://indicate.eculturelab.eu

OpenLayers OSM in Geoserver WMS

**ATHENA**  
Access to cultural heritage networks across Europe

Co-funded by the European Union

**INDICATE**  
International Network for Digital Cultural Heritage e-Infrastructure

**Geocoded Museums Content Prototype**

The map displays a dense distribution of red square markers across Europe, representing geocoded museum locations. The markers are most concentrated in Western and Central Europe, including the British Isles, France, Germany, and Poland. Major cities like London, Paris, Berlin, and Warsaw are clearly visible. The map also shows national and regional boundaries, with labels in various languages such as English, Polish, and Russian. A navigation toolbar is located on the left side of the map, and a scale bar is visible at the bottom left.



## **e-infrastructure**

- Appropriate for grid computing
- Natural language processing (NLP)
- Use of local Gazetteers and other data sources

## Conclusions

- Use geoparsing for upper level of LOD (Level Of Detail)
- If there are small town or villages they were seldom found, inclusion of national register of geographic names is strongly suggested
- Use geoparsing for validation when existings coordinates are correct
- For current locations use exact coordinates of museum or other cultural memory institution instead of geoparsing
- Use geotagging instead of geoparsing where possible
- Assigning the geographic coordinates as part of documentation process where possible

## Knowledge Café

### 1. table: Archaeological / architectural heritage and GIS

Facilitator: Matteo Lorenzini

Expected attendee: archaeological and architectural institutions (10+)

Questions:

- identification of additional use cases (regularly operating, research, planned) (additional form) navigation...
- benefits and weakness of open source
- effective browsing of 3D cities
- geographical coordinate systems
- the needs for grid and cloud computing (restoration, 3D rendering, caching, ...)

### 2. table: Libraries and GIS

Facilitator: Annette Kolly

Expected attendee: libraries and other cultural institutions (10)

Questions:

- identification of additional use cases (regularly operating, research, planned) (additional form) navigation...
- benefits expected from GIS in libraries
- geocoding or geotagging geographical coordinates
- geocoding of the historical maps
- GIS in Europeana
- the need for grid and cloud computing

### 3. table: Museums and cloud computing

Facilitators: Jernej Porenta , Luka Hribar

Expected attendee: museums and other cultural institutions (10)

Questions:

- identification of additional use cases (regularly operating, research, planned) (additional form) navigation...
- comparison the costs (in house server ITC vis-a-vis cloud computing)
- persistent identifiers
- e-infrastructure expected from NRNs