

1 PD-10: Natural Language Processing in  
2 GIScience Applications  
3

4 **Abstract**

5 Natural Language Processing (NLP) has experienced explosive growth in recent  
6 years. While the field has been around for decades, recent advances in NLP  
7 techniques as well as advanced computational resources have re-engaged aca-  
8 demics, industry, and the general public. The field of Geographic Information  
9 Science has played a small but important role in the growth of this domain.  
10 Combining NLP techniques with existing geographic methodologies and knowl-  
11 edge has contributed substantially to many geospatial applications currently in  
12 use today. In this entry, we provide an overview of current application areas  
13 for natural language processing in GIScience. We provide some examples and  
14 discuss some of the challenges in this area.

15 **Keywords**

16 natural language processing, text analytics, toponym disambiguation, topic  
17 modeling, question answering

18 **1 Definitions**

- 19
- Gazetteer: A dictionary or index of geographical names.
  - n-gram: A sequence of n tokens, where n is a number. N-grams typically  
20 range between 1 (uni-gram) and three (tri-gram).  
21
  - Token: The building blocks of natural language. Small units of text that  
22 (e.g., characters, words, combinations of words) that have been split from  
23 a larger document or corpus.  
24
  - Toponym: A place name. Often derived from a topographic feature.  
25

## 26 **2 Natural Language Processing and GIScience**

27 Natural Language Processing (NLP) is as an interdisciplinary research area that  
28 draws from the fields of linguistics, computational sciences, and many other re-  
29 lated disciplines including geography and geographic information science (GI-  
30 Science) that develop methods to analyze human language data. While the field  
31 includes a wide variety of topics it is primarily concerned with applying compu-  
32 tational techniques to analyze human language in a variety of forms. In recent  
33 years, the field has focused on the extraction of patterns and meaning from  
34 large volumes of natural language data such as text and speech audio. Today,  
35 the field is moving towards “understanding” concepts and themes presented in  
36 natural language with the goal of answering questions and informing decision  
37 making.

38 Historically, the domain of natural language processing has focused on the  
39 extraction of structured content from unstructured text. Early *Symbolic* NLP  
40 approaches involved interpreting text and speech through a series of user-defined  
41 rules. In the 1980s and 1990s various statistical inference techniques were de-  
42 vised for identifying and applying these rules to natural language. More recently,  
43 the domain has seen a shift towards the use of machine learning, including deep  
44 learning, *Neural*, methods. These recent approaches do not take a rule-based  
45 approach but rather aim to understand natural language through statistical  
46 methods which can identify linguistic properties of words, sentences, or docu-  
47 ments.

48 Though NLP does not fall solely within the discipline of Geography, a lot  
49 of human language is situated in geographic space and time and might make  
50 reference to inherently geospatial themes such as culture. Natural language  
51 varies by region meaning that GIScientists are well situated to process, identify,  
52 and contextualize patterns in language. Within the field of GIScience, NLP  
53 has been used to better understand a wide variety of geographic phenomena  
54 through identification of places, events, and activities as well as the extraction  
55 of linguistic patterns related to these entities. NLP techniques offer insight into  
56 geographic phenomenon that may not be accessible through traditional spatial  
57 and temporal analysis.

58 GIScientists are also able to leverage much of their existing expertise when  
59 processing natural language. Knowledge of spatial relationships, regional hi-  
60 erarchies and geographic laws & theories when combined with many leading  
61 NLP approaches result in cutting edge applications, many of which are actively  
62 used today. In the section to follow, a number of different NLP techniques are  
63 discussed with a specific focus on applications within the field of GIScience.  
64 The intent is to demonstrate how natural language processing is being used  
65 within GIScience applications today and discuss some of the challenges moving  
66 forward.

## 67 3 Applications of Natural Language Processing 68 in GIScience

69 A number of natural language processing applications exist within GIScience.  
70 This section summarizes a small, but key set of application areas that have  
71 emerged in recent years.

### 72 3.1 Toponym disambiguation

73 Important locations on the Earth are usually given labels or *toponyms* to allow  
74 them to serve in a common reference system. When someone makes a reference  
75 to Montréal, Canada, for example, there is shared understanding of where this  
76 place is located on the Earth as well as what type of place it is, namely a city.  
77 Toponym disambiguation is the process of (a) identifying Montréal as a location,  
78 and (b) differentiating it from any other location labeled as Montréal.

79 To discuss toponym disambiguation in more detail, we must first take a  
80 large step back and discuss some of the building blocks necessary for many  
81 natural language processing tasks. The first step involves deconstructing natural  
82 language to a format that enables computational analysis, through a method  
83 known as tokenization. Tokenization is the process of breaking down natural  
84 language into smaller lexical units which are referred to as *tokens*. Depending  
85 on the task, these units range from individual characters, to words (or sequences  
86 of words known as n-grams), sentences, paragraphs, or documents. The process  
87 of tokenization is easier for some languages than others. For instance, romance  
88 languages often delimit words with spaces whereas some Asian languages, such  
89 as Chinese, do not mark word boundaries with space delimiters making the  
90 process more complex [28].

91 In many languages, people use different inflection forms of words. For in-  
92 stance, *democratic*, *democracy*, *democracies*, and *democratization* all reference  
93 similar concepts, but for grammatical reasons the different words exist. For  
94 many applications these different concept references can be considered the same,  
95 thus it is advantageous to reduce them to a single token. Stemming is a simple  
96 solution to this problem that typically involves dropping the end of words such  
97 as derivational affixes, to reduce them to only those characters that the words  
98 have in common. For instance, a stemming approach to the above terms might  
99 be *Democra*. Lemmatization is a more complex approach that aims to identify  
100 the root term of the series of similar words. Often this root word is a term that  
101 represents a base concept rather than a sequence of common characters. For  
102 instance, a lemmatization of the example above might be *Democracy*. Lemma-  
103 tization and stemming are often done as a first, data cleaning step along with  
104 tokenization.

105 Given these tokens, we come back to our objective of identifying and labeling  
106 these tokens. To achieve this, we use a technique known as *Named Entity*  
107 *Recognition (NER)*. NER is the process of labeling and categorizing lexical units  
108 extracted from unstructured natural language. This is typically an automated

109 process of comparing tokenized entities found in unstructured text to an existing  
110 structured dictionary or determining the category of an entity based on the  
111 context in which the token exists. Pre-defined categories are often entities such  
112 as people, places, organizations, currencies, etc. This is not a trivial process  
113 as natural language can be quite complex and there is often a large amount of  
114 ambiguity in the meaning of words. Consider, for example, the sentence below.

115 *I watched the Chicago Bulls game last night.*

116 In this example, the term *Bulls* is ambiguous on its own as it is most often  
117 used to reference male cattle. It is only through analysis of contextual infor-  
118 mation that one is able to determine that *Bulls* in this instance refers to the  
119 Chicago-based professional basketball team. A state-of-the art NER applica-  
120 tion, such as Apache OpenNLP, would annotate each of the n-gram tokens in  
121 the example text with Chicago being labeled as a city in the United States, and  
122 the *Chicago Bulls* being labeled as a professional sports team. Today, many  
123 leading NER systems provide close to human-level performance in annotating  
124 unstructured text.

125 Even in the simple example above, the importance of geography is appar-  
126 ent. The region in which cattle are found, the city of Chicago, and dominance  
127 of basketball in discourse all relate to geography, and geographic knowledge  
128 can be leveraged in processing and labeling this information. NER is an im-  
129 portant methodology to GIScientists as it is used in the first task of toponym  
130 disambiguation, which is that task of identifying and labeling a token as a ge-  
131 ographic entity. Toponym disambiguation is typically accomplished through a  
132 look-up/matching process involving a geographic dictionary or what is often  
133 referred to as a *digital gazetteer* [13]. For lesser known or local toponyms, iden-  
134 tification based on geographic context may be used. For instance, Hu et al. [16]  
135 use a geospatial clustering approach and contextual information from surround-  
136 ing words to learn and train a machine learning model to identify toponyms  
137 based on unique spatial and linguistic patterns.

138 Once a token is identified as a toponym, the next challenge is differentiating  
139 it from other toponyms. The nature of human language and culture is that  
140 locations are often assigned the same label. For instance, there are at least 88  
141 different locations in the United States with the name *Washington*, including  
142 cities, monuments, and a federal district. Identifying *which* Washington is the  
143 second task in toponym disambiguation. This is often a challenging task and in-  
144 volves examination of the contextual information and descriptive terms through  
145 which the toponym is referenced. In the Chicago Bulls example above, we can  
146 probabilistically identify *Chicago* as a large city in north-eastern Illinois, USA  
147 in a number of ways. First, Chicago, Illinois has the largest population of any  
148 known Chicago, and is therefore more likely to be mentioned in text. Second,  
149 an NER would likely identify the Bulls basketball team as an entity with a *home*  
150 *town* that also linked to the Chicago in Illinois. Leading research in this area  
151 has used a range of approaches that rely on existing geographic methods and  
152 spatial knowledge including graph-based approaches to linking toponyms [8],

153 topic modeling for disambiguation [18], and co-occurrence models [24]. NER in  
154 general, and toponym disambiguation, more specifically, are central to founda-  
155 tional aspects of GIScience such as geocoding [11] and geographic information  
156 retrieval [17].

### 157 **3.2 Spatial relationships in text**

158 Aside from extracting geographic entities from natural language, researchers  
159 and industry professionals are also very interested in understanding the rela-  
160 tionships between geographic (and non-geographic) entities. Natural language  
161 data provides a rich source of relationship information as contributors of text  
162 often describe these relationships with rich detail. For instance if a body of text  
163 discusses the migratory patterns of people between two cities, this information  
164 could be extracted and represented as a geospatial flow between two network  
165 nodes in a GIS application. NLP extraction methods could also be used to  
166 identify mode of travel and quantify number of migrants.

167 As with toponym disambiguation, identifying and extracting relationships  
168 within unstructured natural language can be difficult. It requires us to deter-  
169 mine which descriptors are applied to which words and which actions involve  
170 which actors. In the field of NLP, this process is called *coreference resolution*.  
171 Coreference resolution is the process of identifying which sub-components of  
172 a sentence or document, refer to which other sub-components, or tokens. In  
173 natural language, we often refer to specific entities or concepts through a vari-  
174 ety of different terms and determine which entity is associated with which idea  
175 can be difficult for humans, let along computational model. Take the following  
176 example.

177 *Seattle gets more days of rain than New York City, but it receives*  
178 *less total rainfall per year.*

179 In this case, we have two proper noun city names, Seattle and New York  
180 City as well as some facts about these cities. A coreference resolution task  
181 arises in the use of the pronoun, *it*. Within the context of this statement, *it*  
182 either refers to Seattle or New York City, and determining the correct referent  
183 is important when assigning information to a location. This may be a trivial  
184 task for a human to resolve, but the ambiguity of human language can often be  
185 difficult to represent computationally.

186 There are many ways to resolve ambiguity of coreferences within natural  
187 language and from a geospatial approach, we can leverage existing geographic  
188 knowledge. Early work in this discipline involved developing methods that ap-  
189 plied a set of grammatical rules to natural language. This often meant the  
190 development of *parse trees* which aimed to represent dependency between to-  
191 kens. Over the past couple of decades, techniques have been developed that take  
192 a probabilistic approach to identifying relationships through the construction of  
193 constituency parsing trees. While not all relationships are spatial, identifying  
194 relationships between entities can sometimes involve a spatial component, be it

195 explicitly spatial (e.g., The museum in Montréal), or through regional or cul-  
196 tural context (e.g., The woman used the Algonquian word for fish). For example,  
197 Vasardani et al. [26] extracted mental representations of urban environments for  
198 use in emergency situations from verbal descriptions of places. Spatial hierar-  
199 chies have also been extracted from user-generated text for use in qualitative  
200 spatial reasoning applications [29]. These, and many other processes demon-  
201 strate that spatial relationships can identified and extracted from unstructured  
202 linguistic content.

203 Having a background in GIScience also means that we are not solely reliant  
204 on the information extracted from natural language. We can use NLP techniques  
205 in conjunction with our existing geospatial expertise [22]. For example, Tobler’s  
206 First Law of Geography can be applied in many cases to leverage the similarity of  
207 features in close proximity. Geographical theories such as *Central Place Theory*  
208 can be used to explain the relationships between nearby settlements, and gravity  
209 models can be employed to identify transfer and flow of entities described in text.

### 210 **3.3 Discovering thematic patterns**

211 Another approach to natural language processing is less concerned with labeling  
212 tokens and identifying individual toponyms in text and more interested in the  
213 broader themes or topics represented in natural language. The idea in this  
214 thematic approach to language is to extract groupings of terms that represent  
215 a set of topics on which a document can be characterized. This is important  
216 for representing ideas in documents as a whole as well as comparing themes  
217 across lexical units. The GIScience community has leveraged this approach  
218 to identify thematic patterns within geographic space and observe changes in  
219 patterns over time. One approach to this problem which has seen extensive use  
220 in the field of GIScience aims to extract themes or topics from corpora through  
221 an unsupervised probabilistic approach, called *Topic Modeling* that identifies  
222 the co-occurrence of tokens within documents. For example, applications of  
223 this technique have been used in clustering social media posts [14], location  
224 recommendation services [15], and ad hoc thematic search engines [3]. For  
225 instance, the Pteraform interactive search platform [1] shown in Figure 1 is  
226 built on top of geographically tagged Wikipedia data, and demonstrates how a  
227 topic modeling approach can be used to geographical depict themes over space  
228 and time. Notably, these approaches tend to ignore the sequence of tokens in  
229 a document or corpora and instead take what is commonly referred to as a  
230 *bag-of-words* approach.

231 Characterizing natural language text by themes is a form of classification,  
232 and there are also other ways we can classify a text. Sentiment analysis is the  
233 process of identifying and examining affective states within text and usually  
234 includes characterizing the emotions and attitudes towards a theme or topic.  
235 Techniques for identifying and extracting sentiment range from examining the  
236 *polarity* of individual tokens, to the *emotional state* of a document or grouping  
237 of tokens. Sentiment analysis is a notoriously challenging field of study as it  
238 involves analysis of subjective information and inference of intention by the

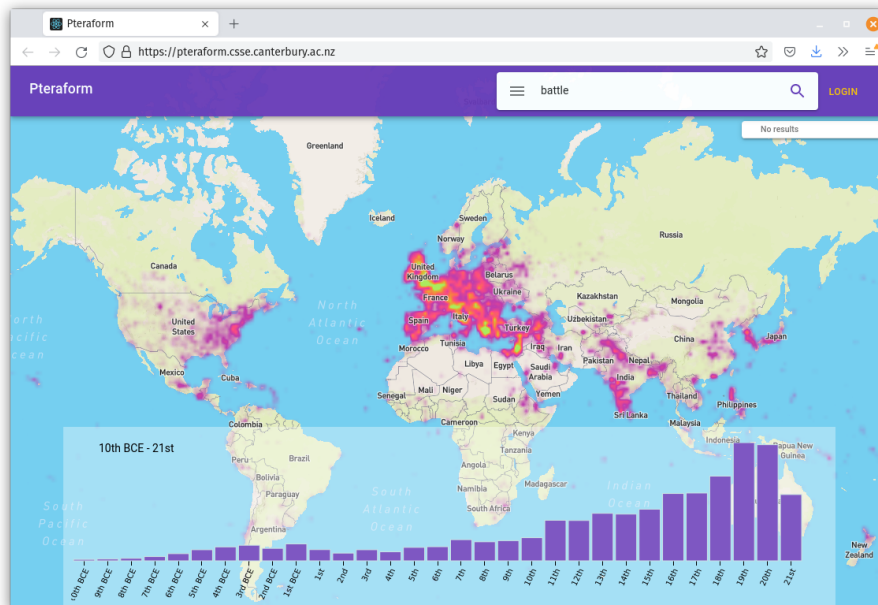


Figure 1: The Pteraform application showing spatial and temporal thematic (keyword: battle) trends over time.

239 language contributor. Applications of sentiment analysis in GIScience have  
 240 included classification of parks through visitor contributions [20], understanding  
 241 disaster response [4], and a plethora of research on attitudes towards travel  
 242 destinations and places of interest [7, 19, 5].

### 243 3.4 Question Answering and Natural Language Generation 244 tion

245 While humans can understand a sentence and the relationship between words  
 246 through reading textual content or verbal communication, computers work in  
 247 the realm of numerical values. Recent advances in NLP have moved towards not  
 248 only representing words as numbers, but also the relationships between words.  
 249 This allows analysts to perform mathematical and logical operations to compare  
 250 terms, extract complex concepts, and better understand the ideas presented in  
 251 natural language. This most often involve assigning a real-value representation  
 252 to a sequence of terms and representing each unit as a numerical vector. Neu-  
 253 ral network-based methods such as word2vec or doc2vec are typically used to  
 254 convert natural language to a series of numerical word vectors or matrices. The  
 255 goal of this approach is to develop *word embeddings*. These encode the meaning  
 256 of words, sentences, and concepts such that words that are closer in meaning  
 257 are also closer in real-value vector space. Essentially, this involves embedding

258 a multi-dimensional concept into a continuous lower-dimensional vector space.  
259 These word embeddings serve as the base unit on which many modern classi-  
260 fication and predictive NLP tasks, including those in the geospatial field, are  
261 performed and often is a key pre-processing step for these other tasks.

262 Other techniques such as recurrent and convolutional neural networks have  
263 been applied to NER tasks with the goal of identifying geographic locations  
264 and places. Adams and McKenzie [2] used a character-level convolutional neu-  
265 ral network to georeference noisy textual content and Cardoso et al. [6] used a  
266 variation on recurrent neural network for toponym resolution in text. Rather  
267 than applying rule-based approaches to identifying the features, deep learning  
268 methods use a representative classification approach to identifying latent fea-  
269 tures in natural language. These models thrive on large training datasets and  
270 the availability of rich and robust training data on which a model can be trained  
271 is critical. Transformer models such as Bidirectional Encoder Representations  
272 from Transformers (BERT) published by Google, have recently emerged. In this  
273 case, a learning model is pre-trained on an exceptionally large, generic dataset  
274 and then fine tuned for a specific task or application area. These *attention-*  
275 *mechanized* transformer models [27] have been shown to improve the accuracy  
276 and relevancy of many NLP-based applications, such as language translation  
277 and document search. These types of models are also being used for geospa-  
278 tial applications such as address validation [30], and identifying the locations of  
279 criminal organizations [23].

280 Question answering is a sub field within natural language processing, infor-  
281 mation retrieval, and artificial intelligence, in which a natural language ques-  
282 tions, typically posed by a human are interpreted by a machine and appropriate  
283 responses are generated. In essence, this a fundamental test for many natu-  
284 ral language processing techniques in that responding to a question requires  
285 comprehension of the concepts presented in the question itself. This approach  
286 involves a high level of automated reasoning. The field of *geographic* question  
287 answering has recently emerged with the goal of identifying and understanding  
288 the relationship between geographic features, places, and people through the  
289 use of many deep learning approaches. The nuances of geospatial concepts in  
290 natural language is unique and designing a system that can interpret and un-  
291 derstand these concepts and relationships can be challenging. Take for example  
292 the question below.

293 *How many people live in the capital of the third largest country on*  
294 *earth?*

295 Not only does the question above require entities to be extracted and labeled  
296 through an NER task or thematically encoded through a neural network, but  
297 it also requires leveraging existing geospatial knowledge such as administrative  
298 boundary hierarchies. For instance a capital is a city, a city exists within state,  
299 and a state with country. The term *largest* is ambiguous here as well as it is  
300 unclear if this is in reference to population volume or physical area. Finally,  
301 *third*, it requires a system to know the populations or areas of all countries,



302 rank them, and extract the third largest. While natural language processing  
303 techniques are increasingly able to learn many of these concepts, understanding  
304 the relationships and answering the question also involves accessing knowledge  
305 graphs, geographic databases, and range of other technologies. This area is  
306 proving to be a burgeoning subfield of GIScience. Scheider et al. [25] discuss  
307 the challenges associated with building a question-based geographic information  
308 system and how existing spatial techniques and technologies can be used within  
309 such a service. Mai et al. [21] demonstrate possibilities and limitations of geo-  
310 graphic question answering through the use of geospatially enabled knowledge  
311 graph embeddings.

312 The complement to question answering is *natural language generation* (NLG).  
313 This approach aims to generate natural language text or speech based on seman-  
314 tically encoded concepts. In many ways, the second part of question answering  
315 demands generating natural language based on the interpreted understanding  
316 of the original question. Applied work in this field has predominantly focused  
317 on automating reports and responses to questions. Within the geographical sci-  
318 ences we see NLG techniques being applied to generating weather reports [12],  
319 descriptions of places and remotely sensed imagery [10], and the broader focus  
320 on chatbots and automated assistants capable of responding to basic questions.

## 321 4 Challenges

322 A number of challenges exist within the domain of natural language processing  
323 and many of them are uniquely spatial. Many of these were mentioned in the  
324 previous sections, but here the challenges are outlined in further detail.

325 Using NLP to interpret fine-grained spatial relationships in text is an active  
326 area of research. While many current NLP approaches are able to identify con-  
327 cepts, ideas, and relationships within natural language, surprisingly few of them  
328 explicitly model spatial relationships. Concepts such as spatial autocorrelation  
329 are fundamental to GIScience, yet very few approaches incorporate this idea in  
330 the process of understanding natural language.

331 Spatial cognition is a branch of cognitive psychology that studies the ways in  
332 which people use spatial information to gain knowledge, self locate, and wayfind.  
333 This field is closely linked with natural language processing in that understand-  
334 ing human-contributed natural language necessitates an understanding of how  
335 humans conceptualize space and communicate those concepts in language [9].  
336 This presents a unique challenge, as how humans conceptualize and commu-  
337 nicate spatial concepts is not fully understood, therefore making it difficult to  
338 train a computational model to represent spatial information in a similar way.

339 While substantial advances have been made in toponym disambiguation and  
340 co-reference resolution within NLP research, it still remains as a challenge.  
341 Given that places are labeled by humans, they tend to change over time, or have  
342 multiple, often localized, names. Humans reference places in different ways and  
343 the ability to identify a single place based on various colloquial references to the  
344 location remains a challenge.

345 Lastly, the automated generation of spatially-aware narratives is a chal-  
346 lenge area that will likely see advances in the coming years. This will involve  
347 the integration of NLP more substantially in location-based systems such as  
348 tourism applications and will leverages geographic knowledge graphs and exist-  
349 ing gazetteers.

## 350 5 Learning Objectives

351 The objective of this chapter is to

- 352 • Explain how natural language processing is being used in geographic in-  
353 formation science applications.
- 354 • Differentiate between some of the key uses of natural language processing  
355 in geography and GIScience.
- 356 • Identify how *spatial is special* in the context of natural language process-  
357 ing.
- 358 • Identify challenges and future directions for applications of NLP in GI-  
359 Science.

## 360 6 Instructional Assessment Questions

- 361 1. What does the field of geography bring to the discussion of natural lan-  
362 guage processing?
- 363 2. What are the two components necessary for toponym disambiguation?
- 364 3. How is geographic question answering different than traditional question  
365 answering?
- 366 4. What is the difference between stemming and lemmatization?

## 367 7 Additional Resources

- 368 • Apache OpenNLP <https://opennlp.apache.org/index.html>
- 369 • Stanford Natural Language Processing Toolkit <https://nlp.stanford.edu/>
- 370 • Python Natural Language Toolkit module <https://www.nltk.org/>
- 371 • R GeoParser package <https://rdr.io/cran/geoparser/>
- 372 • An Extensible and Unified Platform for Evaluating Geoparsers <https://geoai.geog.buffalo.edu/EUPEG/>
- 373 • Creating the Corpus (Spatial Language) <https://geospatiallanguage.massey.ac.nz/creatingthecorpus.htm>
- 374 • EarthLings (Computational Linguistic Atlas) <http://www.earthlings.io/>

## 375 References

- 376 [1] Benjamin Adams. Chronotopic information interaction: integrating tem-  
377 poral and spatial structure for historical indexing and interactive search.  
378 *Digital Scholarship in the Humanities*, 2020.
- 379 [2] Benjamin Adams and Grant McKenzie. Crowdsourcing the character of a  
380 place: Character-level convolutional networks for multilingual geographic  
381 text classification. *Transactions in GIS*, 22(2):394–408, 2018.
- 382 [3] Benjamin Adams, Grant McKenzie, and Mark Gahegan. Frankenplace:  
383 interactive thematic mapping for ad hoc exploratory search. In *Proceedings*  
384 *of the 24th international conference on world wide web*, pages 12–22, 2015.
- 385 [4] Abdullah Alfarrarjeh, Sumeet Agrawal, Seon Ho Kim, and Cyrus Shahabi.  
386 Geo-spatial multimedia sentiment analysis in disasters. In *2017 IEEE In-*  
387 *ternational Conference on Data Science and Advanced Analytics (DSAA)*,  
388 pages 193–202. IEEE, 2017.
- 389 [5] Andrea Ballatore and Benjamin Adams. Extracting place emotions from  
390 travel blogs. In *Proceedings of AGILE*, volume 2015, pages 1–5, 2015.
- 391 [6] Ana Bárbara Cardoso, Bruno Martins, and Jacinto Estima. Using recurrent  
392 neural networks for toponym resolution in text. In *EPIA Conference on*  
393 *Artificial Intelligence*, pages 769–780. Springer, 2019.
- 394 [7] Mario Cataldi, Andrea Ballatore, Ilaria Tiddi, and Marie-Aude Aufaure.  
395 Good location, terrible food: detecting feature sentiment in user-generated  
396 reviews. *Social Network Analysis and Mining*, 3(4):1149–1163, 2013.
- 397 [8] Hao Chen, Stephan Winter, and Maria Vasardani. Georeferencing places  
398 from collective human descriptions using place graphs. *Journal of Spatial*  
399 *Information Science*, 2018(17):31–62, 2018.
- 400 [9] Song Gao, Krzysztof Janowicz, Daniel R Montello, Yingjie Hu, Jiue-An  
401 Yang, Grant McKenzie, Yiting Ju, Li Gong, Benjamin Adams, and Bo Yan.  
402 A data-synthesis-driven method for detecting and extracting vague cogni-  
403 tive regions. *International Journal of Geographical Information Science*,  
404 31(6):1245–1271, 2017.
- 405 [10] Albert Gatt and Emiel Krahmer. Survey of the state of the art in natural  
406 language generation: Core tasks, applications and evaluation. *Journal of*  
407 *Artificial Intelligence Research*, 61:65–170, 2018.
- 408 [11] Daniel W Goldberg, John P Wilson, and Craig A Knoblock. From text  
409 to geographic coordinates: the current state of geocoding. *URISA journal*,  
410 19(1):33–46, 2007.
- 411 [12] Eli Goldberg, Norbert Driedger, and Richard I Kittredge. Using natural-  
412 language processing to produce weather forecasts. *IEEE Expert*, 9(2):45–53,  
413 1994.

- 414 [13] Linda L Hill. Core elements of digital gazetteers: placenames, categories,  
415 and footprints. In *International Conference on Theory and Practice of*  
416 *Digital Libraries*, pages 280–290. Springer, 2000.
- 417 [14] Liangjie Hong and Brian D Davison. Empirical study of topic modeling  
418 in twitter. In *Proceedings of the first workshop on social media analytics*,  
419 pages 80–88, 2010.
- 420 [15] Bo Hu and Martin Ester. Spatial topic modeling in online social media for  
421 location recommendation. In *Proceedings of the 7th ACM conference on*  
422 *Recommender systems*, pages 25–32, 2013.
- 423 [16] Yingjie Hu, Huina Mao, and Grant McKenzie. A natural language process-  
424 ing and geospatial clustering framework for harvesting local place names  
425 from geotagged housing advertisements. *International Journal of Geograph-*  
426 *ical Information Science*, 33(4):714–738, 2019.
- 427 [17] Christopher B Jones and Ross S Purves. Geographical information  
428 retrieval. *International Journal of Geographical Information Science*,  
429 22(3):219–228, 2008.
- 430 [18] Yiting Ju, Benjamin Adams, Krzysztof Janowicz, Yingjie Hu, Bo Yan,  
431 and Grant McKenzie. Things and strings: improving place name disam-  
432 biguation from short texts by combining entity co-occurrence with topic  
433 modeling. In *European Knowledge Acquisition Workshop*, pages 353–367.  
434 Springer, 2016.
- 435 [19] Hanhoon Kang, Seong Joon Yoo, and Dongil Han. Senti-lexicon and im-  
436 proved naïve bayes algorithms for sentiment analysis of restaurant reviews.  
437 *Expert Systems with Applications*, 39(5):6000–6010, 2012.
- 438 [20] Anna Kovacs-Györi, Alina Ristea, Ronald Kolcsar, Bernd Resch, Alessan-  
439 dro Crivellari, and Thomas Blaschke. Beyond spatial proximity—classifying  
440 parks and their visitors in london based on spatiotemporal and sentiment  
441 analysis of twitter data. *ISPRS International Journal of Geo-Information*,  
442 7(9):378, 2018.
- 443 [21] Gengchen Mai, Krzysztof Janowicz, Ling Cai, Rui Zhu, Blake Regalia,  
444 Bo Yan, Meilin Shi, and Ni Lao. Se-kge: A location-aware knowledge graph  
445 embedding model for geographic question answering and spatial semantic  
446 lifting. *Transactions in GIS*, 24(3):623–655, 2020.
- 447 [22] Grant McKenzie and Benjamin Adams. Juxtaposing Thematic Regions De-  
448 rived from Spatial and Platial User-Generated Content. In Eliseo Clemen-  
449 tini, Maureen Donnelly, May Yuan, Christian Kray, Paolo Fogliaroni, and  
450 Andrea Ballatore, editors, *13th International Conference on Spatial Infor-*  
451 *mation Theory (COSIT 2017)*, volume 86 of *Leibniz International Proceed-*  
452 *ings in Informatics (LIPIcs)*, pages 20:1–20:14, Dagstuhl, Germany, 2017.  
453 Schloss Dagstuhl–Leibniz-Zentrum fuer Informatik.

- 454 [23] Javier Osorio and Alejandro Beltran. Enhancing the detection of criminal  
455 organizations in mexico using ml and nlp. In *2020 International Joint*  
456 *Conference on Neural Networks (IJCNN)*, pages 1–7. IEEE, 2020.
- 457 [24] Simon Overell and Stefan R uger. Using co-occurrence models for place-  
458 name disambiguation. *International Journal of Geographical Information*  
459 *Science*, 22(3):265–287, 2008.
- 460 [25] Simon Scheider, Enkhbold Nyamsuren, Han Kruiger, and Haiqi Xu. Geo-  
461 analytical question-answering with gis. *International Journal of Digital*  
462 *Earth*, 14(1):1–14, 2021.
- 463 [26] Maria Vasardani, Sabine Timpf, Stephan Winter, and Martin Tomko. From  
464 descriptions to depictions: A conceptual framework. In *International Con-*  
465 *ference on Spatial Information Theory*, pages 299–319. Springer, 2013.
- 466 [27] Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones,  
467 Aidan N Gomez, Lukasz Kaiser, and Illia Polosukhin. Attention is all you  
468 need. In *Advances in neural information processing systems*, pages 5998–  
469 6008, 2017.
- 470 [28] Jonathan J Webster and Chunyu Kit. Tokenization as the initial phase in  
471 nlp. In *COLING 1992 Volume 4: The 14th International Conference on*  
472 *Computational Linguistics*, 1992.
- 473 [29] Xiaoyu Wu, Jianying Wang, Li Shi, Yong Gao, and Yu Liu. A fuzzy formal  
474 concept analysis-based approach to uncovering spatial hierarchies among  
475 vague places extracted from user-generated data. *International Journal of*  
476 *Geographical Information Science*, 33(5):991–1016, 2019.
- 477 [30] Liuchang Xu, Zhenhong Du, Ruichen Mao, Feng Zhang, and Renyi Liu.  
478 Gsam: A deep neural network model for extracting computational repre-  
479 sentations of chinese addresses fused with geospatial feature. *Computers,*  
480 *Environment and Urban Systems*, 81:101473, 2020.