# Linked Data

1 --One goal

2 – Two types of questions

3 – RDF triples

4 – Four principles

*Learn by Understanding*

5 – Five star LOD

*Learn by Analyzing*

## Hands-on

- Obtain Data Using SPARQL

- Entity extraction

*Learn by Creating*

---

# Digging into unstructured data

How can we improve discovery from unstructured data (i.e., go deeper into the contents) through effective processes?

Instead of being based on the *full-text searching*, *semantic-based analysis*, *extraction*, *mining*, and *annotations* are used to improve the information discovery from unstructured data.

- Automatic semantic analysis, to generate additional access points;

- Text mining;

- Fact mining (taxonomy and ontology-supported)

- User contributed annotations of text and images

Marcia Zeng, 2018 DIS                                         2

1

## An example

Many archival materials are digitized. But they are still unstructured, just like this interesting website:

http://international.loc.gov/intldl/brhtml/br-1/br-1-5-2.html#track1

Now the only name that is obvious is Dom Pedro II.

**Dom Pedro II and America**
**Dom Pedro II e os Estados Unidos**

United States and Brazil Home >> Historical Foundations >> The Imperial period >> Dom Pedro II and America

Historical Foundations

- Timeline
- Navigation and Discovery
- The Colonial Period
- The Movement for Independence
- The Imperial Period
  - Dom Pedro II
  - Dom Pedro II and America
  - Princess Isabel
  - Coffee and Rubber
  - The Triple Alliance War

An admirer of Abraham Lincoln (1809-1865), Dom Pedro II visited the United States in1876 to join President Ulysses S. Grant (1822-1885) in opening the Philadelphia Centennial Exposition, the largest world's fair up to that time. While at the exposition, the emperor met Alexander Graham Bell (1847-1922), who was demonstrating his new invention--the telephone. Dom Pedro II was the first person to buy stock in Bell's company, the Bell Telephone Company. One of the first telephones in a private residence was installed in his palace in Petrópolis, his summer retreat forty miles from Rio de Janeiro.

Dom Pedro II had a deep appreciation for the potential contributions of science and technology to society. He was interested in expanding his country's primary agricultural commodities and in making Brazil a major cotton producer.

The empire of Brazil at the Universal exhibition of 1876 in Philadelphia, 1876. Library of Congress. General Collections.

Dom Pedro II spent two weeks at the Philadelphia Centennial Exposition, where he contacted

Marcia Zeng, 2018 DIS

3

---

*Think: In this document, how many interesting things people would like to know more about…*

An admirer of Abraham Lincoln (1809-1865), Dom Pedro II visited the United States in1876 to join President Ulysses S. Grant (1822-1885) in opening the Philadelphia Centennial Exposition, the largest world's fair up to that time. While at the exposition, the emperor met Alexander Graham Bell (1847-1922), who was demonstrating his new invention--the telephone. Dom Pedro II was the first person to buy stock in Bell's company, the Bell Telephone Company. One of the first telephones in a private residence was installed in his palace in Petrópolis, his summer retreat forty miles from Rio de Janeiro.

We can [manually]make hyperlinks under strings … (machine-processable)

= linking strings

Marcia Zeng, 2018 DIS

4

## Slide 1

From linking <u>strings</u> →to→ linking <u>things</u>

An adm... **agent** Abraham Lincoln (1809-18..) Dom Pedro II visi.. **place** United States in1876 .... President Ulysses S. Grant (1822-1885) **agent** ..ening **event** Philadelphia Centennial Exposition, the largest world's fair up to that time. While at the exposition, the emperor met Alexander Graham Bell **agent** (1847-1922), who was demonstrating his new inve... the telephone. Dom Pedro II **agent** e first **objects** perso.. uy stock in Bell's company, the Bell Telephone Company, **agent** of the first telephones in a private residence was installed in his palace **place** Petrópolis, his summer retreat forty miles from Rio **place** de Janeiro.

**time**

Or, we should think the <u>types of entities</u> involved and make them <u>machine-understandable</u>

=linking <u>things</u>

Marcia Zeng, 2018 DIS                                    5

## Slide 2

THOMSON REUTERS | OPEN CALAIS

# Semantic Analysis Tool #1
## Open Calais

"Open Calais processes the text you submit and returns: Entities, Topic codes, Events, Relations and SocialTags. Where applicable, Open Calais will link to permanent entity identifiers, served by our new beta initiative:"

http://www.opencalais.com/



3

Demo: Using Open Calais (free version) tool to run through the text of that paragraph showing previously.

1. Go to: http://www.opencalais.com/opencalais-demo/



---

2. Copy the text into the box. "Submit"

I used the same text used in my example. If you use another text, do not use private or unclassified material.



http://www.opencalais.com/opencalais-demo/

# Behind the theme

Let's check out what this looks like at the back-end



What are the structured data behind these?
(Coded with RDF/XML syntax. [You do not need to understand the codes at the moment.])

DOCUMENT VIEW    Upload Again    View RDF

An admirer of Abraham Lincoln (1809-1865), Dom Pedro II visited

Clicking on "View RDF" brings us to the machine-processable datasets.

data:text/xml,%3C%3Fxml version%3D"1.0" encoding%3D"UTF−8"%3F%3E%0A%3C!–

This XML file does not appear to have any style information associated with it. The document tree is shown below.

These green lines show the automatically extracted entities, topics, relations, etc.

```
<!--
   Use of the Calais Web Service is governed by the Terms of Service located at http://www.opencalais.      results of the service yo
-->
<!--
   Relations: Acquisition, PersonCareer, PersonTravel
   City: Mobile, Philadelphia, Rio de Janeiro, São Paulo
   Company: Bell, Bell Telephone Company
   Country: Brazil, United States
   Facility: Confederate cemetery, port of Rio de Janeiro
   Holiday: 4th of July
   Person: Abraham Lincoln, Alexander Graham Bell, Bill Clinton, Civil War, Jimmy Carter, Pedro II, Rosalyn Carter, Ulysses S. Grant, William Hutchinson Norris
   Position: Colonel, President, emperor, major cotton producer
   ProvinceOrState: Alabama, Arkansas
   Technology: Dom
-->
<rdf:RDF>
 − <rdf:Description c:calaisRequestID="51d83d4b-5c39-6bf3-1599-a67f5a31af30" c:id="http://id.opencalais.com/1B*Pc4vYVs73fl-F-bP3oQ" c:ontology="http://mdaas-virtual-
   onecalais.int.thomsonreuters.com/owlschema/10.2/onecalais.owl.allmetadata.xml" rdf:about="http://d.opencalais.com/dochash-1/4ad50bb4-16a4-350f-9e85-82c08b6d6bfa">
     <rdf:type rdf:resource="http://s.opencalais.com/1/type/sys/DocInfo"/>
   − <c:document>
      <Document><Source></Source><Body>An admirer of Abraham Lincoln (1809-1865), Dom Pedro II visited the United States in1876 to join President Ulysses S. Grant (1822-1885) in
```

Marcia Zeng, 2018 DIS                                                         12

**Slide 1**

THOMSON REUTERS | OPEN CALAIS

– **<rdf:Description rdf:about**="http://d.opencalais.com/er/person/ralg-pa1/3441832
  **<rdf:type rdf:resource**="http://s.opencalais.com/1/type/er/Person"/>
  **<c:docId rdf:resource**="http://d.opencalais.com/dochash-1/4ad50bb4-16a4-350
  **<c:name**>John Alexander Bell**</c:name**>
  **<c:personid**>2157887**</c:personid**>
  **<c:paid**>34418327859**</c:paid**>
  **<c:officerid**>N/A**</c:officerid**>
  **<c:commonname**>John Alexander Bell**</c:commonname**>
  **<c:score**>0.6591502**</c:score**>
  *<!--Alexander Graham Bell-->*
  **<c:subject rdf:resource**="http://d.opencalais.com/pershash-1/89ff7fbd-7150-3225-be2d-a15d6fce1d34"/>
**</rdf:Description>**

> If you search "Bell" on this code page, you can find several description sets, including this one, coded his name, officer–ID, etc.

> Remember on the human-readable page, we saw this data already?

largest world's fair up to that time. While at the exposition, the emperor met Alexander Graham Bell (1847-1922), who was

| | |
|---|---|
| persontype | N/A |
| nationality | N/A |
| confidencelevel | 0.985 |
| firstname | Alexander |
| middlename | Graham |
| lastname | Bell |
| commonname | John Alexander Bell |
| personid | 2157887 |
| paid | 34418327859 |

Pedro II was
Telephone
nce was
forty miles

☑ Civil War 20%
☑ James M. C… 20%
☑ John Alexan… 20%
☑ Pedro II 20%
☑ Rosalyn Car… 20%
…es S. … 20%

Position
☑ Province Or State

interested in expanding his country's primary agricultural

Marcia Zeng, 2018 DIS          13

---

**Slide 2**

THOMSON REUTERS | OPEN CALAIS

– **<rdf:Description rdf:about**="http://d.opencalais.com/er/geo/city/ralg-geo1/06db416f-aa27-f883-8bbc-e038dbfbde9d">
  **<rdf:type rdf:resource**="http://s.opencalais.com/1/type/er/Geo/City"/>
  **<c:docId rdf:resource**="http://d.opencalais.com/dochash-1/4ad50bb4-16a4-350f-9e85-82c08b6d6bfa"/>
  **<c:name**>Philadelphia,Pennsylvania,United States**</c:name**>
  **<c:shortname**>Philadelphia**</c:shortname**>
  **<c:latitude**>39.9522**</c:latitude**>
  **<c:longitude**>-75.1642**</c:longitude**>
  **<c:containedbystate**>Pennsylvania**</c:containedbystate**>
  **<c:containedbycountry**>United States**</c:containedbycountry**>
  *<!--Philadelphia-->*
  **<c:subject rdf:resource**="http://d.opencalais.com/genericHasher-1/4b474cf2-e8df-3193-86a4-8e858a5c5e37"/>
  **<c:rcscode**>G:22H**</c:rcscode**>
**</rdf:Description>**

> Same for places, e.g., "Philadelphia".
> Can you identify its latitude and longitude?

> Again, remember in the human-readable page, we saw this data already?

**CITY**
**Philadelphia,Pennsylvania,United States**

20%   Relevance
      **20%**

| | |
|---|---|
| shortname | Philadelphia |
| latitude | 39.9522 |
| longitude | -75.1642 |
| containedbystate | Pennsylvania |
| containedbycountry | United States |

☑ City
☑ Philadelphia… 20%
☑ Rio de Jane… 20%

Lincoln (1809-1865)
to join President U
the Philadelphia Cer
that time. While at
Graham Bell (1847-
demonstrating his new invention--the teleph
the first person to buy stock in Bell's compa
Company. One of the first telephones in a

Marcia Zeng, 2018 DIS          14

How can we use this data? (Technical processes).

1. If you know programming, you can convert this data into a table.
2. The RDF/XML file can be converted into a different format, and stored.

**Refine** OPEN

3. For individual or multiple sets of data, use OpenRefine http://openrefine.org/ to clean up and reconcile.

E.g., merge synonyms of an entity, such as "Kent State", "KSU", "Kent State University", "Kent State U.", etc.

After merging, the named entity can go through the reconcile process; OpenRefine also enables you to obtain URIs from name authorities and thesauri through API.

*Let me know if you want an example from my research of using OpenRefine.*

| | | | | | |
|---|---|---|---|---|---|
| | 3. | John D. Cimperman | | 0.661 | Cleveland L State Unive |
| | | ☑ ☑ Create new topic | | | |
| | | Search for match | | | |
| | 29. | Wayne Woodrow Hayes | | 0.648 | Wayne W. W |
| | | ☑ ☑ Create new topic | | | |
| | | Search for match | | | |
| | 7. | Charles Adams | edit | 0.646 | Hartford Me University o |
| | | ☑ ☑ Charles Adams Platt (0.684) | | | |
| | | ☑ ☑ Charles Adams Blakely (0.619) | | | |
| | | ☑ ☑ Charles Adams (Colorado) (0.542) | | | |
| | | ☑ ☑ Create new topic | | | |
| | | Search for match | | | |
| | 21. | Dayton Fic | | 0.645 | Guide to the |
| | | ☑ ☑ Create new topic | | | |
| | | Search for match | | | |
| | 12. | Frank S. Box | | 0.629 | Wisconsin Historical S |
| | | ☑ ☑ Create new topic | | | |
| | | Search for match | | | |
| | 22. | Ernst Bacon | | 0.616 | Ernst Bacon |
| | | Choose new match | | | |
| | 25. | Barton Warren Evermann | | 0.609 | Barton Warr Sciences.tx |
| | | Choose new match | | | |
| | 27. | Barton Warren Evermann | | 0.609 | Alice Eastw |

Marcia Zeng, 2018 DIS                    15

---

# Additional Information

- Cogito Intelligence http://www.intelligenceapi.com/demo/ for for Intelligence, Terrorism, Cyber Crime, Crime, Geographic domains, etc.
- NIH's MetaMap: https://metamap.nlm.nih.gov/ for biomedical and health related domains.

## Slide 19

# Behind the theme

- **Five (5) specific taxonomies** of terms (in over 1,000 different categories):
  - •Intelligence,
  - •Terrorism,
  - •Cyber Crime,
  - •Crime, and
  - •Geographic domains
- **A domain ontology** (updated regularly) with a wide range of diverse topics

### COGITO® Intelligence API

LIVE DEMO  HOME  PRODUCT  DEVELOPERS  NEWS ROOM  CONTACTS  GET A KEY

#### Home

Cogito Intelligence API provides full semantic processing features—text mining (with time references (alpha version), semantic reasoning and inferential entities), categorization, semantic tagging, emotions, sentiment, fact mining, writeprint, and extraction relationships between entities that developers can easily integrate into their analysis platforms and applications for faster evaluation and analysis of documents, web pages, social media data or any big data sets or real-time information streams. The API comes in both SOAP XML-based and RESTful JSON-based flavors, and the features include:

- **5 specific taxonomies** of terms (in over 1,000 different categories) for Intelligence, Terrorism, Cyber Crime, Crime and Geographic domains
- A domain ontology **(updated regularly)** with a wide range of diverse topics, for example: weapons, crimes, cyber attacks, points of interest, chemical weapons, controlled substances, terrorist groups, critical infrastructure, world leaders, public companies and more

Semantic Technology for Intelligence
from Expert System
supports analysis and even those without technical know-how
in all phases of the intelligence cycle

**GO TO OUR LIVE DEMO!**
no registration required
**GET A FREE KEY!**
Watch the Cogito Intelligence API video

Copyright © 2014-2016 Expert System S.p.A. | P.I. 02608970360 Privacy Policy | Terms of Service | Expert System

19

## Slide 20

# How can I use these tools in the future?

- For each of the tools, you can apply for a free "key" so that you can conduct analysis of longer text and patch process for folders and files.
- These tools might not be the best for humanities domains. But they (machine) learn quickly.
- Any text submitted will be kept in their corpus, so, **do not** submit any private or unclassified documents!

20

**Anther tool:** NIH's MetaMap: https://metamap.nlm.nih.gov/ for biomedical and health related domains.

Supported by *Metathesaurus*, a large biomedical thesaurus that is organized by concept or meaning, and it links similar names for the same concept from nearly 200 different vocabularies from the world.