

# Metadata Topic Harmonization and Semantic Search for Linked-Data-Driven Geoportals -- A Case Study Using ArcGIS Online

Yingjie Hu<sup>1</sup>, Krzysztof Janowicz<sup>1</sup>, Sathya Prasad<sup>2</sup>, and Song Gao<sup>1</sup>

<sup>1</sup> STKO Lab, Department of Geography, U.C. Santa Barbara

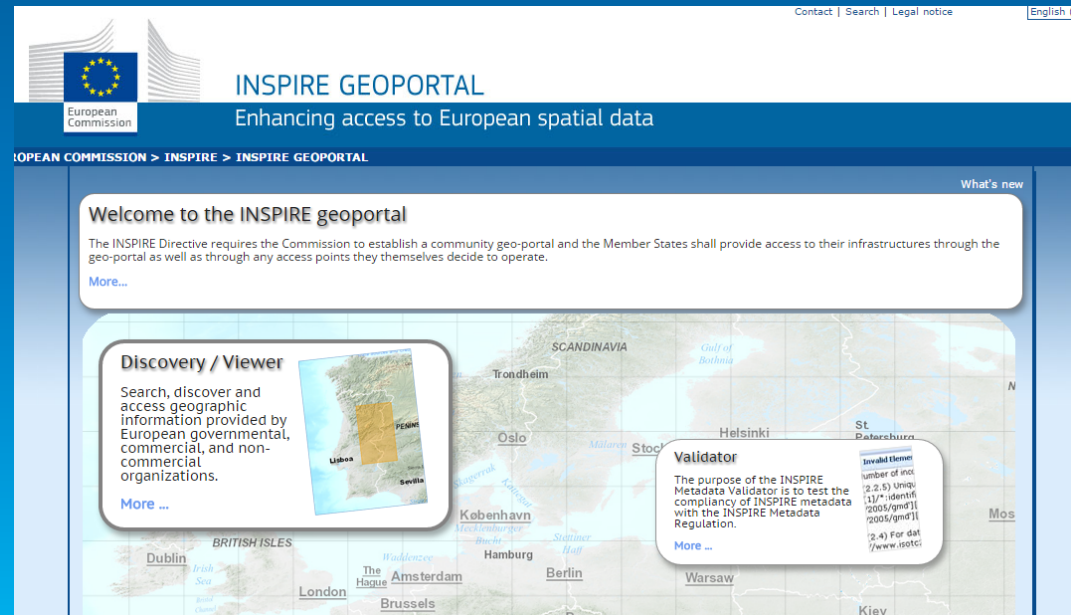
<sup>2</sup> Applications Prototype Lab, Esri Inc.

# Outline

- **Introduction and motivation**
- **Metadata Topic Harmonization**
- **Semantic Search on Linked Data**
- **Experiments and Evaluation**
- **Conclusions and Future Work**

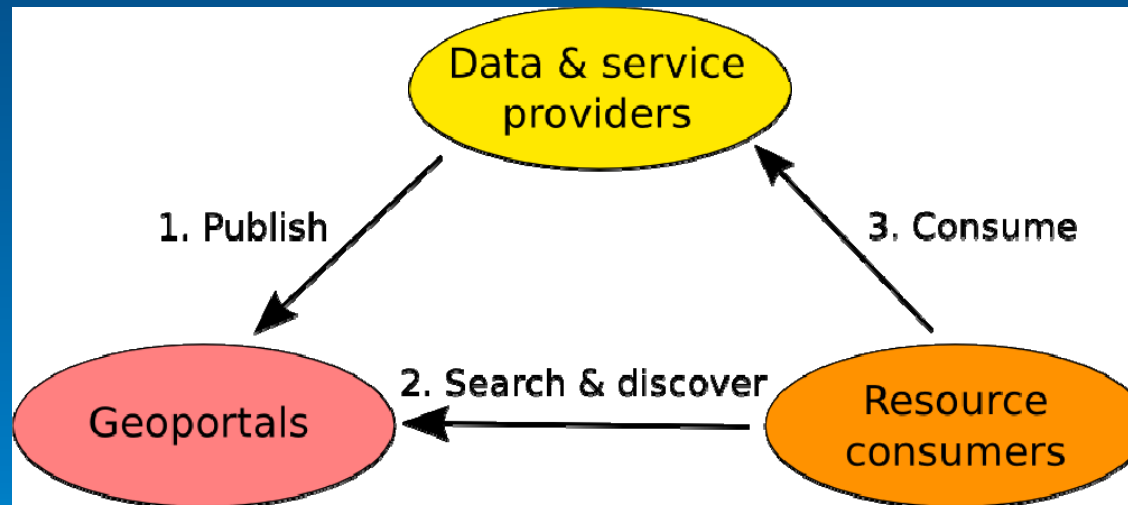
# Introduction

- Geoportals are Web gateways that provide integrated access to geospatial resources
- Geoportals are key components of Spatial Data Infrastructures (SDI)
- Existing geoportals:
  - Data.gov
  - INSPIRE
  - California geoportal
  - ...



# Introduction

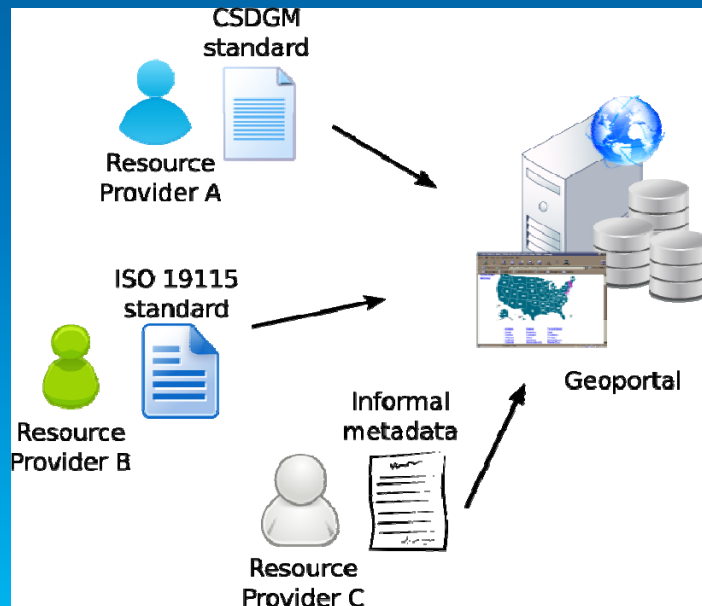
- A typical *publish-find-bind* pattern has been used by many geoportals



- Two main factors that can influence resource discovery:
  - Quality of metadata
  - Search functionality

# Introduction

- Quality of metadata
  - Multiple standards have been established to ensure the metadata quality, e.g., FGDC's CSDGM and ISO 19115
  - However, data contributed to the same geoportal may be in different standards



The screenshot shows the FGDC (Federal Geographic Data Committee) website. The page title is 'Geospatial Metadata Standards'. The content discusses the transition from CSDGM to ISO metadata standards. A red box highlights the text: 'developed standards. Since ISO 19115 and the associated standards are endorsed by the FGDC, federal agencies are encouraged to transition to ISO metadata as their agencies are able to do so. While the selection of'.

you are here: home → metadata → geospatial metadata standards

## Geospatial Metadata Standards

Most NSDI stakeholders have long utilized the Content Standard for Digital Geospatial Metadata (CSDGM), which will continue to have a legacy for many years. International geospatial metadata standards are emerging in the community. FGDC policy states that non-Federally authored standards that are endorsed by the FGDC have the same status as FGDC developed standards. Since ISO 19115 and the associated standards are endorsed by the FGDC, federal agencies are encouraged to transition to ISO metadata as their agencies are able to do so. While the selection of appropriate standards is dependent on the nature of your metadata collection and publication process, ISO metadata should be considered an option now. It's recognized that the transition to ISO metadata will be occurring over the next few years.

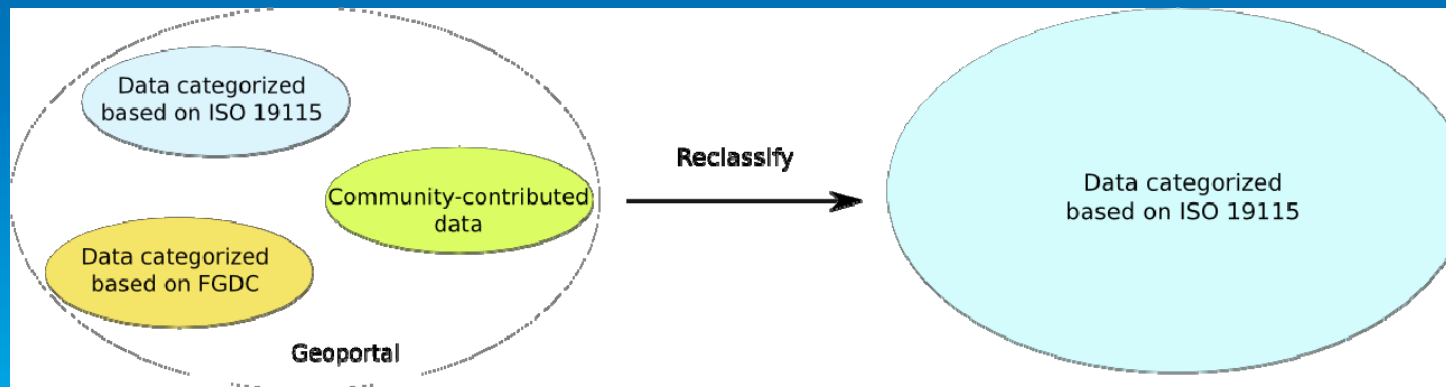
- FGDC Endorsed ISO Metadata Standards
  - Which ISO Metadata Standards Do Organizations Utilize?
  - Should Our Agency be using ISO Metadata?
  - ISO Metadata Implementation is Happening
- The Content Standard for Digital Geospatial Metadata (CSDGM)
  - CSDGM Resources
  - CSDGM Profiles and Extensions

FGDC Endorsed ISO Metadata Standards

The Office of Management and Budget's Circular A-119 establishes policies on Federal use and development of voluntary consensus standards. The FGDC

# Introduction

- **Quality of metadata**
  - **How to harmonize metadata in different standards?**
    - Some elements can be automatically mapped using, e.g., NOAA's metadata transformation tool
    - Some others have to be transformed manually, e.g., the topics



# Introduction

- **Search functionality**
  - **Traditional keyword-based search**
    - Based on keyword matching
    - E.g., A search of “natural disaster” can only return maps which contain words “natural” or “disaster”
  - **Semantic search**
    - Find maps based on the meaning of input query
    - E.g., A search of “natural disaster” can return different disasters, such as wildfire, hurricane, earthquake...

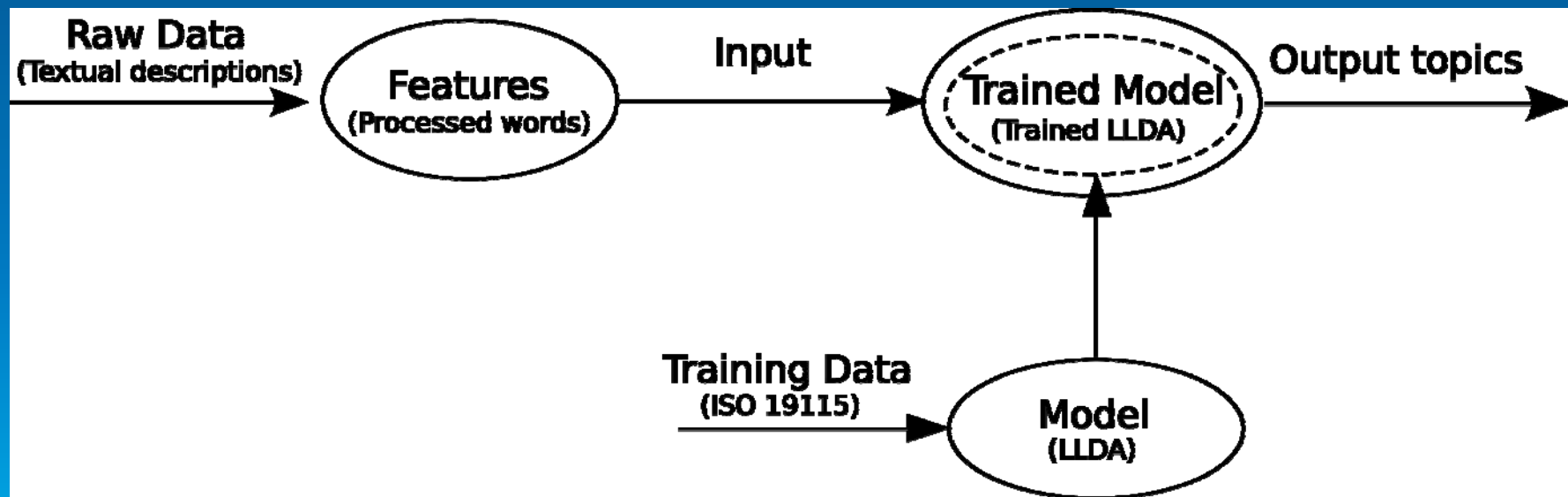
# Introduction

- **Search functionality**
  - **The emergence of Linked-Data-driven Geoportals**
    - Accommodate heterogeneous data using RDF data model
    - Graph-based data storage and browsing
    - Help discover the links hidden in data
  - **Semantic search functionality for RDF data is not available for Linked-Data-driven geoportals**



# Metadata Topic Harmonization

- A machine learning approach
  - A multi-label classification problem
  - One metadata can be associated with multiple topics
  - Based on titles and descriptions of each metadata entry



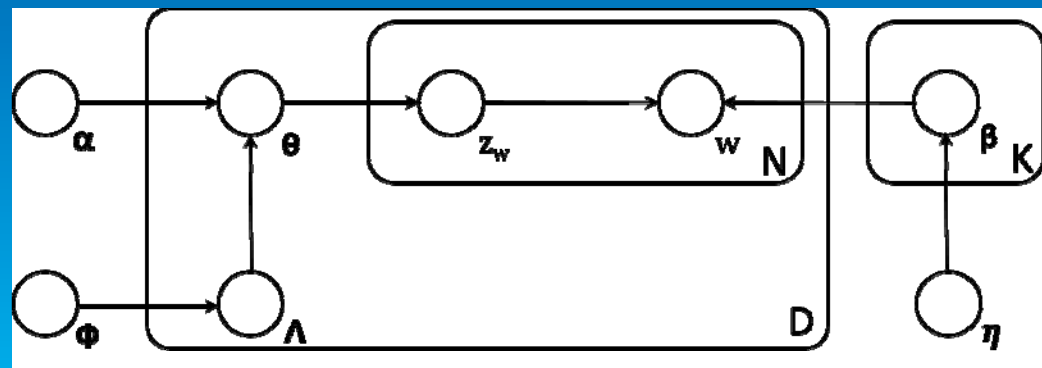
# Metadata Topic Harmonization

- **Naïve Bayesian model has been commonly used in text-based classification, but it has limitations**
  - **Assuming only one topic for each record; not suitable for multi-label classification**
  - **Biased estimation for prior probability**
  - **May result in overfitting for long texts**

$$P(t_i|d) \propto \prod_{j=1}^N P(w_j|t_i) \times P(t_i)$$

# Metadata Topic Harmonization

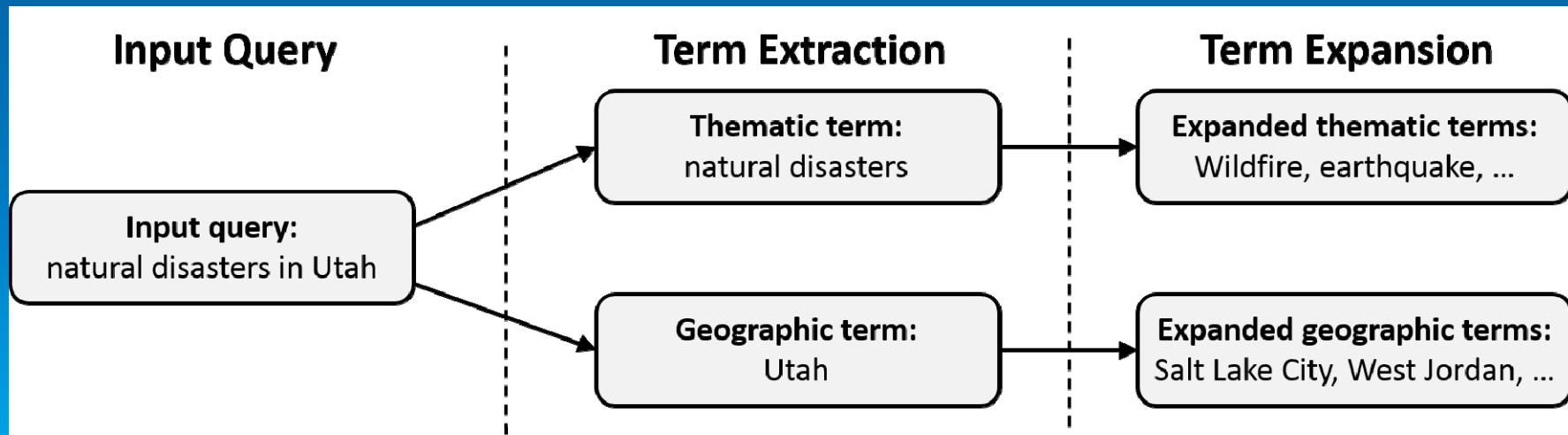
- **LLDA (Labeled Latent Dirichlet Allocation)**
  - An extension of LDA by adding a component of supervised learning
- **Advantages of LLDA**
  - Considers each document as a mix of multiple topics
  - Robust estimation for prior probabilities of topics
  - Avoid overfitting for long descriptions



# Semantic Search for Linked Data

- Query expansion

- Extracting concepts and entities from the input query
- Expanding them using related concepts and entities
  - Thematic concepts: Latent Semantic Analysis (LSA) and Wordnet
  - Geographic entities: Gazetteer service (Geonames)



# Semantic Search for Linked Data

- **Constructing Matching Features**

- Is this matching happens in title or in description?
- Is this matching a thematic matching or geographic matching?
- Is this an exact matching or a similar matching?
- Resulted in 8 matching features (2 x 2 x 2)

---

Title Thematic Exact match (TTE)

Title Geographic Exact match (TGE)

Snippet Thematic Exact match (STE)

Snippet Geographic Exact match (SGE)

---

Title Thematic Similar match (TTS)

Title Geographic Similar match (TGS)

Snippet Thematic Similar match (STS)

Snippet Geographic Similar match (SGS)

---

# Semantic Search for Linked Data

- **Constructing Matching Features**

- **An additional feature: Thematic-Geo Interaction (TGI)**

$$TGI = (TTE + TTS + STE + STS) \times (TGE + TGS + SGE + SGS)$$

- **Rationale for introducing this interaction feature:**
  - Thematic or geo matching alone cannot determine the relevance
  - E.g., Searching “Crime in California”
  - “Crime in Florida” or “Waterbody in California” may not be what users want
  - “Robberies in Los Angeles” may be relevant

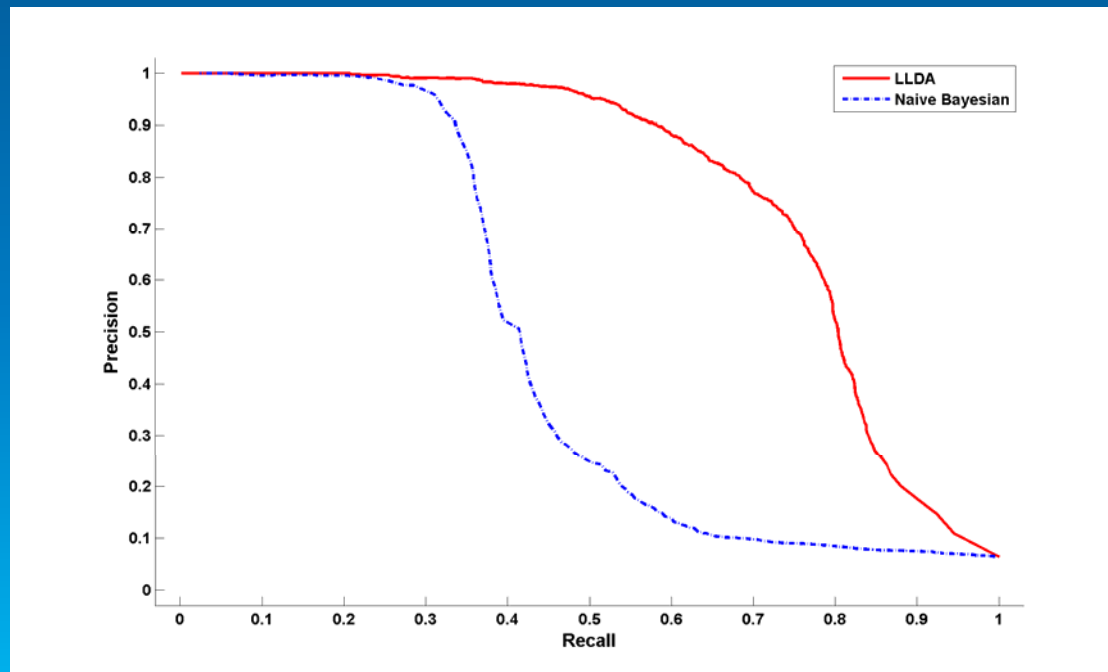
$$R(q, m) = \lambda_1 TTE + \lambda_2 TTS + \lambda_3 TGE + \lambda_4 TGS + \lambda_5 STE + \lambda_6 STS + \lambda_7 SGE + \lambda_8 SGS + \lambda_9 TGI$$

# Experiments and Evaluation

- **Experimental data:**
  - 26, 917 metadata records from Data.gov in ISO 19115
  - 10, 201 metadata records from ArcGIS Online
- **Experiment procedure:**
  - Use metadata from Data.gov to train the LLDA model
  - Apply the trained LLDA to the unstandardized ArcGIS Online data
  - Train the ranking model using a human participant experiment

# Experiments and Evaluation

- Metadata topic harmonization
- Comparing the performances of LLDA and naïve Bayesian
  - Data from Data.gov (with known ground truth)
  - Ten-fold cross validation
  - Precision and recall curves

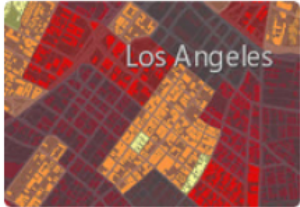




# Experiments and Evaluation

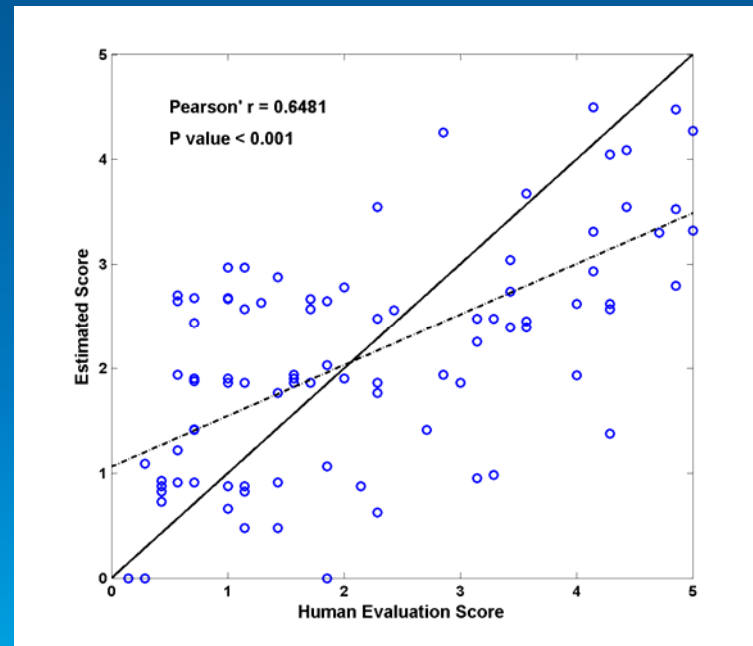
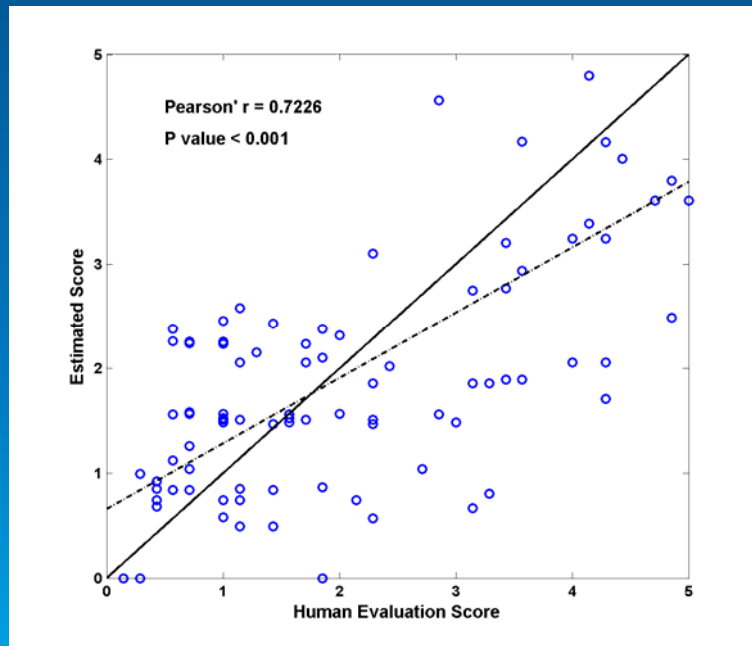
- Semantic search and ranking
- Human participant experiment
  - 7 human participants
  - Each person evaluate 10 queries and each query has 10 candidate maps
  - For each query and candidate, provide a score [0, 5]

## Query 3: "california population density"

Map	Link of the Map
<p>3.1</p>  <p><b>Los Angeles Population Density</b> This map emphasizes areas with the highest population density (more than 50,000 persons per square kilometer).</p>	<p><a href="http://www.arcgis.com/home/webmap/viewer.html?webmap=4971065a7a734e31a7079ace59a19f27">http://www.arcgis.com/home/webmap/viewer.html?webmap=4971065a7a734e31a7079ace59a19f27</a></p>

# Experiments and Evaluation

- Semantic search and ranking
- Ten-fold cross validation using Pearson's  $r$ 
  - With v.s. without the interaction variable



# Semantic Search for Linked Data

- Embedding the semantic search to a geoportal
  - A SPARQL query to implement the regression model

```
SELECT ?item (COUNT(?titleThematicExact) AS ?TTE
(COUNT(?titleThematicSimilar) AS ?TTS)
(COUNT(?titleGeoExact) as ?TGE)
(COUNT(?titleGeoSimilar) as ?TGS)
(COUNT(?snipThematicExact) as ?STE)
(COUNT(?snipThematicSimilar) as ?STS)
(COUNT(?snipGeoExact) as ?SGE)
(COUNT(?snipGeoSimilar) as ?SGS)
(((?TTE+?TTS+?STE+?STS)*(?TGE+?TGS+?SGE+?SGS)) as ?TGI)
((  $\lambda_1$ *?TTE +  $\lambda_2$ *?TTS +  $\lambda_3$ *?TGE +  $\lambda_4$ *?TGS +  $\lambda_5$ *?STE +  $\lambda_6$ *?STS +
 $\lambda_7$ *?SGE +  $\lambda_8$ *?SGS +  $\lambda_9$ *?TGI) as ?ranking)
WHERE {
  OPTIONAL {
    ?item :hasTitleThematicTerm ?titleThematicExact .
    FILTER ( ?titleThematicKey = :exactThematicTerm ) }
  OPTIONAL {
    ?item :hasTitleThematicTerm ?titleThematicSimilar .
    FILTER ( ?titleThematicSimilar = :expandedThematicTerm ) }
  OPTIONAL {
    ?item :hasTitleGeoTerm ?titleGeoExact .
    FILTER ( ?titleGeoExact = :exactGeoTerm ) }
  OPTIONAL {
    ?item :hasTitleGeoTerm ?titleGeoSimilar .
    FILTER ( ?titleGeoSimilar = :expandedGeoTerm ) }
```


# Interactive prototype

- <http://stko-exp.geog.ucsb.edu/linkedportal>

Linked-Data-driven Geoportal for ArcGIS Online

## Semantic Search

Semantic search based on a sample of map data from ArcGIS Online.  
For efficiency, maximum 1,000 maps are returned for each query.




natural disaster | [Search](#)

Identified Thematic Concept: [natural disaster](#)


ISO 19115 topics:

- [all Topics \(354\)](#)
- [transportation \(0\)](#)
- [elevation \(2\)](#)
- [oceans \(3\)](#)
- [geoscientific \(26\)](#)
- [imagery \(10\)](#)
- [environment \(34\)](#)
- [inlandWaters \(83\)](#)
- [boundaries \(3\)](#)
- [location \(16\)](#)
- [utilities \(0\)](#)
- [health \(167\)](#)
- [structure \(1\)](#)
- [planning \(14\)](#)
- [society \(10\)](#)
- [biota \(2\)](#)
- [economy \(1\)](#)
- [farming \(0\)](#)
- [climatology \(26\)](#)
- [intelligence \(0\)](#)




**MIDLAND TORNADO - 2010**  
Destruction path of the Midland Tornado in the summer of 2010

[Open Map](#)




**Hurricanes striking the Country of Belize**  
for the past few years, Belize has been hit by some hurricanes. These hurricanes have caused vast destruction and to some extent caused deaths.

[Open Map](#)




**DRAFT Flood Study for Cedar Creek & Tribs in Washington Coun...**  
DRAFT Flood Study for Cedar Creek & Tribs in Washington County

[Open Map](#)




**DRAFT Flood Study for Cedar Creek & Tribs in Washington Coun...**  
DRAFT Flood Study for Cedar Creek & Tribs in Washington County

[Open Map](#)



**DRAFT Flood Study for Cedar Creek & Tribs in Washington Coun...**  
DRAFT Flood Study for Cedar Creek & Tribs in Washington County

[Open Map](#)



**DRAFT Flood Study for Cedar Creek & Tribs in Washington Coun...**  
DRAFT Flood Study for Cedar Creek & Tribs in Washington County

[Open Map](#)

## Conclusions and Future Work

- **Geoportals provide integrated access to geospatial resources**
- **The quality of metadata and the capability of the search function are two major factors affecting resource discovery**
- **We present a LLDA-based approach for harmonizing metadata topics, as well as enabled semantic search for RDF data**
- **Limitations and future work**
  - **Small scale human participants test need to be expanded**
  - **Increase the response efficiency of semantic search**

# Thank you!

Yingjie Hu  
yingjiehu@umail.ucsb.edu  
<http://www.geog.ucsb.edu/~hu/>