

Building a Corpus of Spatial Relational Expressions Extracted from Web Documents

Jan Oliver Wallgrün
Human Factors in GIScience
Lab
Pennsylvania State University
University Park, PA 16802,
USA
wallgrun@psu.edu

Alexander Klippel
Department of Geography,
GeoVISTA Center
Pennsylvania State University
University Park, PA 16802,
USA
klippel@psu.edu

Timothy Baldwin
Department of Computing and
Information Systems
The University of Melbourne
Victoria 3010, Australia
tb@ldwin.net

ABSTRACT

Spatial language, despite decades of research, still poses substantial challenges for automated systems, for instance in geographic information retrieval or human-robot interaction. We describe an approach to building a corpus of natural language expressions extracted from web documents for analyzing and modeling spatial relational expressions (SRE). The unique characteristic of this corpus is that it is built around georeferenced triplets, with each triplet containing two entities (including their latitude/longitude coordinates) related by a spatial expression such as *near*. While the approach is still experimental, our first results are promising, in that we believe they will form the foundation for a comprehensive contextualized model for interpreting spatial natural language expressions. For the time being, we are focusing on a single domain, hotel reviews. This domain restriction allowed us to implement a proof-of-concept that this approach, with advances in natural language technologies, will indeed deliver a comprehensive corpus. The potential to collect larger corpora, and associated challenges, is discussed.

Categories and Subject Descriptors

H.3.3 [Information Storage and Retrieval]: Information Search and Retrieval

General Terms

Experimentation, Human Factors, Languages

Keywords

spatial relations, proximity, information retrieval, corpus building, spatial language

1. INTRODUCTION

In human language, spatial relationships between objects are usually specified as spatial relational expressions (SRE); Permission to make digital or hard copies of all or part of this work for personal or classroom use is granted without fee provided that copies are not made or distributed for profit or commercial advantage and that copies bear this notice and the full citation on the first page. Copyrights for components of this work owned by others than ACM must be honored. Abstracting with credit is permitted. To copy otherwise, to republish, to post on servers or to redistribute to lists, requires prior specific permission and/or a fee. Request permissions from Permissions@acm.org.
SIGSPATIAL'14, November 04-07 2014, Dallas/Fort Worth, TX, USA
Copyright 2014 ACM 978-1-4503-3135-7/14/11...\$15.00
<http://dx.doi.org/10.1145/2675354.2675702>

prepositions, adverbial phrases, surrogate terms), such as *at*, *near*, *to the left of*, *to the south of*. In addition to the main spatial aspect underlying the respective relation (e.g. proximity of the involved objects in the case of *near*), it has been demonstrated that the meaning of spatial relations is strongly influenced by various kinds of contextual information (e.g., object type/dimension, domain) [4, 10, 2, 12].

Understanding and formally modeling the semantics of SRE, like the ones listed above, is critical for a wide variety of applications, including information retrieval [23, 6], human-robot interaction [20], and automatic text and speech processing and generation [17, 1]. As a step in this direction, the work in this paper is directed toward the long-term goal of constructing and analyzing a large corpus of human-produced expressions containing SREs alongside georeferenced locations of the involved objects, and the context in which these expressions have been produced.

The great utility of the Web in conducting studies on human language use and related topics has been widely recognized [18, 24, 26]. However, there are many challenges to achieving the goal of collecting spatial expressions, together with geometric information about the involved objects as well as other contextual information. Objects in spatial relationships within a given expression need to be unambiguously identified, a challenging task for the following reasons.

- Names of objects are often not unique (*Springfield is close to Hartford*).
- Names are often shortened, otherwise modified, or described using complex expressions (*New York* for *New York City*, or *the southern part of Pennsylvania*).
- Vague pronouns are used to represent objects, and thus a fair amount of deduction is required to determine the pronominal antecedents (a complex topic on its own).
- Text, in particular on the Web, often does not follow simple patterns like `<object> <relation> <object>` triples (e.g., *Hotel Sunshine is close to Central Park*), but uses complex syntactic structures, including fragments, or spans multiple sentences (e.g., *Hotel Sunshine - great hotel! Very close to Central Park btw.*).

Once the objects involved have been unambiguously identified, geometric information about their spatial extension and location needs to be determined, a process termed *georeferencing* or *geocoding*. Geographic gazetteers or worldwide spatial data collections such as Geonames and Open-

StreetMap can help in this step, but due to naming ambiguity, it is not easy to find the correct objects in these datasets. Furthermore, these sources typically contain only certain types of objects, are otherwise incomplete, are restricted to point information (for example, latitude and longitude coordinates in the WGS1984 coordinate system), or are difficult to handle because of their sheer size. Places related by spatial prepositions that are vaguely described (*the center of Pittsburgh*) cannot be looked up in any existing dataset.

Extracting additional contextual information is also a challenging task. Certain kinds of contextual information in which one might be interested, for instance, the political bias of the producer, are even impossible to obtain. Other factors that complicate the construction of a suitable corpus include the problem of identifying when a preposition is used to denote an actual spatial relation and the fact that spatial prepositions are often heavily modified, such as by linguistic hedges (e.g. *very / quite / somewhat / pretty close*) and negation (*not even close*).

As a result of these challenges, this kind of corpus-building tends to require sophisticated natural language processing and interpretation capabilities. There have been only a few attempts to analyze and model the meaning of spatial relations, and the resultant impact of different kinds of contextual information, based on information garnered from the web in the literature [23, 6] (see Section 2 for details). In the remainder of this paper, we describe a general approach that we designed to build a corpus consisting of geo-referenced triplets (GRT; see also [16]) which contain a located object (**lo**), spatial relation (**rel**), and reference object (**ro**) (e.g., *Baltimore is located close to Washington D.C.*). In our description of the approach, we cover a range of problems, from the generation of suitable objects, to the search for and extraction of SREs from web pages, to the identification of important contextual information (including the actual geometry of the involved objects). We also report on results from a first case study performed within this framework that focused on proximity relations (*near, close, and next to*). The corpus built in this case study consists of expressions in which hotels are set into a proximity relation with an attraction, or special building or other entity in the same city. Results are collected from typical hotel booking web portals, which provide reviews of the respective hotels. As we will discuss, restrictions made in this case study allow us to bring the mentioned challenges down to a manageable level. Results from a first analysis of the collected data and insights gained are presented and discussed.

The paper is structured as follows: Section 2 discusses work from the literature that is related to our approach. In Section 3, we present the corpus collection framework and our design decisions. In Section 4, we present the mentioned case study and data analysis. In Section 5 we close with conclusions and an outlook on future research activities.

2. RELATED WORK

Substantial efforts have been made to understand and model factors that change the interpretation of proximity terms, both from the perspective of experimental studies and from the perspective of building formal models. In comparison to controlled experiments, which naturally address rather specific questions, the use of crowd science to model and understand spatial language is rather sparse. Hence, we

only briefly summarize some of the main results obtained through behavioral studies and assumptions made in formal models.

Distance is the most obvious choice for understanding linguistically expressed proximity, e.g., *near* [10, 14, 9]. While metric distances are relatively easy to obtain and certainly do play a role, distance by itself cannot explain how the same SRE (e.g., *near*) is applicable, for example, to both the bowl near the teapot as well as the moon near earth. Additional factors (either theoretically or in controlled experiments) that have been identified as contextually relevant for the interpretation of proximity terms are: The size of the entities themselves as well as their relative size [10, 3]; the same entities placed in different contexts [10, 5, 19]; whether or not an entity is movable [25]; semantic context [21]; the presence of other entities and, for example, relative distance between them [10, 13]; mode of transportation [28]; and familiarity [28].

To account for factors other than distance, researchers have developed both theoretical and formal frameworks that address contextual effects in the interpretation of SREs. Among the most prominent theoretical frameworks is the *extrageometric framework*, also referred to as the *functional geometric framework* [4]. Formal frameworks include fuzzy set theories [7, 22], the definition of impact areas of entities based on a specific widths to the boundaries of common convex shapes [15], the calculation of a size/distance ratios [11], Voronoi diagrams [8], and contextualized proximity measures [2], to name just a few.

The use of passive crowdsourcing (e.g., by crawling web documents) for replacing controlled behavioral experiments or theoretical considerations is a rather recent phenomenon and at a very early stage [6, 27, 22]. Our case study on proximity relations reported in Section 4 uses a scenario that is similar to the approach reported in [22] but puts a stronger emphasis on the extraction of context information.

3. CORPUS COLLECTION APPROACH

The goal of the general approach presented in this section is to build a corpus of spatial natural language expressions, extracted from the web, in which two objects, the located object (**lo**) and the reference object (**ro**), are related by a SRE, using a spatial relation (**rel**) in the form of a geo-referenced triplet (GRT) **lo** * **rel** * **ro** where each * represents a small number (≤ 8 in the study reported in Section 4) of arbitrary filler words (cf. Figure 2 for an example).

Our framework consists of four main modules (see Figure 1). Collected web pages and intermediate results are stored in the central corpus repository. The approach allows several instances of the modules to be run in parallel, while avoiding any duplicate processing or downloading of the same web document. The approach makes use of several query tools and web resources that build the reference object list, search for occurrences of the desired pattern, geocode addresses, and determine contextual information. Next, we will briefly discuss the different modules.

3.1 Reference Object Generator

The first component of the corpus building approach is responsible for automatically constructing lists of objects to be used as located or reference entities in a GRT (for simplicity referred to as *reference object lists* from now on), by drawing from suitable web-based geographical data sources. Our cur-

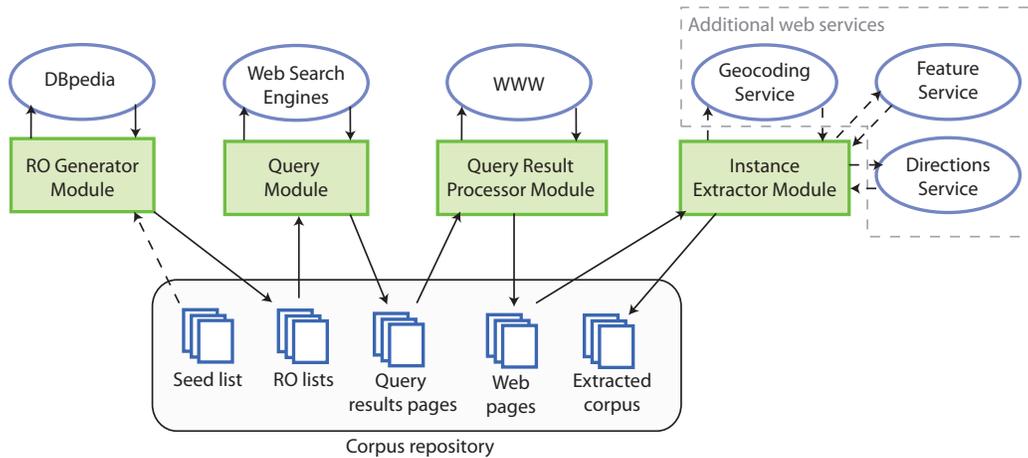


Figure 1: Illustration of the general corpus building approach.

rent object lists consist of monuments, tourist attractions, or other special buildings or entities extracted from DBpedia¹ and selected based on their type. We extract the names and WGS1984 coordinates of the selected objects, and store them in the corpus repository.

3.2 Query Module

This module operates on the reference object lists and uses existing web search engines to query a pattern of **lo**, **rel**, and **ro** (**lo * rel * ro**), where **rel** is replaced by one of the hard-coded spatial relations, and **lo** and **ro** are replaced by names of objects from the reference object lists in the repository, or by other keywords (see Section 4 for more details). Exact-phrase queries that include wildcards yield results that contain the three key elements **lo**, **rel**, and **ro**, in that order, but potentially with fill words in between. Obtained result web pages up to a specified page limit are stored in the corpus repository.

3.3 Query Result Processor

The Query Result Processor downloads the web pages listed on the result pages collected in the previous step. This module is currently set to download and store all page results, although in our case study (Section 4) we processed only a small fraction of them. It is our expectation that progress on natural language processing and geoparsing will make more of these results usable in the future.

3.4 Instance Extractor

The Instance Extractor has the task of processing the web documents in which a search engine has determined that the triplet pattern (**lo * rel * ro**) appears. This web page processing component has several submodules that are executed in sequential order.

3.4.1 Extracting the Phrases

The first step is to extract the natural language expressions that match our "**lo * rel * ro**" pattern, using html tag removal and regular expressions. To keep a bit more context surrounding the matching phrase, we store the character sequence starting m characters before the **lo** and ending m characters after the **ro**, so for example:

¹<http://DBpedia.org>

and felt so liberated! That was in Baltimore, near Washington, DC, which doesn't get as cold as

3.4.2 Extracting Disambiguation Information

Any additional information needed to unambiguously identify the objects in the expression, or to verify that these are indeed the reference objects used in our query, is extracted.

3.4.3 Geocoding and Derivation of Main Spatial Information

If both objects are reference objects used to formulate the queries, their WGS 1984 coordinates are already known. If not, the related objects have to be geocoded using a special geocoding service. When dealing with proximity information, it is important to know the distance between the two involved objects. Therefore, the line-of-sight distance is computed from the coordinates of both the **lo** and the **ro**. In the case of spatial relations other than distance (e.g. direction, shape, size), the specific nature of these relations needs to be derived which is typically more involved than in the case of distance.

3.4.4 Extraction of Other Context Information

In addition to geometric information (location, spatial extension) about the involved object, it is desirable to extract other information that could have an impact on when and how humans employ SREs, such that the nature of this impact can be analyzed and modeled. Additional web-based (or knowledge-based) services may have to be queried. For instance, in the case of proximity relations, we are not only interested in line-of-sight distance but also travel distance or time for different modes of travel (driving, walking, etc.).

3.4.5 Storing the New Instance

As a last step, the Instance Extractor module stores a tuple for the new instance in the corpus database with the following information:

- the URL of the web page and the name under which it is stored in the repository;
- the extracted phrase;
- **lo** information (name, coordinates/geometry, context information);

- the spatial relation;
- **ro** information (name, coordinates/geometry, context information);
- additional spatial information, e.g. line-of-sight distance;
- optionally: additional context information.

In the next section, we put this approach into practice in a case study about the usage of proximity relations.

4. CASE STUDY: PROXIMITY RELATIONS

The goal of the case study presented in this section was to create a first corpus of spatial natural language expressions extracted from the web, in which one of the three proximity relations *close*, *near*, and *next to* is used to spatially relate a located object (**lo**) to a reference object (**ro**) in a triplet of the form **lo** * **rel** * **ro** where, as mentioned, each * represents a small number of filler words. While *Baltimore is located close to Washington D.C.* would be an example of such an expression, we make the additional restriction that the **lo** is always the word *hotel*, and the **ro** is the name of an object from the automatically generated list of reference objects which are all attractions, special buildings or other spatial entities with (at least locally) unique names from different cities around the world. The expressions are extracted from hotel booking and reviewing web pages (such as TripAdvisor² and Expedia³) that contain reviews of a particular hotel together with address information. Because the reference objects are already identified and have known WGS 1984 coordinates, and because the hotels appearing as **lo** come with address information which is relatively easy to geocode, the object identification and geocoding problems are greatly simplified. Given the reference object coordinates, it is also straightforward to derive the line-of-sight distance between the two related objects. Figure 2 shows an example fitting the pattern just described. For each found example, we collect and store a tuple consisting of the following information:

- URL of the web page that contains the expression
- all characters starting 20 characters before the **lo** and ending 20 characters after the **ro**
- name, address, and coordinates of the **lo** (hotel)
- name, coordinates, and polygonal geometry of the **ro** (attraction, special building)
- line-of-sight distance between **lo** and **ro** computed from the two point locations
- travel distances and times between the two objects for driving and walking, respectively.

Our explanation below roughly follows the modules of our framework described in Section 3.

²<http://www.tripadvisor.com>

³<http://www.expedia.com>

4.1 Reference Object List Generation

The list of objects for the **ro** are created automatically in a two-step process based on linked data from DBpedia.org. First a list of 882 cities is generated as a seed list. Second, for each of the cities, a list of suitable attractions and special buildings is constructed and stored in the corpus repository. To generate the seed list of cities, the Reference Object Generator module uses a SPARQL query for DBpedia resources of cities (class `dbpedia-owl:city`) in the world with a population larger than 500,000. It then reads the resource for each of the cities and follows the `dbpedia-owl:location_of` relations to the DBpedia resources of entities located in that city. Based on the entity type, suitable objects are selected, their names and WGS1984 coordinates extracted, and the results added to the reference object lists stored as data sets (one for each city) in the corpus repository. For instance, the list created for Paris contains the following entry for the Eiffel Tower:

```
entity name: Eiffel Tower
city: Paris
country: France
latitude: 48.85815
longitude: 2.2945
```

4.2 Querying and Result Page Processing

The Query module searches for query phrases such as *hotel * close * Eiffel Tower*. Result pages are stored and processed by the Query Result Processor module, and the listed web pages are downloaded.

4.3 Extracting the Hotel Name and Address

The Instance Extractor then extracts phrases like the following: *... only recommend it. Hotel Jardins d'Eiffel is located very close to the Eiffel Tower. The rooms are very...* from the downloaded web pages. Address extractors are written specifically for the individual types of web pages (e.g. a user review web page from TripAdvisor). Each address extractor returns the specific address information for a hotel on a particular page, i.e., street address, zip code, city, and country. We make the implicit assumption that the hotel referred to in the **lo** position is the hotel that the page is about. While this could lead to errors, we include some sanity checks to minimize the number of erroneously identified located objects. An example output of this module is:

```
name: Hotel Jardins d'Eiffel
street address: 8 rue Amelie
zip code: 07 Arr., 75007
city: Paris
country: France
```

4.4 Geocoding and Distance Calculation

We use Google's geocoder service⁴ to determine the WGS 1984 coordinates of the hotel, based on the address extracted in the previous step. Returned are the hotel's latitude and longitude in decimal degrees:

```
lo latitude: 48.8591309
lo longitude: 2.3070443
```

At this point, we have the coordinates of both the **lo** and the **ro**. Figure 3 shows the **lo** and **ro** from our running

⁴<https://developers.google.com/maps/documentation/geocoding/>

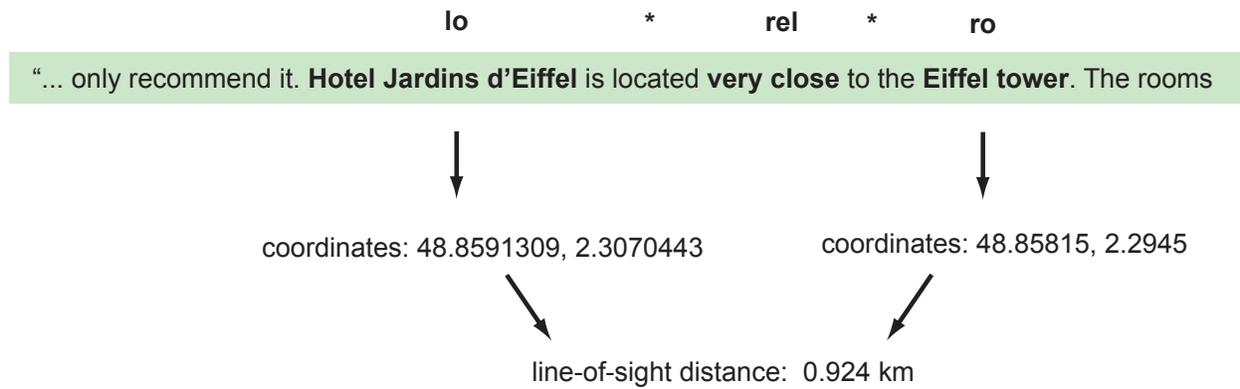


Figure 2: Example of a lo * rel * ro triple with coordinates and derived line-of-sight distance. Such instances (with additional information) make up the corpus built in this case study.

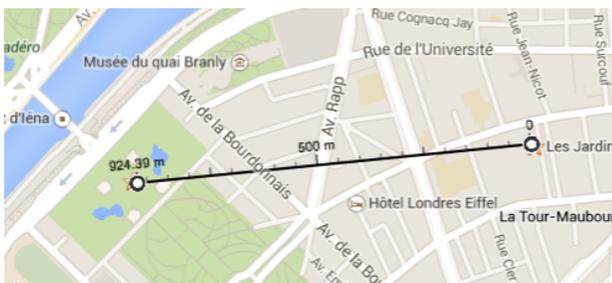


Figure 3: ro and lo from the example shown on Google Maps. The line-of-sight distance is simple to compute from the coordinates of both entities [source: Google Maps (<https://www.google.com/maps>)].



Figure 4: Polygonal geometry obtained from OSM for the Eiffel Tower reference object [source: OpenStreetMaps (<http://www.openstreetmap.org>)].

example on Google Maps. Given the coordinates, the line-of-sight distance between the two entities is computed resulting in:

line-of-sight distance: 0.924 km

4.5 Geometry Derivation

To study how it affects human use of SREs, one piece of spatial contextual information we collect is a geometric description of the boundary of the ro. In most cases, DBpedia.org provides a point geometry representing the WGS 1984 coordinates of the object's centroid. To obtain a polygonal representation of the object's boundary, our system uses a different source, namely OpenStreetsMap⁵ (OSM). We access the OSM data via the Nominatim⁶ search interface, querying the name of the reference object which in turn yields a list of candidate entities. The best match from that list is determined by a sanity check on whether the centroid coordinates of the ro from DBpedia.org are contained in the bounding box of the respective OSM entity, and are very close to this entity's centroid. The best match that satisfies our criteria is assumed to be the ro, and we maintain its polygonal geometry. Figure 4 shows the geometry obtained for our example with the Eiffel Tower as the ro.

⁵<http://www.openstreetmap.org/>
⁶<http://nominatim.openstreetmap.org/>

4.6 Travel Distance / Time Derivation

Last, we obtain the travel distances and times between the lo and ro to do a comparison between their travel distance and line-of-sight distance, employing the Google Directions API web service.⁷ The service allows us to query directions between the point coordinates of lo and ro. The provided route directions include both the distance and time for the route. Moreover, Google's service returns results for different modes of travel. Our system currently collects travel distances and times for driving and walking. For our example (see Figure 5), the travel distances are:

driving distance: 2.191 km
 driving time: 352 s
 walking distance: 1.102 km
 walking time: 830 s

4.7 Storing the New Instance

As the final step, the new instance is stored in the corpus. The tuple consists of the url and name of the web page, the extracted before-and-after text, lo name, lo address, lo coordinates, ro name, ro coordinates, ro boundary polygon, line-of-sight distance, travel distances, and travel times for both driving and walking.

⁷<https://developers.google.com/maps/documentation/directions/>



Figure 5: Driving (top) and walking (bottom) routes for Hotel Jardins d’Eiffel and the Eiffel Tower [source: Google Maps (<https://www.google.com/maps>)].

5. CORPUS ANALYSIS

By using our general approach in this way, we collected (so far) about 330 instances. While our approach contains measures to prevent errors, it can—for the reasons we discussed in the introduction—not be expected that all instances indeed contain a SRE and that the involved objects (hotels) with their coordinates are always identified correctly. Indeed, manually going through the instances and checking the stored information for plausibility was required and led to the removal of 50 instances. The main reasons were:

- relation occurs with negation (e.g., *is not close to*); while it will be interesting to include negations in the analysis, we currently do not deal with this case;
- *near*, *close*, or *next to* are used in a non-spatial sense or as a spatial adjective (e.g., *near perfection*, *in close proximity*);
- the reference object in the phrase is not the one from our reference list; this happens when the actual name appears as part of longer name; for instance in one case *Notre Dame Metro Station* occurs in the phrase, while the actual reference object on our list is *Notre Dame*;
- the word *hotel* in the **lo** position clearly does not refer to the hotel the review is about (e.g., *There is another hotel right next to Central Station*).

While it seems possible to identify and filter out some of these cases automatically, advanced natural language processing techniques will be required to reach a level where such errors become negligible. Therefore, we expect that manual inspection will be required for the foreseeable future. Of the remaining 280 instances, 170 involved the relation *near*, 82 the relation *close*, and 28 the relation *next to*. In addition to determining travel distances and times for both walking and driving for these instances, our current approach of deriving polygonal geometries for the reference objects was able to do so for 66.1% of the instances. While the results are clearly preliminary and more instances will be needed, in particular to investigate the effects of different kinds of context information, we will in the following show some preliminary results from analyzing the data.

5.1 Histograms and Basic Statistics

We created histograms over the line-of-sight distances between the involved objects. Figure 6 shows the histograms for *near*, *close*, and *next to*. While *near* and *close* have very similar mean values (1.562km for *near* and 1.560km for *close*) compared to *next to* (1.046km), *near* has a smaller standard deviation (2.585km) compared to *close* (3.434km). The standard deviation of *next to* is 1.215.

5.2 Line-of-sight vs. Travel Distance

In addition to generating the basic statistics for the three relations based on line-of-sight distance, we also did so for the travel distances and times based on both driving and walking, as discussed in Section 4.6. Table 1 summarizes the results by showing mean values and standard deviations for all combinations. We see a consistent increase in the mean values from *next to* to *close* to *near* in all cases and for the most part, *next to* has a much smaller standard deviation with a notable exception being the driving times. *Near* and *close*, in general, have very similar means and standard deviations. Given the relatively small amount of instances, we leave the question of which of these properties has the larger influence on or is best suited to modeling the human usage of the three proximity relations for future analyses.

5.3 Size of the Reference Object

We also looked more closely at the instances for all three relations that have rather large distances (see again Figure 6). We noticed quickly that these are typically reference objects with a significant size such as airports, stadiums, or parks. It is clear that in these cases the centroid alone is not a good approximation. As a next step, we therefore plan to employ the collected actual polygonal geometries to compute the line-of-sight distance not between the centroids but between the centroid of the hotel and the closest point of the reference object, and compare the results to what we are currently getting. While this demonstrates that the actual spatial extension of the involved object is an important piece of information in this kind of analysis, there is other context information that may play a role such as the location of entry/exit points which may be more difficult to collect.

6. CONCLUSIONS & OUTLOOK

While we are currently only using a small fraction of the collected web pages (those with hotel reviews), our database includes 106,000 documents collected over two weeks in a process that can still be sped up drastically. This indicates a huge potential for scaling up the corpus. However, as noted in the introduction, such a large-scale analysis will require advanced language processing, interpretation, and geoparsing techniques. Most likely, we envision an incremental process in which the restrictions of the current study will be dropped step by step, yielding more and more results in wider realms. Once enough instances have been collected, the analysis will focus on understanding the impacts of different contextual effects, and modeling these and the relations themselves using a probabilistic approach (cmp. [23]).

7. ACKNOWLEDGMENTS

We thank the reviewers for valuable feedback and suggestions, and Elaine Guidero for proofreading the manuscript.

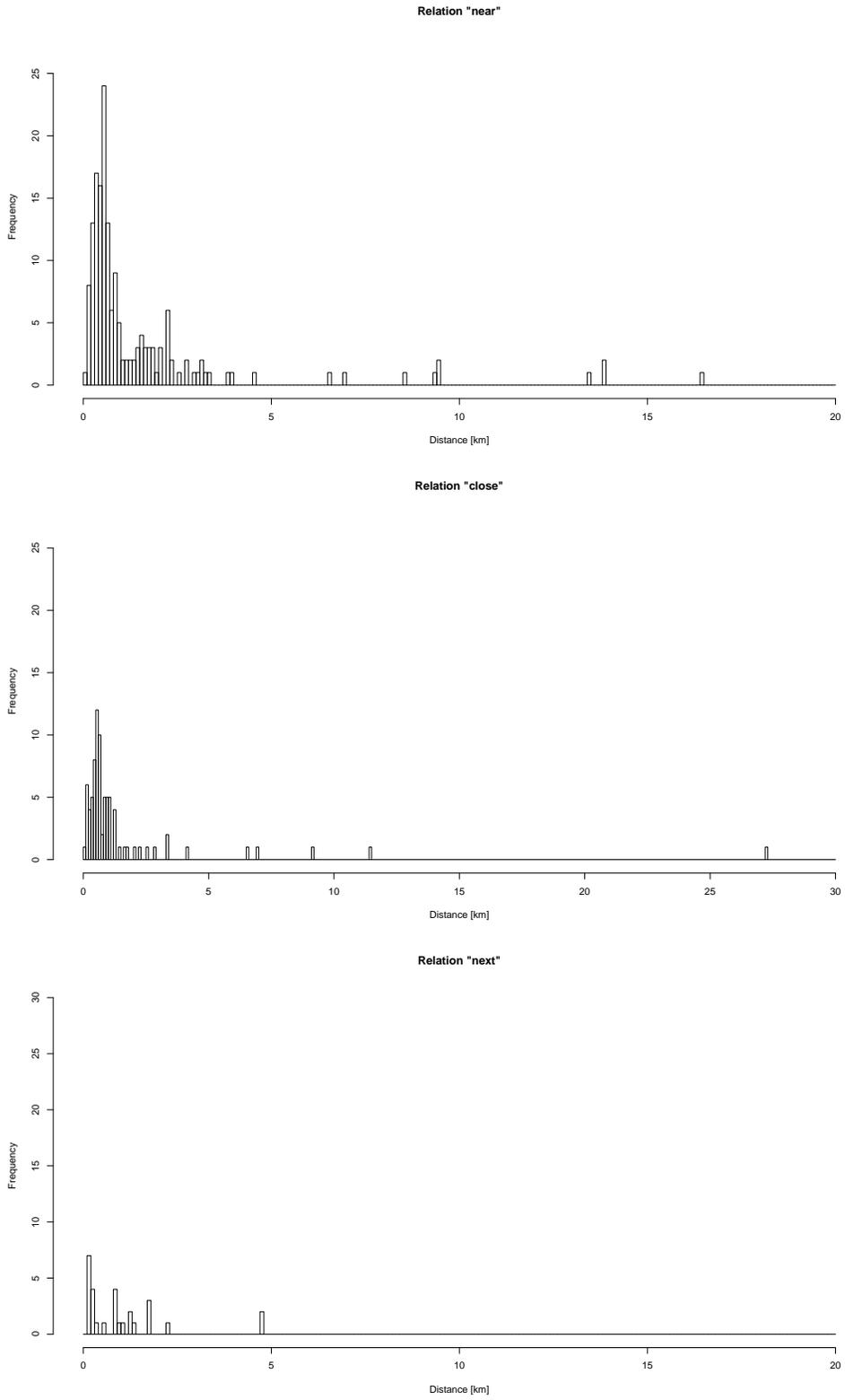


Figure 6: Histograms of instances over distance.

	line-of-sight distance [km]	travel distance [km] (driving)	travel time [s] (driving)	travel distance [km] (walking)	travel time [s] (walking)
near	m=1.562, sd=2.585	m=2.739, sd=4.385	m=362.8, sd=752.3	m=1.985, sd=3.114	m=1467.6, sd=2274.9
close	m=1.560, sd=3.434	m=2.415, sd=3.958	m=313.0, sd=260.4	m=1.929, sd=3.841	m=1440.2, sd=2835.2
next to	m=1.046, sd=1.215	m=1.946, sd=1.889	m=310.4, sd=276.1	m=1.411, sd=1.640	m=1043.9, sd=1204.0

Table 1: Mean values (m) and standard deviations (sd) over the collected instances *near*, *close*, and *next to*.

8. REFERENCES

- [1] J. A. Bateman. Language and space: a two-level semantic approach based on principles of ontological engineering. *International Journal of Speech Technology*, 13:29–48, 2010.
- [2] J. Brennan and E. Martin. Spatial proximity is more than just a distance measure. *International Journal of Human-Computer Studies*, 70(1):88–106, 2012.
- [3] M. Burigo and K. Coventry. Context Affects Scale Selection for Proximity Terms. *Spatial Cognition & Computation*, 10(4):292–312, 2010.
- [4] K. R. Coventry and S. Garrod. *Saying, Seeing, and Acting: The Psychological Semantics of Spatial Prepositions*. Psychology Press, Hove, 2004.
- [5] K. R. Coventry and S. Garrod. Towards a classification of extra-geometric influences on the comprehension of spatial prepositions. In L. A. Carlson and E. van der Zee, editors, *Functional features in language and space*. Oxford UP, 2004.
- [6] C. Derungs and R. Purves. Where’s near? Using web n-grams to explore spatial relations. In *GIScience 2014*, to appear.
- [7] S. Dutta. Qualitative spatial reasoning: A semi-quantitative approach using fuzzy logic. In A. Buchmann, O. Günther, T. R. Smith, and Y. F. Wang, editors, *Design and implementation of large spatial databases*, pages 345–364. Springer, 1990.
- [8] G. Edwards and B. Moulin. Toward the simulation of spatial mental images using the voronoi model. In P. Olivier and K.-P. Gapp, editors, *Representation and processing of spatial expressions*, pages 163–184. Lawrence Erlbaum, 1998.
- [9] A. U. Frank. Qualitative spatial reasoning about distances and directions in geographic space. *Journal of Visual Languages and Computing*, 3:343–371, 1992.
- [10] M. Gahegan. Proximity operators for qualitative spatial reasoning. In A. U. Frank and W. Kuhn, editors, *Spatial Information Theory*. 1995.
- [11] K.-P. Gapp. Object localization: Selection of optimal reference objects. In A. U. Frank and W. Kuhn, editors, *Spatial Information Theory*. Springer, 1995.
- [12] M. M. Hall and C. B. Jones. Cultural and language influences on the interpretation of spatial prepositions. In *Proceedings of the GI Forum 2012*, 2012.
- [13] A. Herskovits. *Language and Spatial Cognition: An Interdisciplinary Study of the Representation of the Prepositions in English*. Cambridge University Press, Cambridge and England, 1986.
- [14] D. Kemmerer. ”Near” and ”far” in language and perception. *Cognition*, 73(1):35–63, 1999.
- [15] D. Kettani and B. Moulin. A spatial model based on the notions of spatial conceptual map and of object’s influence area. In C. Freksa and D. M. Mark, editors, *Spatial Information Theory. Cognitive and Computational Foundations of Geographic Information Science*, pages 401–416. Springer, Berlin, 1999.
- [16] A. Khan, M. Vasardani, and S. Winter. Extracting spatial information from place descriptions. In *Proc. of the First ACM SIGSPATIAL International Workshop on Computational Models of Place*, 2013.
- [17] P. Kordjamshidi, J. Hois, M. van Otterlo, and M.-F. Moens. Learning to interpret spatial natural language in terms of qualitative spatial relations. In T. Tenbrink, J. Wiener, and C. Claramunt, editors, *Series Explorations in Language and Space*. 2013.
- [18] M. Lapata and F. Keller. The web as a baseline: Evaluating the performance of unsupervised web-based models for a range of NLP tasks. In *Proceedings of the Human Language Technology conference of the North American Chapter of the Association for Computational Linguistics*, pages 121–128.
- [19] G. A. Miller and P. N. Johnson-Laird. *Language and perception*. Cambridge University Press, 1976.
- [20] R. Moratz, K. Fischer, and T. Tenbrink. Cognitive modelling of spatial reference for human-robot interaction. In *IEEE International Conference on Tools with Artificial Intelligence*, pages 222–228. 2000.
- [21] D. G. Morrow and H. H. Clark. Interpreting words in spatial descriptions. *Language and Cognitive Processes*, 3(4):275–291, 1988.
- [22] S. Schockaert, M. D. Cock, and E. Kerre. Location approximation for local search services using natural language hints. *International Journal of Geographical Information Science*, 22:315–336, 2008.
- [23] G. Skoumas, D. Pfoser, and T. K. Anastasios. On quantifying qualitative geospatial data: a probabilistic approach. In *Proceedings of the 2nd GEOGROWD Workshop*, pages 71–78, 2013.
- [24] K. Stock and C. Cialone. Universality, Language-Variability and Individuality: Defining Linguistic Building Blocks for Spatial Relations. In M. Egenhofer, N. Giudice, R. Moratz, and M. Worboys, editors, *Spatial Information Theory. 10th International Conference, COSIT 2011. Proceedings*, pages 391–412. Springer, 2011.
- [25] L. Talmy. Force dynamics in language and cognition. *Cognitive Science*, 12(1):49–100, 1988.
- [26] F. Twaroch, R. Purves, and C. Jones. Stability of qualitative spatial relations between vernacular regions mined from web data. In *Proceedings of Workshop on Geographic Information on the Internet*, 2009.
- [27] S. Xu and A. Klippel. Developing nearness models from geocoding spatial entities in a news corpus. In *GIScience 2012, extended abstracts*. 2012.
- [28] X. Yao and J.-C. Thill. How far is too far? – A statistical approach to context-contingent proximity modeling. *Transactions in GIS*, 9(2):157–178, 2005.