

# MIRACLE at GeoCLEF Query Parsing 2007: Extraction and Classification of Geographical Information

Sara Lana-Serrano<sup>1,3</sup>, Julio Villena-Román<sup>2,3</sup>, José Miguel Goñi-Menoyo<sup>1</sup>

<sup>1</sup> Universidad Politécnica de Madrid

<sup>2</sup> Universidad Carlos III de Madrid.

<sup>3</sup> DAEDALUS - Data, Decisions and Language, S.A.

slana@diatel.upm.es, jvillena@daedalus.es, josemiguel.goni@upm.es

## Abstract

This paper describes the participation of MIRACLE research consortium at the Query Parsing task of GeoCLEF 2007. Our system is composed of three main modules. First, the Named Geo-entity Identifier, whose objective is to perform the geo-entity identification and tagging, i.e., to extract the “where” component of the geographical query, should there be any. This module is based on a gazetteer built up from the Geonames geographical database and carries out a sequential process in three steps that consist on geo-entity recognition, geo-entity selection and query tagging. Then, the Query Analyzer parses this tagged query to identify the “what” and “geo-relation” components by means of a rule-based grammar. Finally, a two-level multiclassifier first decides whether the query is indeed a geographical query and, should it be positive, then determines the query type according to the type of information that the user is supposed to be looking for: map, yellow page or information. According to a strict evaluation criterion where a match should have all fields correct, our system reaches a precision value of 42.8% and a recall of 56.6% and our submission is ranked 1<sup>st</sup> out of 6 participants in the task. A detailed evaluation of the confusion matrixes reveal that some extra effort must be invested in “user-oriented” disambiguation techniques to improve the first level binary classifier for detecting geographical queries, as it is a key component to eliminate many false-positives.

## Categories and Subject Descriptors

**H.3 [Information Storage and Retrieval]:** H.3.1 Content Analysis and Indexing; H.3.2 Information Storage; H.3.3 Information Search and Retrieval; H.3.4 Systems and Software; H.3.7 Digital libraries. **H.2 [Database Management]:** H.2.5 Heterogeneous Databases; H.2.8 Database Applications - *Spatial databases and GIS*.

## Keywords

Linguistic Engineering, classification, geographical IR, geographic entity recognition, gazetteer, semantic expansion, Wordnet.

## 1. Introduction

The goal of Geographical Information Retrieval (GIR) is to deal with those information retrieval problems that contain some kind of spatial awareness, i.e., that include geographical references (georeferences) which are essential for the meaning of the query, for example: “find me nice and cheap hotels near Madrid”. It is a complex task because of its strong dependence on geographical information resources, which tend to be incomplete and inexact. Moreover, geographical information is mainly arranged in a tree-like hierarchy, therefore queries usually imply a multilevel search (for example: “give me documents about villages in Northern Spain”). Finally, additional translation problems arise when dealing with multiple languages, due to the lack of specific and specialized translation resources in a worldwide domain.

GeoCLEF is the cross-language geographic retrieval track that runs as part of the Cross Language Evaluation Forum (CLEF) campaign, whose aim is to provide with the necessary framework in which to evaluate GIR systems for search tasks involving both spatial and multilingual aspects. This year, apart from the traditional task, GeoCLEF 2007 offered the Query Parsing task.

A geographic query is usually composed of three components, “what”, “geo-relation” and “where”. The keywords in “what” indicate what users want to find; “where” refers to their geographic area of interest; and

“geo-relation” stands for the relationship between “what” and “where”. For instance, in the first example, “what” would be “nice and cheap hotels”, “where” would be “Madrid”, and “geo-relation” would be “NEAR”. Note that “Madrid” is itself ambiguous and can refer not only to the capital of Spain or the autonomous region where the city of Madrid is located, but also other cities or administrative divisions in United States, Philippines, Mexico, Argentina, Equatorial Guinea, Colombia, Dominican Republic, Sweden...

The key problem for GIR is to understand how to parse and extract those key components from the queries. This is the objective of the Query Parsing task. Participants were given a set of 800,000 untagged queries and had to detect whether each query was a geographical query or not, and, should the result be positive, had to extract the three components: “where” (with its corresponding latitude/longitude), “geo-relation” (normalized into a predefined relation type such as IN, NEAR, FROM, TO, NORTH\_OF...) and “what” (categorized into a type of request: “map”, “yellow page” or “information”).

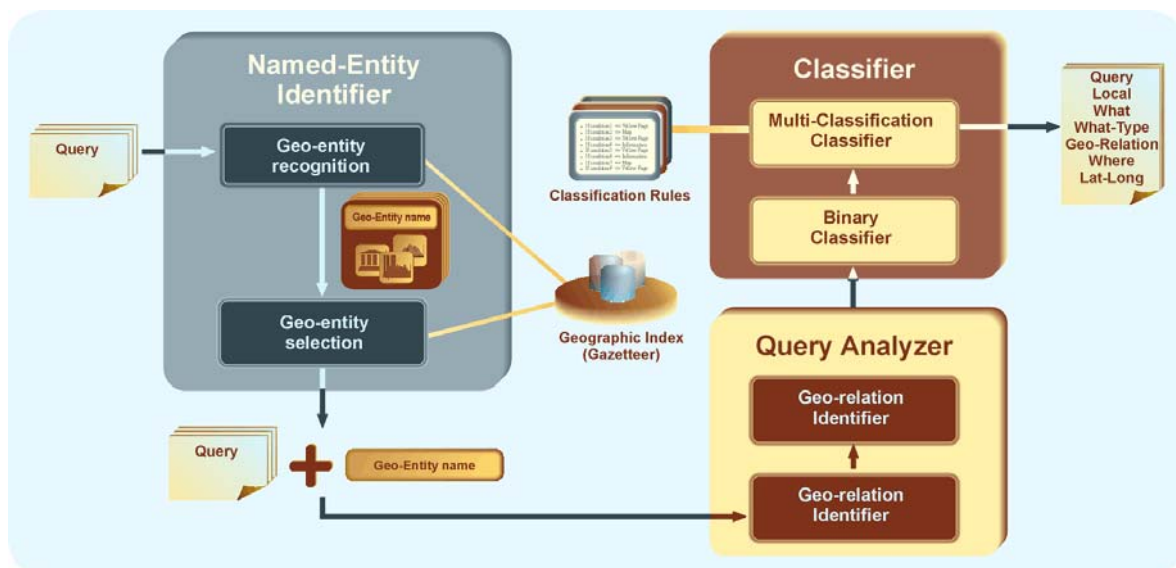
The MIRACLE team is a research consortium formed by research groups of three different universities in Madrid (Universidad Politécnica de Madrid, Universidad Autónoma de Madrid and Universidad Carlos III de Madrid) along with DAEDALUS, a small/medium size enterprise (SME) founded in 1998 as a spin-off of two of these groups and a leading company in the field of linguistic technologies in Spain. MIRACLE has taken part in CLEF since 2003 in many different tracks and tasks, including the main bilingual, monolingual and cross lingual tasks [4] as well as in ImageCLEF, Question Answering, WebCLEF and GeoCLEF [5] [6] tracks.

This paper describes the MIRACLE participation at the Query Parsing task of GeoCLEF 2007. In the following sections, we will first give an overview of the architecture of our system. Afterwards we will elaborate on the different modules. Finally, the results will be presented and analyzed.

## 2. System Description

**Figure 1** presents the system architecture. Observe that the approach consists of three sequential tasks executed by independent modules:

- **Named Geo-entity Identifier:** performs geo-entity identification and query expansion.
- **Query Analyzer:** identifies the “what” and “geo-relation” components of a geographical query.
- **Query Type Classifier:** determines the type of geographical query.



**Figure 1.** Overview of the system.

### 2.1. Named Geo-entity Identifier

The objective of this module is to perform the geo-entity identification and tagging, i.e., to extract the “where” component of the query, should there be any. It is composed of two main components: a geo-entity parser based on a gazetteer, i.e. a database with geographical resources that constitutes the knowledge base of the system.

Our gazetteer is built up from the Geonames geographical database [2], available free of charge for download under a creative commons attribution license. It contains over 8 million geographical names with more than 6.5 million unique features about 2.2 million populated places and 1.8 million alternate names. Those features include a unique identifier, the resource name, alternative names (in other languages), county/region, administrative divisions, country, continent, longitude, latitude, population, elevation and timezone. All features are categorized into one out of 9 feature classes and further subcategorized into one out of 645 feature codes.

Geonames integrates geographical data (such as names of places in various languages, elevation or population) from various sources, mainly the Geonet Names Server (GNS) [9] gazetteer of the National Geospatial Intelligence Agency (NGA), the Geographic Names Information System (GNIS) [8] gazetteer of the U.S. Geographic Survey, the GTOPO30 [3] digital elevation model for the world developed by United States Geological Survey (USGS) and Wikipedia, among others.

For our purposes, all data was loaded and indexed in a MySQL database, although not all fields (such as time zone or elevation) were to be used: the relevant fields are UFI (unique identifier), NAME\_ASCII (name), NAME\_ALTERNATE (alternate names), COUNTRY, ADM1 and ADM2 (administrative region where the entity is located), FEATURE\_CLASS, FEATURE\_TYPE, POPULATION, LATITUDE and LONGITUDE. To simplify the queries, each row is complemented with the expansion of country codes (ES→Spain) and/or state codes (NC→North Carolina) –when applicable. The final database uses 865KB.

The geo-entity parser carries out the following tasks:

- **Geo-entity recognition:** identifies named geo-entities [6] using the information stored in the gazetteer, looking for candidate named entities matching any substring of one or more words [1] included in the query and not included in a stopword (or stop-phrase) list [7].

The stopword list is mainly automatically built by extracting those words that are both common nouns and also georeference entities, assuming that the user is asking for the common noun sense (for example, “Aguilera” –for Christina Aguilera– or “tanga” –thong). Specifically we have used lexicons for English, Spanish, French, Italian, Portuguese and German, and have selected words that appear at least with a certain frequency in the query collection. The final stopword list contains 1712 entries.

- **Geo-entity selection:** The selected named geo-entity will be the one with the longest number of words and, if the same, the one with higher score. The score is computed according to the type of geographic resource (country, region, county, city...) and its population, as shown in the following table.

**Table 1.** Entity score.

Feature type	Code	Score
Capital and other big cities	PPLA, PPLC, PPLG	Population+100,000,000
Political entities	PCL, PCLD, PCLF, PCLI, PCLIX, PCLS	Population+10,000,000
Countries	A	Population+1,000,000
Other cities	PP, STLMT	Population+100,000
Other	*	Population
For all cities, if country/state name/code is also in the query	PP, STLMT	Score += 100,000,000

Those values were arbitrarily chosen after different manual executions and subsequent analysis.

- **Query tagging:** expands the query with information about the selected entity: name, country, longitude, latitude, and type of geographic resource.

The output of this module is the list of queries in which a possible named geo-entity has been detected, along with its complete tagging. For example:

Query| score|ufi|entity|state (code)|country (code)|latitude|longitude|feature\_class|feature\_type  
 airport {{alicante}} car rental week|2693959|2521976|Alicante||Spain (ES)|38.5|-0.5|A|ADM2  
 bedroom apartments for sale in {{bulgaria}}|10000000|732800||Bulgaria (BG)|43.0|25.0|A|PCLI  
 hotels in {{south lake tahoe}}|123925|5397664|South Lake Tahoe|California (CA)|United States (US)|38.93|-119.98|P|PPL  
 helicopter flight training in southwest {{florida}}|100100000|4920378|Florida|Indiana (IN)|United States (US)|40.16|-85.71|P|PPL

Observe that the geo-entity is specifically marked in the original query, enclosed between double curly brackets, to help the following module to identify the rest of the components of the geographical query.

## 2.2. Query Analyzer

This module parses each previously tagged query to identify the “what” and “geo-relation” components of a geographical query, sorting out the named geo-entity detected by the previous module, enclosed between curly brackets.

It consists of two subsystems:

- **Geo-relation identifier:** identifies and qualifies spatial relationships supported by a regular expression rules based. Its output is the input list of queries expanded with information related to the identified “geo-relation”.

For instance, continuing with the previous examples, the output would be the following:

```
Query/geo-relation/entity/state/country/country (code)/latitude/longitude/feature_class/feature_type  
airport {{alicante}} car rental week|NONE|Alicante|Spain|ES|38.5|-0.5|A|ADM2  
bedroom apartments for sale #@IN#@# {{bulgaria}}|IN|Bulgaria|Bulgaria|BG|43.0|25.0|A|PCLI  
hotels #@IN#@# {{south lake tahoe}}|IN|South Lake Tahoe|California|United States|US|38.93|-  
119.98|P|PPL  
helicopter flight training in #@SOUTH_WEST_OF#@# {{florida}}|SOUTH_WEST_OF|Florida|  
Indiana|United States|US|40.16|-85.71|P|PPL
```

Observe that the geo-relation is also marked in the original query.

- **Concept identifier:** analyses the output of the previous step and extracts the “what” component of a geographical query applying manually defined grammar rules based on the identified “where” and “geo-relation” components.

## 2.3. Query Type Classifier

Finally, the last step is to decide whether the query is indeed a geographical query and, should it be positive, to determine the type of query, according to the type of information that the user is supposed to be looking for:

- Map type: users are looking for natural points of interest, like rivers, beaches, mountains, monuments...
- Yellow page type: businesses or organizations, like hotels, restaurants, hospitals, etc.
- Information type: users are looking for text information, like news, articles, blogs, and so on.

The process is carried out by a two level classifier:

1. **First level:** a binary classifier to determine whether a query is a geographical or a non-geographical query. This simple classifier is based on the assumption that a query is geographical if the “where” component is not empty.
2. **Second level:** a multi-classification rule-based classifier to determine the type of geographical query. The multi-classifier treats the tagged queries as a lexicon of semantically related terms (words, multi-words and query parts).

The classification algorithm applies a knowledge base that consists on a set of manually defined grammar rules, including nouns and grammatically related part-of-speech categories as well as the type of geographical resource. The different valid lemmas are unified using Wordnet synsets.

## 3. Results

For the evaluation, multiple human editors labeled 500 queries that were chosen to represent the whole query set. Then all the submitted results were manually compared to those queries following a strict criterion where a match should have all fields correct.

**Table 2** shows the evaluation results of our submission, using the well-known evaluation measures of precision, recall and F1-score.

**Table 2.** Overall results.

Precision <sup>(1)</sup>	Recall <sup>(2)</sup>	F1-score <sup>(3)</sup>
0.428	0.566	0.488

$$^{(1)} \text{precision} = \frac{\text{correctly\_tagged\_queries}}{\text{all\_tagged\_queries}}$$

$$^{(2)} \text{recall} = \frac{\text{correctly\_tagged\_queries}}{\text{all\_relevant\_queries}}$$

$$^{(3)} \text{F1-score} = \frac{2 * \text{precision} * \text{recall}}{\text{precision} + \text{recall}}$$

According to the task organizers, our submission achieved the best performance out of the 8 submissions of this year, which was a good reward for our hard work.

As participants in the task were provided with the evaluation data set, we have further evaluated our submission to separately study the results for each component of the geographical queries and also the performance level-by-level of the final classifier.

**Table 3** shows the individual analysis of the classifier per each extracted field. The first-level classifier achieves a precision of 75.40%. However, the second-level classifier reduces this value to 56.20% for the WHAT-TYPE feature. According to a strict evaluation criterion, this would be the precision of the overall experiment. If evaluated only over well-classified (geographical/non geographical) queries, the precision arises to 74.53%.

**Table 3.** Individual analysis per component.

	LOCAL		WHAT		WHAT-TYPE		WHERE		ALL	
	Total	%	Total	%	Total	%	Total	%	Total	%
<b>All topics</b>	377	75.40	323	64.60	281	56.20	321	64.20	259	51.80
<b>Well-classified</b>	377	100.00	323	85.67	281	74.53	321	85.15	259	68.70

The confusion matrix for the first-level classifier is shown in **Table 4**.

**Table 4.** Confusion matrix for the binary classifier.

	LOCAL YES	LOCAL NO	Precision <sup>(1)</sup>	Recall <sup>(2)</sup>	Accuracy <sup>(3)</sup>
	<b>ASSIGNED YES</b>	297			
<b>ASSIGNED NO</b>	12	80			

$$^{(1)} \text{precision} = \frac{\text{TP}}{\text{TP} + \text{FP}} \quad ^{(2)} \text{recall} = \frac{\text{TP}}{\text{TP} + \text{FN}} \quad ^{(3)} \text{accuracy} = \frac{\text{TP} + \text{TN}}{\text{TP} + \text{TN} + \text{FP} + \text{FN}}$$

**Table 5** (a, b, c) presents the confusion matrixes for the multiclassifier, individualized per class and calculated over all topics.

**Table 5a.** Confusion matrix for the multiclassifier, “Yellow Page” class.

	Yellow-Page YES	Yellow-Page NO	Precision	Recall	Accuracy
	<b>ASSIGNED YES</b>	142			
<b>ASSIGNED NO</b>	7	161			

**Table 5b.** Confusion matrix for the multiclassifier, “Map” class.

	Map YES	Map NO	Precision	Recall	Accuracy
	<b>ASSIGNED YES</b>	45			
<b>ASSIGNED NO</b>	41	398			

**Table 5c.** Confusion matrix for the multiclassifier, “Information” class.

	Information YES	Information NO	Precision	Recall	Accuracy
ASSIGNED YES	14	1	0.93	0.20	0.88
ASSIGNED NO	57	428			

**Table 6** (a, b, c) presents the same confusion matrixes per class, but calculated only over topics which are correctly classified by the first level binary classifier.

**Table 6a.** Confusion matrix for the multiclassifier, “Yellow Page” class.

	Yellow-Page YES	Yellow-Page NO	Precision	Recall	Accuracy
ASSIGNED YES	142	92	0.61	0.99	0.75
ASSIGNED NO	2	141			

**Table 6b.** Confusion matrix for the multiclassifier, “Map” class.

	Map YES	Map NO	Precision	Recall	Accuracy
ASSIGNED YES	14	0	0.92	0.55	0.89
ASSIGNED NO	54	309			

**Table 6c.** Confusion matrix for the multiclassifier, “Information” class.

	Information YES	Information NO	Precision	Recall	Accuracy
ASSIGNED YES	14	0	1.00	0.21	0.86
ASSIGNED NO	54	309			

#### 4. Conclusions and Future Work

We have however some disagreements with the evaluation data provided by the organizers. Although some of them may be actual errors, most are due to the complexity and ambiguity of the queries. **Table 7** shows some examples of queries that have been classified as geographical by our system but have been evaluated as false-positives. In fact, we think that it would be almost impossible to reach a complete agreement in the parsing or classification for every case among different human editors. The conclusion to be drawn from this is that the task to analyze and classify queries is very hard without a previous contact and without the possibility of interaction and feedback with the user.

**Table 7.** Some examples of ambiguities.

QueryNo	Query	Extracted “where”	Why not?
113501	calabria chat	calabria, Italy	chat rooms about the region of Calabria?
443245	Machida	machida, Japan	Hiroko Machida (actress), Kumi Machida (artist) or the city of Machida?
486273	montserrat reporter	montserrat, Montserrat	online newspaper or reporters in Montserrat?

The analysis of the confusion matrixes for the multiclassifier that are calculated over the topics correctly classified by the first level classifier shows that the probability that a geographical query is classified as “Yellow Page” is very high. This could be related to the uneven distribution of topics (almost 50% of the geographical queries belong to this class). In addition, “Information” type queries have a very low recall. These combined facts point out that the classification rules have not been able to establish a difference between both classes. We will focus on this issue in future participations. Moreover, we will try to incorporate some “user-oriented”

disambiguation techniques to improve the first level binary classifier, as it is a key component to eliminate many false-positives.

## Acknowledgements

This work has been partially supported by the Spanish R&D National Plan, by means of the project RIMMEL (Multilingual and Multimedia Information Retrieval, and its Evaluation), TIN2004-07588-C03-01; and by the Madrid's R&D Regional Plan, by means of the project MAVIR (Enhancing the Access and the Visibility of Networked Multilingual Information for the Community of Madrid), S-0505/TIC/000267.

## References

- [1] Charniak, Eugene. A Maximum-Entropy-Inspired Parser. In Proceedings of NAACL-2000, 2000.
- [2] Geonames geographical database. On line <http://www.geonames.org> [Visited 14/08/2007].
- [3] Global 30 Arc-Second Elevation Data Set. On line <http://eros.usgs.gov/products/elevation/gtopo30.html> [Visited 14/08/2007].
- [4] Goñi-Menoyo, J.M.; González-Cristóbal, J.C.; and Villena-Román, J. MIRACLE at Ad-Hoc CLEF 2005: Merging and Combining without Using a Single Approach. Accessing Multilingual Information Repositories: 6th Workshop of the Cross Language Evaluation Forum 2005, CLEF 2005, Vienna, Austria, Revised Selected Papers (Peters, C. et al., Eds.). Lecture Notes in Computer Science, vol. 4022, Springer, 2006.
- [5] Lana-Serrano, S.; Goñi-Menoyo, J.M.; and González-Cristóbal, J.C. MIRACLE at GeoCLEF 2005: First Experiments in Geographical IR. Accessing Multilingual Information Repositories: 6th Workshop of the Cross Language Evaluation Forum 2005, CLEF 2005, Vienna, Austria, Revised Selected Papers (Peters, C. et al., Eds.). Lecture Notes in Computer Science, vol. 4022, pp. 920-923. Springer, 2006.
- [6] Goñi-Menoyo, J.M.; González-Cristóbal, J.C.; Lana-Serrano, S.; Martínez-González A. MIRACLE's Ad-Hoc and Geographic IR approaches for CLEF 2006: 7th Workshop of the Cross Language Evaluation Forum 2006, CLEF 2006, Alicante, Spain, Revised Selected Papers (Peters, C. et al., Eds.). Lecture Notes in Computer Science.
- [7] University of Neuchatel. Page of resources for CLEF (Stopwords, transliteration, stemmers ...). On line <http://www.unine.ch/info/clef> [Visited 18/07/2006].
- [8] U.S. Geological Survey. On line <http://www.usgs.gov> [Visited 14/08/2007].
- [9] U.S. National Geospatial Intelligence Agency. On line <http://www.nga.mil> [Visited 14/08/2007].