# GIR Experimentation

Andogah Geoffrey

Computational Linguistics Group

Centre for Language and Cognition Groningen (CLCG)

University of Groningen

Groningen, The Netherlands

g.andogah@rug.nl, annageof@yahoo.com

**Abstract**

Geographic Information Retrieval (GIR) community has generally accepted the thesis that both thematic and geographic aspect of documents are useful for GIR. This paper describes a preliminary experiment exploring this thesis by seperately indexing/searching geographical relevant-terms (place names, geo-spatial relations, geographic concepts and geographic adjectives) extracted from reference document collection. Two indexes were created one for extracted geographic relevant-terms (i.e. document footprint) and one for reference document collections. Geo-Score and Thematic-Score against document collection footprint and reference document collection respectively were combined through a linear interpolation to obtained the final score for document relevance ranking. We used several freely available geographic resources – Wikipedia, World-Gazetteer, GEOnet Name Server (GNS), and WordNet. Apache Lucene was used as an indexing and search platform while Alias-I LingPipe was used to detect geographic named entities (GNEs), and other geo-relevant concepts and terms in documents. We submitted runs for monolingual English task, and our system achieved mean average precision (MAP) of 0.1690 to 0.2194. No significant improvement was observed through geographic query expansion.

## Categories and Subject Descriptors

H.3 [**Information Storage and Retrieval**]: H.3.1 Content Analysis and Indexing; H.3.3 Information Search and Retrieval; H.3.4 Systems and Software; H.3.7 Digital Libraries; H.2.3 [**Database Managment**]: Languages—*Query Languages*

## General Terms

system architecture, performance, experimentation

## Keywords

geographic information retrieval, geographic query expansion, geographic named entity tagging, document footprints, geographic knowledge base, thematic score, geographic score, linear interpolation, relevance ranking

## 1 Introduction

Geographic Information Retrieval (GIR) concerns the retrieval of information involving some kind of spatial awareness. Geographic information pervades many documents, and therefore, geographic

references may be important for Information Retrieval (IR). Additionally, many documents contain geographic references expressed in multiple languages which may or may not be the same as the query language [Gey et al, 2005].

To perform GIR both thematic (non-geographic aspect) and geographic aspect of documents need consideration. In order to approach this thesis we derive for each document in the collection a corresponding document footprint containing place names (e.g. Uganda), geo-spatial relations (e.g. west of), geographic concepts (e.g. country) and geographic adjectives (e.g. Ugandan). The document footprint and reference document collection provided were separately indexed and searched. Geo-Score and Thematic-Score against document collection footprint and reference document collection were combined through a linear interpolation to obtained the final score to perform document relevance ranking. Queries were performed for geographic relevant terms identified in topics against document collection footprints and reference document collection provided to investigate impart of geo-references and geo-relevant terms for GIR.

Freely available geographic resources (from: Wikipedia[1], World-Gazetteer[2], GEOnet Name Server[3] (GNS), WordNet[4]) were consulted for query geographic reference expansion. Apache Lucene[5] was used as an indexing and search platform while Alias-I LingPipe[6] was used to detect geographic named entities (GNEs) and other geo-relevant concepts and terms in documents.

# 2   GeoCLEF 2006

GeoCLEF evaluation track was run for the first time at CLEF 2005 to evaluate retrieval of multilingual documents with an emphasis on geographic search [Gey et al, 2005]. As GeoCLEF 2005, GeoCLEF 2006 outline the following challenges to GIR in a multilingual environment: (1) translation of locations (e.g. Uganda (EN) to Oeganda (NL)), (2) resolution of geographic reference ambiguities (e.g. "Jack *London*" the author not a place; South Yorkshire and S. Yorks refer to the same place), (3) resolution of spatial ambiguity (e.g. Sheffield in UK or USA), (4) finding or creating suitable multilingual geographic knowledge base, and (5) combining both text and spatial retrieval methods. The specific aims for GeoCLEF 2006 are: (1) compare methods of query translation, (2) query expansion, (3) translation of geographical references, (4) use of text and spatial retrieval methods separately or combined, and (5) retrieval models and indexing methods.

GeoCLEF 2006 consists of document collections in English, German, Portuguese and Spanish, and 25 search topics in these languages. The tasks for GeoCLEF 2006 are: (1) monolingual retrieval – retrieval where the topic and document languages are the same, and (2) bilingual retrieval – cross-language retrieval where the topic language is different from the document language, i.e. $X \rightarrow \{DE, EN, ES, PT\}$. For each document language, participants may submit the results of up to 10 runs: 5 monolingual and 5 bilingual. Two of these runs are required: (1) Title-Description – where the search queries are created using only the contents of the Title and Desc tags of the topic, and (2) Title-Description-Narrative – where the search queries are created using the contents of the Title, Desc and Narr tags from the topic. The Narrative tag contains a more comprehensive description of the information request defined by the topic, including specifics about the geography of the topic such as a list of desired cities, states, countries, rivers or latitudes and longitudes. An example search topic is depicted below:

> *<top>*
> *<num>GC027</num>*
> *<EN-title>Cities within 100km of Frankfurt</EN-title>*
> *<EN-desc>Documents about cities within 100 kilometers of the city of Frankfurt in Western Germany</EN-desc>*

---

[1]http://www.wikipedia.org
[2]http://www.world-gazetteer.com
[3]http://earthinfo.nga.mil/gns/html
[4]http://wordnet.princeton.edu
[5]http://jakarta.apache.org/lucene
[6]http://alias-i.com/lingpipe

*<EN-narr>Relevant documents discuss cities within 100 kilometers of Frankfurt am Main Germany, latitude 50.11222, longitude 8.68194. To be relevant the document must describe the city or an event in that city. Stories about Frankfurt itself are not relevant</EN-narr>*
*</top>*

# 3    Our appraoch

We are participating in GeoCLEF evaluation track at CLEF 2006 for the first time. The main motivation for our participation is to experiment with both thematic and geographic aspect of a document for GIR. In this section we describe our approach and resources used. Our appraoch borrows techniques from (Larson [2005], Ferres et al [2005], Hughes [2005], Buscaldi et al [2005], Gey and Vivien [2005], Leidner [2005]) with few exceptions such as the creation of an index of document collection footprint along side the index of reference document collection, and thereby combining query results of the two index searches using linear interpolation.

## 3.1    Resources

### 3.1.1    Geographic Knowledge Base

We used the World Gazetteer, GEOnet Names Server (GNS), Wikipedia and WordNet as the bases for our Geographic Knowledge Base (GKB) for several reasons – free availability, multilingual (English, Germany, Portuguese and Spanish), most popular and major places, etc. Volcano active region, European river, Atlantic Ocean ports/coast and European Wine processing region information were specifically gathered from the Wikipedia.

### 3.1.2    GeoTagger

Alias-I LingPipe was used to detect named entities (location, person and organisation), geographic concepts (continent, region, country, city, town, village, etc.), spatial relations (near, in, south of, north west, etc.) and locative adjectives (e.g. Ugandan).

### 3.1.3    GeoCoder

We used a simple appraoch to geo-code identified geographic named entities (GNEs) presented in CLIN 2005 [Andogah, 2005]. The approach exploits location type (e.g. city, mountain) and hierarchy information integrated in GKB to ground GNEs.

### 3.1.4    Lucene Search Engine

Apache Lucene is a high-performance, full-featured text search engine library written entirely in Java. It is a technology suitable for nearly any application that requires full-text search, especially cross-platform. Lucene's default similarity measure is based on vector space model[7] (VSM).

## 3.2    Document Pre-processing

Documents were pre-processed using the Alias-I LingPipe to detect place names (e.g. Kampala), geographic concepts (e.g. city), spatial relations (e.g. west of) and adjectives referring to things or people or language connected to a place (e.g. Ugandan).

Candidate locations for detected place names is obtained from GKB. Place names are resolved to their respective locations using a simple geo-coding approach exploiting location type and hierarchical information present in GKB. The preliminary experimental result of geo-coding approach

---

[7]The vector space model (VSM) is an algebraic model used for information filtering and information retrieval. It represents natural language documents in a formal manner by the use of vectors in a multi-dimensional space. $http://en.wikipedia.org/wiki/Vector\_space\_model$

used here was reported in [Andogah, 2005][8]. However, due to time limitation geo-coding task was not experimented as planned, instead we assume that all geo-relevant terms detected in a document will some-how relate or point to a specific geographic region/scope or geographic concept in the discourse.

## 3.3 Indexing document collection

Footprint document collection repository derived from document collection was created. Footprint documents contain geo-relevant terms such as place name, geographic concepts, spatial relations, locative adjectives plus their respective term frequency as depicted below.

> *<GeographicTermFrequency docid="GH950102-000006">*
> *<GT name="east" tf="2" gtt="SPR" />*
> *<GT name="america" tf="1" gtt="LOC" />*
> *<GT name="new york" tf="2" gtt="LOC" />*
> *<GT name="buffalo" tf="3" gtt="LOC" />*
> *<GT name="orlando magic" tf="1" gtt="LOC" />*
> *<GT name="american" tf="4" gtt="GAD" />*
> *<GT name="texas" tf="1" gtt="LOC" />*
> *<GT name="buffalo jills" tf="1" gtt="LOC" />*
> *<GT name="city" tf="1" gtt="GCO" />*
> *</GeographicTermFrequency>*

Derived footprint documents were indexed using Lucene along side index of reference document collection provided for the experiment (see [Table 1] for details).

Table 1: Footprint document index structure

| Field | Lucene Type | Description |
|---|---|---|
| nm | Field.Keyword | Geo-relevant term e.g. Kampala, city, west, Ugandan |
| tf | Filed.Keyword | Geo-relevant term frequency |
| gtt | Field.Keyword | Geo-relevant term type e.g. LOC (location), GCO (concept), SPR (spatial relation), GAD (geo-adjective) |
| docid | Field.Keyword | Document unique identification/number |

Geo-relevant term frequency is factored into the index by adding the same geo-relevant-term to the index the number of times it occurences in the document(e.g. american in the above sample footprint document is added 4 times during indexing).

Reference document collection provided for experimentation were indexed using Lucene. Document HEADLINE and TEXT contents were combined to created document content for indexing (see [Table 2] for details).

## 3.4 Querying document collection

Mandatory runs 1 and 2 queries were formulated by topic TITLE-DESC (CLCGGeoEE1) and TITLE-DESC-NARR (CLCGGeoEE2) contents respectively. These queries were submitted to search reference document index (Lucene field content [Table 2]). The mandatory queries perform general-purpose search of Lucene index returning the top 1,000 documents retrieve.

Our third run query was formulated by topic TITLE-DESC only (CLCGGeoEE5). The query was submitted to search footprint document index (Lucene field nm [Table 1]) and reference

---

[8]http://www.science.uva.nl/events/CLIN2005/Program/Abstracts/abstract-andogah.html

Table 2: Reference document index structure

| Field | Lucene Type | Description |
|-------|-------------|-------------|
| docid | Field.Keyword | Document unique identification/number |
| content | Field.Unstored | Combination of HEADLINE and |
| | | TEXT tag content |

document index (Lucene field content [Table 2]). Search results were combined through linear interpolation ( see Equation 1) for final relevance rank (retrieving top 1,000 relevant documents).

Our fourth run (CLCGGeoEE10) combine run 2 query result with result of querying footprint document index (Lucene field nm [Table 1]) for geo-relevant-terms extracted from topic TITLE-DESC-NARR. To combine the result of run 2 with result of querying footprint document index we used the linear interpolation (see Equation 1) [Leidner, 2005].

$$FinalScore = \lambda_T \, ThematicScore(d_R, q_T) + \lambda_G \, GeoScore(d_F, q_G) \qquad (1)$$
$$\lambda_T + \lambda_G = 1 \qquad (2)$$

where $d_R$ is reference document, $d_F$ the footprint document, $q_T$ the thematic (non-geographic plus geographic terms) query, $q_G$ the geographic (geo-relevant terms) query, $\lambda_T$ the thematic interpolation factor and $\lambda_G$ the geographic interpolation factor. For this experiment $\lambda_T = \lambda_G = 0.5$.

Our fifth run (CLCGGeoEE11) is similar to run four (CLCGGeoEE10) except that geo-revelant-terms extracted from topic TITLE-DESC-NARR were augmented with geo-references obtain from GKB. For example, topic *G033* geo-relevant-terms were augmented with the names of major cities/towns/places within Ruhr area of Germany – Bochum, Bottrop, Dortmund, Duisburg, Essen, Gelsenkirchen, Hagen, Hamm, Herne, Mlheim, Oberhausen, Recklinghausen, Ennepe-Ruhr, Unna, Wesel, Mlheim an der Ruhr, Mulheim an der Ruhr.

# 4 Evaluation and future work

## 4.1 Official GeoCLEF results

Tables 3 show the result of our official runs. Though both CLCGGeoEE1 and CLCGGeoEE5 use topic TITLE-DESC (querying different document collection content), CLCGGeoEE5 performed better. CLCGGeoEE5, CLCGGeoEE10 & CLCGGeoEE11 schemes use linear interpolation (with $\lambda_T$ and $\lambda_G$ set to 0.5) to combine result of query against reference document collection and document collection footprint indexes. We note that CLCGGeoEE10 performed poorly while CLCGGeoEE11 performed better.

Table 3: Individual Run Performance as measured by Mean Average Precision and R-Precision

| | CLCGGeoEE1 | CLCGGeoEE2 | CLCGGeoEE5 | CLCGGeoEE10 | CLCGGeoEE11 |
|---|---|---|---|---|---|
| MAP | 0.1730 | 0.2163 | 0.1757 | 0.1690 | **0.2194** |
| R-Precision | 0.1983 | **0.2194** | 0.1777 | 0.1762 | 0.2144 |

Several factors might have influenced the performance:

- predominance of geographic concepts and spatial-relationship qualifiers such as country, city, southern, west, etc. both in the query and document footprints at expense of place names, and thereby shifting query result in wrong direction propagating irrelevant documents to

the top; among the top 20 geo-relevant-terms – `country, city, street, west, north, south, east, southern, district, town`.

- value of 0.5 assigned to $\lambda_T$ and $\lambda_G$ in linear interpolation [Equation (1)] above might have tilted result by asigning higher scores to documents retrieved from reference document collection or vice versa, and thereby propagating irrelevant documents to the top in the final rank

- not all documents were indexed as our adopted geographic named entity tagger (Alias-i Lingpipe) reported content error for certain files while processing reference collection files. Why? Missing `<?xml version = "1.0" encoding = "UTF-8"?>` in files: `la071794.xml` (LA94), and `9511[16-18,23-30].xml` and `9512[01-30].xml` (GH94). As a result 51,525 Glasgow Heralds documents were indexed out of 56,472 and 112,552 LA Times documents were indexed out of 113,005. This might have imparted query results negatively as 5,400 documents (which might have contained relevant documents) were left out.

## 4.2 Post GeoCLEF results

After the release of official GeoCLEF 2006 results a prelimenary experiment was done to confirm performance influencing factors mentioned in the previous sub-section. The results of the experiment are shown in Table 4. Scheme 3 performed better across all runs.

Table 4: Individual run perfronacne as measured by Mean Average Precision and R-Precision

| Run | $\lambda_G = \lambda_T = 0.5$ Scheme 1[a] | | $\lambda_G = \lambda_T = 0.5$ Scheme 2[b] | | $\lambda_G = 0.35, \lambda_T = 0.65$ Scheme 3[c] | |
|---|---|---|---|---|---|---|
| | MAP | R-Prec | MAP | R-Prec | MAP | R-Prec |
| CLCGGeoEE1 | 0.1743 | 0.2019 | 0.1966 | 0.2302 | 0.1966 | 0.2302 |
| CLCGGeoEE2 | 0.2165 | 0.2123 | 0.2386 | 0.2312 | 0.2386 | **0.2312** |
| CLCGGeoEE5 | 0.1822 | 0.1883 | 0.1976 | 0.1919 | 0.2102 | 0.2109 |
| CLCGGeoEE10 | 0.1663 | 0.1772 | 0.2069 | 0.1919 | 0.2388 | 0.2256 |
| CLCGGeoEE11 | **0.2170** | **0.2192** | **0.2397** | **0.2369** | **0.2390** | 0.2307 |
| CLCGGeoEE8[d] | 0.1401 | 0.1334 | 0.1135 | 0.1227 | 0.1135 | 0.1227 |

[a]All documents processed [GH94 & LAT95], geo-query against geo-relevant-terms [place name, geo-concept, geo-relationship, geo-adjective] (footprint index)

[b]All documents processed, geo-query against place names (footprint index)

[c]All documents processed, geo-query against place names (footprint index)

[d]Query expansion on TITLE-DESC-NARR, query footprint document index only

## 4.3 Future work

The results of our submitted runs raised several pertinent questions for future investigation:

- extend to which geographic aspect of document influence GIR result: (1) querying topic geographic aspect against reference document collection, (2) querying topic non-geographic aspect against reference document

- an appropriate value for $\lambda$ in linear interpolation [Equation (1)] above for GIR

- an appropriate document collection footprint indexing strategy

- improve geographic named entity recognition, classification and real world resolution

- geographic query expansion strategies – blind feedback, addition of place names, expansion through hierarchical information mining, distance based query expansion.

# 5    Concluding remarks

We employed a strategy of separately indexing document footprint along side index of reference document, and combine query results of the two indexes through linear interpolation. Our approach yielded an average result as compared to overall GeoCLEF 2006 result on monolingual English task. A number of pertinent questions were raised for future investigation which we hope to address and integrate in our system. Analysis of individual topic performance to give further insight in our approach is under way.

# 6    Acknowledgements

# References

Gey et al [2005] Fredric Gey, Ray Larson, Mark Sanderson, Hideo Joho, Paul Clough and Vivien. *GeoCLEF: the CLEF 2005 Cross-Language Geographic Information Retrieval Track Overview.* At GeoCLEF 2005 in CLEF 2005 Workshop, 21 - 23 September 2005, Vienna, Austria.

Gey and Vivien [2005] Fredric Gey and Vivien Petras. *Berkeley2 at GeoCLEF: Cross-Language Geographic Information Retrieval of German and English Documents.* At GeoCLEF 2005 in CLEF 2005 Workshop, 21 - 23 September 2005, Vienna, Austria.

Larson [2005] Ray R. Larson. *Cheshire II at GeoCLEF: Fusion and Query Expansion for GIR.* At GeoCLEF 2005 in CLEF 2005 Workshop, 21 - 23 September 2005, Vienna, Austria.

Andogah [2005] Geoffrey Andogah. *Is Groningen referring to the city or the province? Geographic Named Entity Disambiguation.* Presentation at CLIN 2005, The 16th Meeting of Computational Linguistics in the Netherlands Amsterdam, December 16, 2005.

Leidner [2005] Jochen L. Leidner. *Preliminary Experiments with Geo-Filtering Predicates for Geographic IR.* At GeoCLEF 2005 in CLEF 2005 Workshop, 21 - 23 September 2005, Vienna, Austria.

Ferres et al [2005] Daniel Ferres, Alicia Ageno, and Horacio Rodriguez. *The GeoTALP-IR System at GeoCLEF-2005: Experiments Using a QA-based IR System, Linguistic Analysis, and a Geographical Thesaurus.* At GeoCLEF 2005 in CLEF 2005 Workshop, 21 - 23 September 2005, Vienna, Austria.

Hughes [2005] Baden Hughes. *NICTA i2d2 at GeoCLEF 2005.* At GeoCLEF 2005 in CLEF 2005 Workshop, 21 - 23 September 2005, Vienna, Austria.

Buscaldi et al [2005] Davide Buscaldi, Paolo Rosso, Emilio Sanchis Arnal. *A WordNet-based Query Expansion method for Geographical Information Retrieval.* At GeoCLEF 2005 in CLEF 2005 Workshop, 21 - 23 September 2005, Vienna, Austria.