

# PABench: Designing a Taxonomy and Implementing a Benchmark for Spatial Entity Matching

**B. Berjawi, F. Duchateau, F. Favetta, M. Miquel, R. Laurini**

Laboratoire d'InfoRmatique en Image et Systèmes d'information

**GEOProcessing 2015**  
**Lisbon, Portugal**



UNIVERSITÉ  
LUMIÈRE  
LYON 2



# Motivation

## ■ Multiplication of Points of Interest (POI) and data sources

- Several Location-Based Services (LBS) providers
- Incomplete, inconsistent, inaccurate, wrong information

## ■ Integration of multiple sources

- Similarity measures
- Probability measures
- Learning-based methods

## ■ How to evaluate and compare spatial integration methods?

# Related Work



## Ontology matching:

- Ontology Alignment Evaluation Initiative (OAEI) [1]



## Schema matching:

- XBenchMatch [2]
- STBenchmark [3]



## Entity matching:

- EMBench [4]

[1]: Ontology Alignment Evaluation Initiative,” URL: <http://oaei.ontologymatching.org>

[2]: F. Duchateau and Z. Bellahsene, “Designing a benchmark for the assessment of schema matching tools,” in Open Journal of Databases (OJDB), vol. 1, no. 1. RonPub, Germany, 2014, pp. 3–25.

[3]: B. Alexe, W. C. Tan, and Y. Velegrakis, “Stbenchmark: towards a benchmark for mapping systems,” Proceedings of the VLDB, vol. 1, no. 1, 2008, pp. 230–244.

[4]: E. Ioannou, N. Rassadko, and Y. Velegrakis, “On generating benchmark data for entity matching,” Journal on Data Semantics, vol. 2, no. 1, 2013, pp. 37–56.

# Related Work

## Spatial Entity matching:

■ Geoddupe [5]

■ Random-spatial-dataset generator [6]

**Need for a spatial entity matching benchmark**

**PABench: Point of Interest Alignment Benchmark**

[5]: H. Kang, V. Sehgal, and L. Getoor, “Geoddupe: A novel interface for interactive entity resolution in geospatial data,” in International Conference on Information Visualisation, 2007, pp. 489–496.

[6]: C. Beeri, Y. Doytsher, Y. Kanza, E. Safra, and Y. Sagiv, “Finding corresponding objects when integrating several geo-spatial datasets,” in ACM International Workshop on Geographic Information Systems, 2005, pp. 87–96.

# Contributions

## ■ Taxonomy of LBS

- Describe the context of LBS providers
- Compare the LBS providers
- Characterize the differences that occur between LBS providers

## ■ Benchmark

- Construct PABench based on the taxonomy characterization
- Generate a characterized training dataset using real data

# Outline

- Introduction
- Related Work
- Taxonomy of LBS
  - Preliminary definitions
  - Differences
- Benchmark
  - Benchmark construction
  - Datasets
- Conclusion and Future Work

# Taxonomy - Preliminary definitions

■ **POI: geographical object described by a set of properties**

$POI = (name, type, coordinates, shape)$

■ **Schema of provider: structure of entities offered by the provider**

I: Internal identifier

L: Spatial attributes

A: Primary terminological

B: Secondary terminological

■ **Entity of POI: instance of a schema and refers to one real-world POI**

$e = \{(id_k:label, id_k:val), (LATITUDE_k:label, LATITUDE_k:val), \dots\}$

# Taxonomy - Preliminary definitions

- Association function  $f$ : returns the POI described by a given entity

$$f: E \rightarrow P$$
$$e \rightarrow f(e) = p$$

- Corresponding entities: two entities from two distinct providers refer to the same POI ( $e_1 \equiv e_2$ )

$$\exists p \in P \setminus f(e_1) = f(e_2) = p$$

- Corresponding attributes: two attributes from two distinct schemas represent the same concept ( $\text{att}_1 \equiv \text{att}_2$ )

# Taxonomy - Differences

Category	Difference
Schema	Attribute Heterogeneity
	Different structure
Terminology	Semantic Different Data (SEM)
	Syntactic Different Data (SYN)
	Missing Data (MD)
	Similar Data (SD)
Spatial	Different locations (DL)
	Equipollent Positions (EP)
	Superposition (SUP)
Availability	Not found POI
	Duplicate Entities

 Differences of corresponding entities

 Differences of non-corresponding entities

# Taxonomy - Example

Entity x (offered by provider 1)	Entity y (offered by provider 2)
EntityID: 51190385	id: fd0cfb424bbd79bf28a832e1764f1c2
Latitude: 48.858606 Longitude: 2.293971 	geometry: { location : { lat : 48.85837, lng: 2.294481}} 
DisplayName: Tour Eiffel EntityTypeID: 7999	name: Eiffel Tower types: establishment
Phone: 0892701239 CountryRegion: FRA Locality: Paris PostalCode: 75007 AddressLine: Champ De Mars, Avenue Anatole France ...	formatted phone number: +33892701239 website: http://www.tour-eiffel.fr formatted address: Champ de Mars, 5 Avenue Anatole France, 75007 Paris, France

**Attribute Heterogeneity**  
 $(\text{att}_i = \text{att}_j) \wedge (\text{att}_i.\text{label} \neq \text{att}_j.\text{label} \vee \text{att}_i.\text{type} \neq \text{att}_j.\text{type})$

# Taxonomy - Differences

Category	Difference
Schema	Attribute Heterogeneity Different structure
Terminology	Semantic Different Data (SEM) Syntactic Different Data (SYN) Missing Data (MD)
	Similar Data (SD)
Spatial	Different locations (DL) Equipollent Positions (EP) Superposition (SUP)
Availability	Not found POI Duplicate Entities



Differences of corresponding entities



Differences of non-corresponding entities

# Taxonomy - Example

Entity x (offered by provider 1)	Entity y (offered by provider 2)
EntityID: 51190385	id: fd0cfb424bbd79bf28a832e1764f1c2
Latitude: 48.858606 Longitude: 2.293971	geometry: { location : { lat : 48.85837, lng: 2.294481}}
DisplayName: Tour Eiffel EntityTypeID: 7999	name: Eiffel Tower types: establishment
Phone: 0892701239 CountryRegion: FRA Locality: Paris PostalCode: 75007 AddressLine: Champ De Mars, Avenue Anatole France ...	formatted phone number: +33892701239 website: <a href="http://www.tour-eiffel.fr">http://www.tour-eiffel.fr</a> formatted address: Champ de Mars, 5 Avenue Anatole France, 75007 Paris, France

## Different Structure

$$\text{att}_i = (\text{att}_1, \text{att}_2, \dots) \vee (\text{att}_1, \text{att}_2, \dots) = \text{att}_j$$

# Taxonomy - Differences

Category	Difference
Schema	Attribute Heterogeneity
	Different structure
Terminology	Semantic Different Data (SEM)
	Syntactic Different Data (SYN)
	Missing Data (MD)
	Similar Data (SD)
Spatial	Different locations (DL)
	Equipollent Positions (EP)
	Superposition (SUP)
Availability	Not found POI
	Duplicate Entities



Differences of corresponding entities



Differences of non-corresponding entities

# Taxonomy - Example

Entity x (offered by provider 1)	Entity y (offered by provider 2)
EntityID: 51190385	id: fd0cfb424bbd79bf28a832e1764f1c2
Latitude: 48,858606 Longitude: 2,293971	geometry: { location : { lat : 48.85837, lng: 2.294481}}
DisplayName: Tour Eiffel EntityTypeID: Touristic place	name: Eiffel Tower types: Landmark - attraction
Phone: 0892701239 CountryRegion: FRA Locality: Paris PostalCode: 75007 AddressLine: Champ De Mars, Avenue Anatole France ...	formatted phone number: +33892701239 website: http://www.tour-eiffel.fr formatted address: Champ de Mars, 5 Avenue Anatole France, 75007 Paris, France

## Semantic and Syntactic Different Data

$$\exists \text{ att}_i \in A_1 \cup B_1, \exists \text{ att}_j \in A_2 \cup B_2 \setminus \\ e1 = e2 \wedge (e1.\text{att}_i = e2.\text{att}_j) \wedge (e1.\text{att}_i.\text{val} = e2.\text{att}_j.\text{val})$$

# Taxonomy - Differences

Category	Difference
Schema	Attribute Heterogeneity
	Different structure
Terminology	Semantic Different Data (SEM)
	Syntactic Different Data (SYN)
	Missing Data (MD)
	Similar Data (SD)
Spatial	Different locations (DL)
	Equipollent Positions (EP)
	Superposition (SUP)
Availability	Not found POI
	Duplicate Entities



Differences of corresponding entities



Differences of non-corresponding entities

# Taxonomy - Example

Entity x (offered by provider 1)	Entity y (offered by provider 2)
EntityID: 51190385	id: fd0cfb424bbd79bf28a832e1764f1c2
Latitude: 48.858606 Longitude: 2.293971	geometry: { location : { lat : 48.85837, lng: 2.294481}}
DisplayName: Tour Eiffel EntityTypeID: Touristic place	name: Eiffel Tower types: Landmark - attraction
Phone: 0892701239 CountryRegion: FRA Locality: Paris PostalCode: 75007 AddressLine: Champ De Mars, Avenue Anatole France	formatted phone number: +33892701239 website: http://www.tour-eiffel.fr formatted address: Champ de Mars, 5 Avenue Anatole France, 75007 Paris, France

## Missing Data

$$\exists \text{ atti} \in A1 \cup B1, \exists \text{ attj} \in A2 \cup B2 \setminus (\text{atti} = \text{attj}) \wedge (\text{e1.atti.val} = \text{NULL} \vee \text{e2.attj.val} = \text{NULL})$$

# Taxonomy - Differences

Category	Difference
Schema	Attribute Heterogeneity
	Different structure
Terminology	Semantic Different Data (SEM)
	Syntactic Different Data (SYN)
	Missing Data (MD)
	Similar Data (SD)
Spatial	Different locations (DL)
	Equipollent Positions (EP)
	Superposition (SUP)
Availability	Not found POI
	Duplicate Entities



Differences of corresponding entities



Differences of non-corresponding entities

# Taxonomy - Example

Entity x (offered by provider 1)	Entity y (offered by provider 2)
EntityID: 51190385	id: fd0cfb424bbd79bf28a832e1764f1c2
Latitude: 48.858606 Longitude: 2.293971	geometry: { location : { lat : 48.85837, lng: 2.294481}}
DisplayName: Tour Eiffel EntityTypeID: Touristic place	name: Eiffel Tower types: Landmark - attraction
Phone: 0892701239 CountryRegion: FRA Locality: Paris PostalCode: 75007 AddressLine: Champ De Mars, Avenue Anatole France	formatted phone number: +33892701239 website: http://www.tour-eiffel.fr formatted address: Champ de Mars, 5 Avenue Anatole France, 75007 Paris, France

## Different Location

$$\begin{aligned} e1 = e2 \wedge (e1.\text{LATITUDE}.val \neq e2.\text{LATITUDE}.val \vee \\ e1.\text{LONGITUDE}.val \neq e2.\text{LONGITUDE}.val) \end{aligned}$$

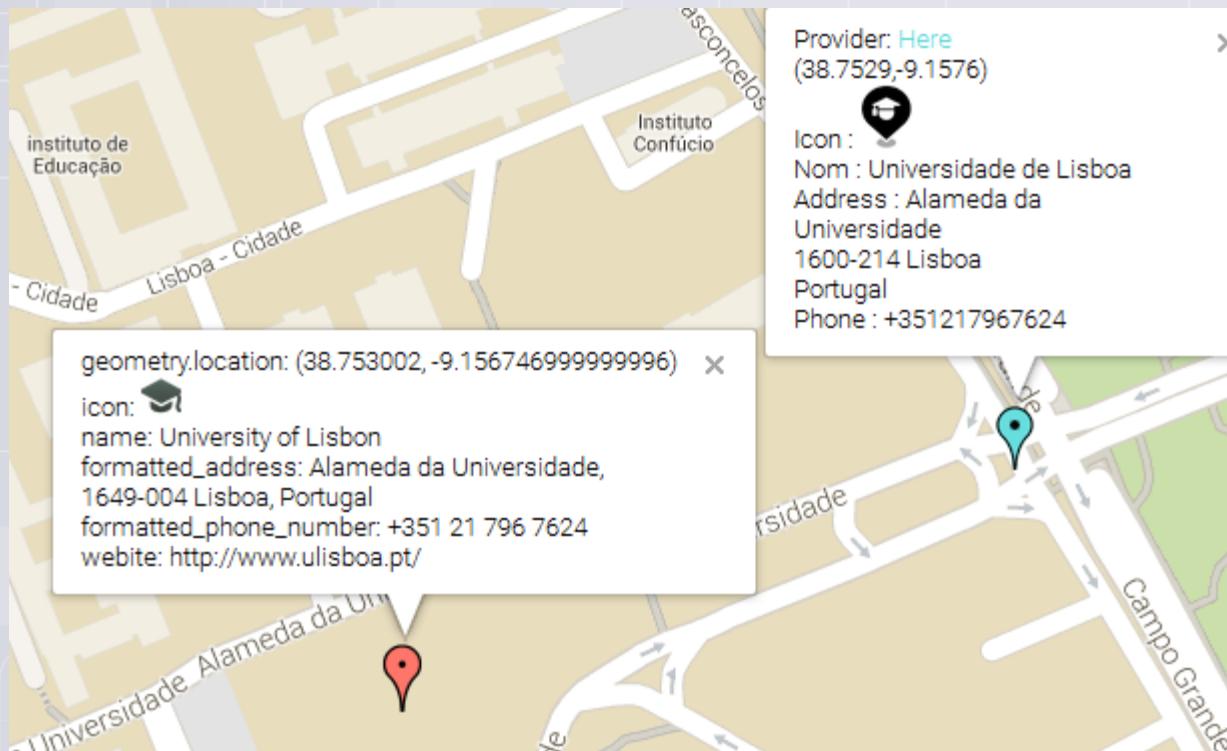
# Taxonomy - Differences

Category	Difference
Schema	Attribute Heterogeneity
	Different structure
Terminology	Semantic Different Data (SEM)
	Syntactic Different Data (SYN)
	Missing Data (MD)
	Similar Data (SD)
Spatial	Different locations (DL)
	Equipollent Positions (EP)
	Superposition (SUP)
Availability	Not found POI
	Duplicate Entities

 Differences of corresponding entities

 Differences of non-corresponding entities

# Taxonomy - Example



## Equipollent Positions

$(e_1 = e_2) \wedge (e_1.\text{LATITUDE}, e_1.\text{LONGITUDE}) \subset p.\text{coordinates} \wedge (e_2.\text{LATITUDE}, e_2.\text{LONGITUDE}) \subset p.\text{coordinates} \wedge (e_1.\text{LONGITUDE}.val \neq e_2.\text{LONGITUDE}.val) \wedge (e_1.\text{LATITUDE}.val \neq e_2.\text{LATITUDE}.val)$

# Taxonomy - Differences

Category	Difference
Schema	Attribute Heterogeneity
	Different structure
Terminology	Semantic Different Data (SEM)
	Syntactic Different Data (SYN)
	Missing Data (MD)
	Similar Data (SD)
Spatial	Different locations (DL)
	Equipollent Positions (EP)
	Superposition (SUP)
Availability	Not found POI
	Duplicate Entities



Differences of corresponding entities



Differences of non-corresponding entities

# Benchmark - Construction

## Differences concerning corresponding entities:

Level	Attributes	Set of possible differences
Spatial	Location	$\emptyset$ , DL, EP
Primary Terminological	Name and Type	$\emptyset$ , SEM, SYN, {SEM, SYN}
Secondary Terminological	Phone, Address, Site, etc.	$\emptyset$ , MD, SEM, SYN, {SEM, SYN, MD}, {SEM, SYN}, {SEM, MD}, {SYN, MD}

96 (3x4x8) distinct situations of differences

Generate a test case for each of the 96 situations

Example of a situation:  $s = \{DL, \{SEM, SYN\}, MD\}$

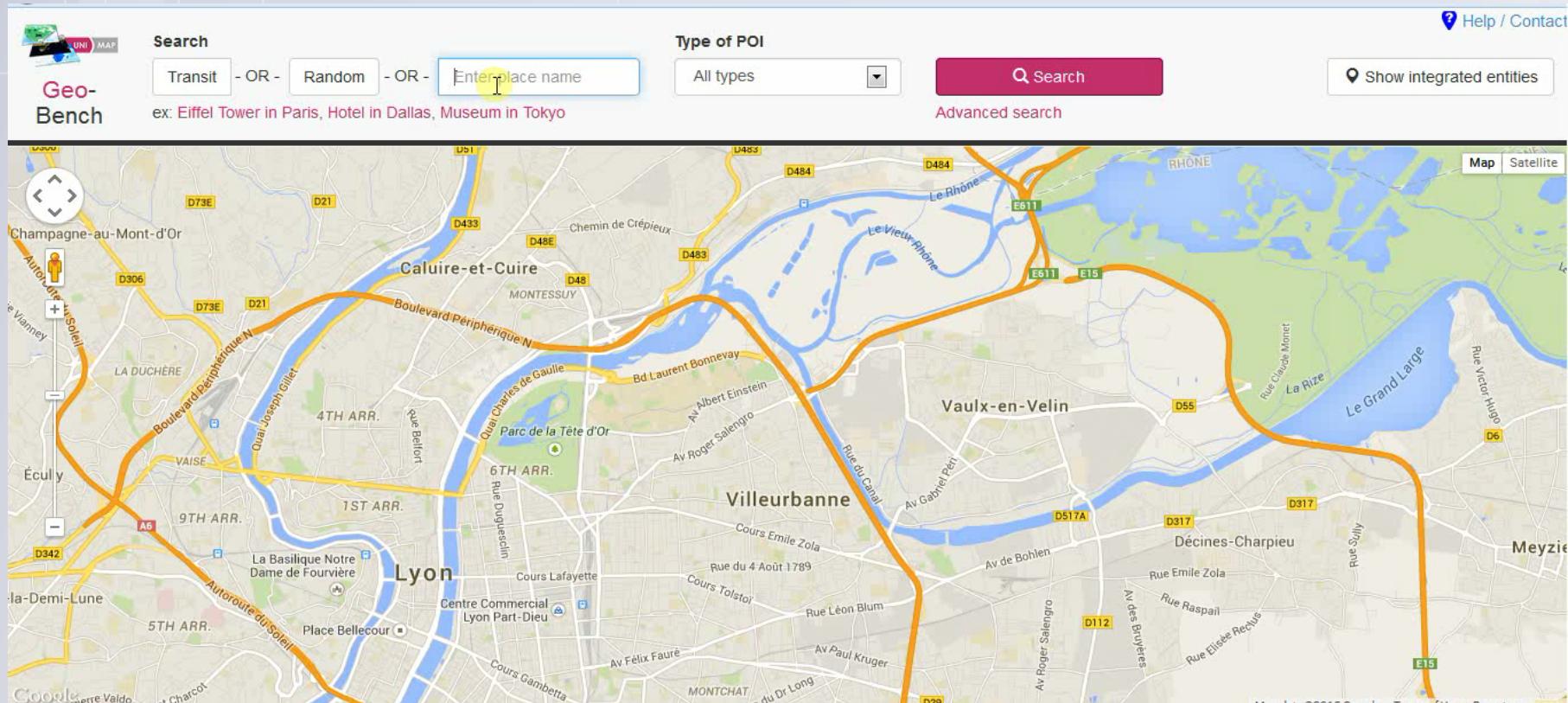
Test\_case( $s$ )= (Source dataset, Target dataset, Ground\_truth)

Remaining differences will be used to add noise

Superposition, Similar Data, Not found POI

# Benchmark - Datasets

## Create a characterized dataset - GeoBench tool [7]



[7]: G. Morana, T. Morel, B. Berjawi, and F. Duchateau, "GeoBench: a Geospatial Integration Tool for Building a Spatial Entity Matching Benchmark (Demo)," in ACM SIGSPATIAL International Conference on Advances in Geographic Information Systems, Dallas, Texas, USA, 4-7 November, 2014, pp. 533-536. <http://tinyurl.com/p3dbmpj>

# Benchmark - Datasets

## Test cases generator

This tool allows you to create a test cases to evaluate an entity matching systems that consists in matching two geographical datasets that contain punctual geographical object. Test cases are generated base on three datasets recovered from three LBS providers.

Number of entities of the first Dataset (E1)	849	Total number of correspondences between E1 and E2	673
Number of entities of the second Dataset (E2)	687	Total number of correspondences between E1 and E3	286
Number of entities of the third Dataset (E3)	314	Total number of correspondences between E2 and E3	277
Total number of entities	1850	Total number of correspondences	1236

A test case consists in evaluating an entity matching system according to a specific situation of differences between corresponding entities. (To see more details about the situations of differences please check this [document](#))

Situations list

<http://tinyurl.com/nc4rurr>

# Conclusion and Future Work

## Contribution:

- Taxonomy that describes LBS context
- Necessary specifications to design PABench
- GeoBench tool to create a characterized dataset and a test case generator

## Future Work:

- Extend PABench by adding more entities
- Create a survey that compares and evaluates existing approaches using PABench
- Extend the taxonomy to cover complex objects



Thank you for  
your attention

Bilal Berjawi  
LIRIS, INSA de Lyon, France  
[bberjawi@liris.cnrs.fr](mailto:bberjawi@liris.cnrs.fr)  
<http://unimap.liris.cnrs.fr>

# Benchmark – Datasets statistics

Dataset	Number of Entities
E1	846
E2	685
E3	314
<b>Total</b>	<b>1845</b>

Number of correspondences	
E1, E2	671
E1, E3	286
E2, E3	277
<b>Total</b>	<b>1234</b>

Situations of differences	Number of correspondences
{EP, SYN, {SYN, MD}}	147
{DL, SYN, {SYN, MD}}	93
{EP, SYN, SYN}	71
{Ø, SYN, {SYN, MD}}	70
{EP, Ø, {SYN, MD}}	63

## Example of differences

Positioning

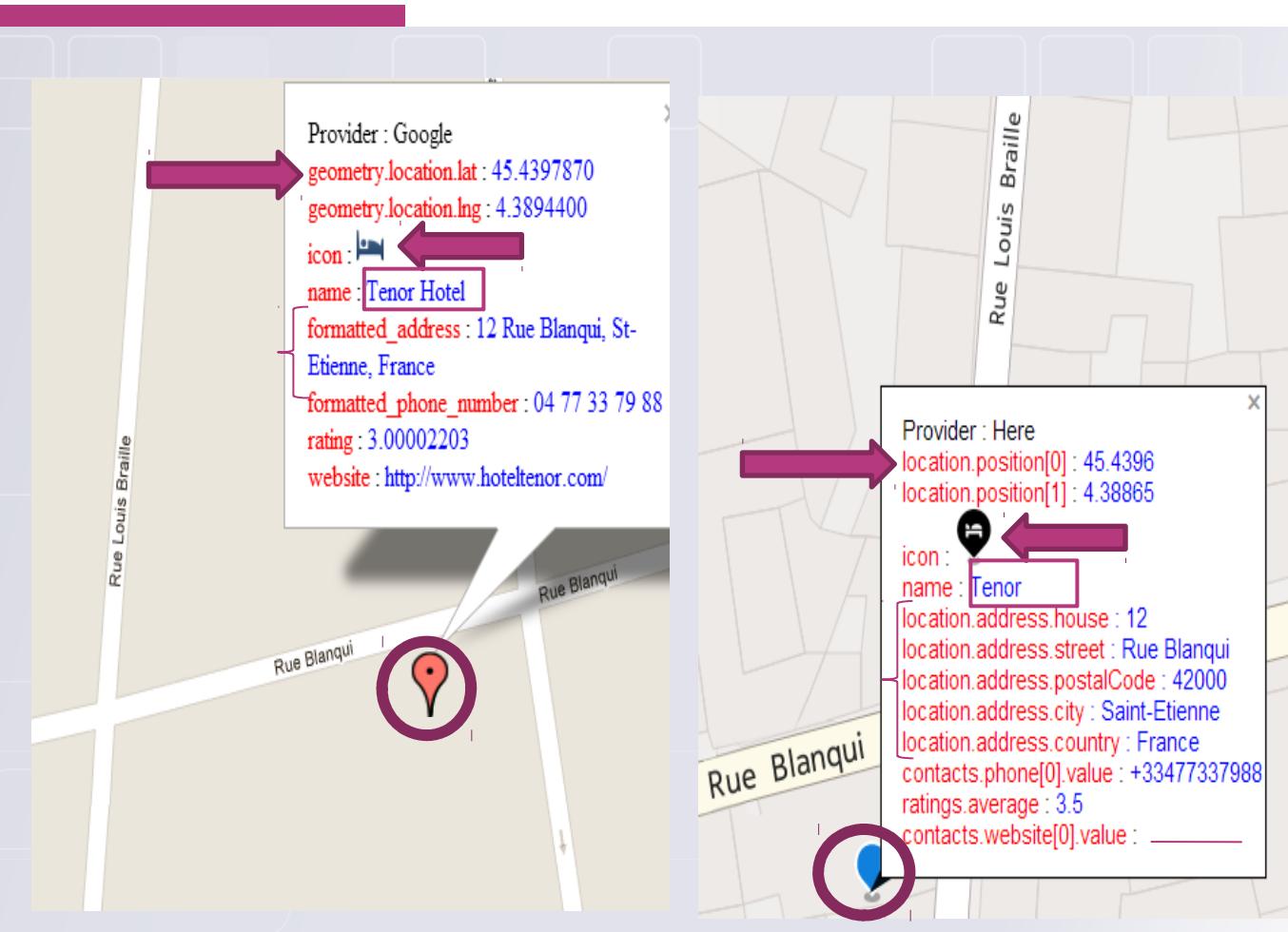
Attributes names

Structure

Legends

Different values

Missing values



# Benchmark - Datasets

## Statistics and Test Cases generator



This tool allows you to punctual geographical

**Statistics**

- Number of entities
- Number of entities
- Number of entities
- Total number of e

**Test Cases**

**Test Case Generator**

You have selected this situation:

Spatial differences: 0  
Primary Terminological differences: 0  
Secondary Terminological differences: SYN ^ MD  
# of Available Correspondences: 39

**Test Case configurations**

Format: SQL \*  
NB of correspondences: 35 (<= 39)\*  
% of Noise entities: 10 %\*

Generate the test case

**or**

Two geographical datasets that contain S providers.

Correspondences between E1 and E2	671
Correspondences between E1 and E3	286
Correspondences between E2 and E3	277
Others	1234

A test case consists in evaluating an entity matching system according to a specific situation of differences between corresponding entities. (To see more details about the situations of differences please check this [document](#))

<http://tinyurl.com/nc4rurr>