

1 An empirical study on the names of points of 2 interest and their changes with geographic 3 distance

4 **Yingjie Hu**

5 GSDA Lab, Department of Geography, University of Tennessee, Knoxville, USA
6 yhu21@utk.edu

7 **Krzysztof Janowicz**

8 STKO Lab, Department of Geography, University of California, Santa Barbara, USA
9 jano@ucsb.edu

10 — Abstract —

11 While Points Of Interest (POIs), such as restaurants, hotels, and barber shops, are part of
12 urban areas irrespective of their specific locations, the names of these POIs often reveal valuable
13 information related to local culture, landmarks, influential families, figures, events, and so on.
14 Place names have long been studied by geographers, e.g., to understand their origins and relations
15 to family names. However, there is a lack of large-scale empirical studies that examine the
16 *localness* of place names and their changes with geographic distance. In addition to enhancing our
17 understanding of the coherence of geographic regions, such empirical studies are also significant
18 for geographic information retrieval where they can inform computational models and improve
19 the accuracy of place name disambiguation. In this work, we conduct an empirical study based on
20 112,071 POIs in seven US metropolitan areas extracted from an open Yelp dataset. We propose
21 to adopt term frequency and inverse document frequency in geographic contexts to identify local
22 terms used in POI names and to analyze their usages across different POI types. Our results
23 show an uneven usage of local terms across POI types, which is highly consistent among different
24 geographic regions. We also examine the decaying effect of POI name similarity with the increase
25 of distance among POIs. While our analysis focuses on urban POI names, the presented methods
26 can be generalized to other place types as well, such as mountain peaks and streets.

27 **2012 ACM Subject Classification** H.2.8 Spatial databases and GIS; H.3.1 Linguistic processing.

28 **Keywords and phrases** Place names; points of interest; geographic information retrieval; se-
29 mantic similarity; geospatial semantics.

30 **Digital Object Identifier** 10.4230/LIPIcs.GIScience.2018.23

31 **1** Introduction

32 People name the environment that surrounds them to communicate about it. Almost every
33 aspect of geographic space that can be described and depicted can be named. It has been
34 suggested that place names, or toponyms, play a key role in stabilizing the otherwise un-
35 wieldy space into more manageable textual inscriptions [38, 25, 42]. From a perspective
36 of *space* and *place* [45], the creation of a place name signifies the important moment when
37 people explicitly integrate human experience with space.

38 Place names, made available via digital gazetteers, power GIS, geographic information
39 retrieval (GIR), and modern search engines and recommender systems more broadly [20, 13,
40 47]. After all, people communicate using place names not coordinates. Interestingly, and
41 in difference to human geography, most GIR research simply uses place names as identifiers
42 instead of examining how those names were formed and how similar they are to nearby



© Yingjie Hu and Krzysztof Janowicz;
licensed under Creative Commons License CC-BY

10th International Conference on Geographic Information Science (GIScience 2018).

Leibniz International Proceedings in Informatics



LIPICs Schloss Dagstuhl – Leibniz-Zentrum für Informatik, Dagstuhl Publishing, Germany

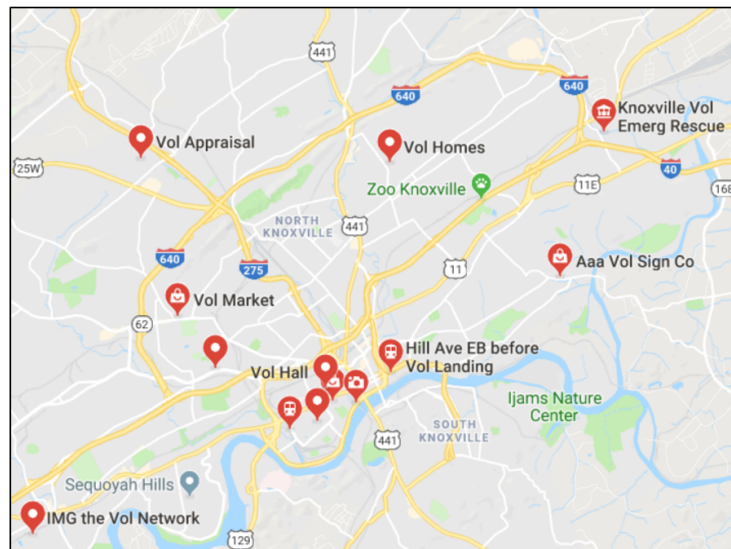
43 names. This is understandable since we are often interested in questions such as *What are*
 44 *the best Italian restaurants within 10 minutes driving?* instead of the specific names of these
 45 restaurants or what they reveal about the history of a region, such as immigration trends.

46 Place names have long been studied in human geography with a traditional focus on
 47 etymology and place taxonomies [52, 40]. For example, the place name *Las Vegas* means *The*
 48 *Meadows* in Spanish and points to the former abundance of wild grasses and desert springs,
 49 both of which were crucial information for travelers and led to the descriptive place name.
 50 While such studies provide in-depth explanation of place names, they are often limited to
 51 case-by-case examinations with qualitative descriptions. This could include studies focusing
 52 on specific regions, names, places types, and so forth.

53 In contrast, this work is based on more than 110,000 place names of different types
 54 distributed across seven metropolitan areas within the US. Our focus is on uncovering term
 55 usage patterns and their relations with geographic locations, e.g., as modeled by a decaying
 56 influence or local names with increasing distance. There are several reasons for performing
 57 such a large-scale, data-driven study. First, place names reveal many social and cultural
 58 characteristics, and can help us understand various aspects of urban areas. Previous research
 59 in human geography has considered place names, such as street names, as *city-text* embedded
 60 in the cityscape [6, 7]. A systematic examination on these city-texts, can help expand
 61 our knowledge of the studied regions. Second, large-scale empirical research examining
 62 place names can aid in discovering common principles in place naming and relations to
 63 environments. This can be distinguished from case-by-case place name studies in which the
 64 discovered knowledge often cannot be generalized to other names or geographic areas. Third,
 65 such studies can facilitate the development of computational models for places. We can
 66 integrate the discovered common principles, socio-cultural characteristics, and local terms
 67 into computational models, e.g., via an implemented knowledge base, to better support tasks
 68 such as place name disambiguation [4, 27, 37, 17]. This last point is a key strength of this
 69 work. Our results can act as a quantitative foundation for the localness of identifiers *per*
 70 *place*, which will enable future research to push the envelop on place name disambiguation.
 71 In fact, our previous *Things and Strings* place disambiguation method [22] has demonstrated
 72 the usefulness and need for combining geographic and linguistic information.

73 The names of Points Of Interest (POIs), such as restaurants, hotels, grocery stores, and
 74 auto repairs, are examined in this work. These POI names are from an open dataset released
 75 by Yelp, a company that provides search services for local businesses. POIs play important
 76 roles in supporting many aspects of our daily life [33, 36, 51]. One reason we select POI names
 77 for this study is that these names reflect more of the diverse views of the general public,
 78 since the business owners can decide on names themselves. This can be differentiated from
 79 other place names, such as street names, which often result from political and administrative
 80 decisions [7, 1, 41]. In addition, the names of POIs often contain local information, such
 81 as city or state names, natural or man-made geographic features, vernacular names, local
 82 families (e.g., a family-owned business), language patterns, local cultural differences, and
 83 others. Figure 1 shows an example of searching for the word “Vol” in the city of Knoxville,
 84 Tennessee, USA using Google Maps. It returns many places which use this term as part of
 85 their names, as “Vol” is the local nickname of the popular football team “Volunteer”. The use
 86 of American sports team names in toponyms was also noted in previous human geography
 87 research [8]. In GIR and place name disambiguation, understanding the link between “Vol”
 88 and the city of Knoxville can help locate related place names more accurately.

89 More specifically, we aim to answer the following questions in this work: 1) what are the
 90 local terms that are used in POIs in different geographic areas? 2) how are these local terms



■ **Figure 1** An example of POIs in Knoxville, TN, USA that use “Vol” as part of their names.

91 used in different types of POIs, such as restaurants, hotels, and barber shops? and 3) how
 92 do POI names change with geographic distance? **The contributions of this paper are**
 93 **as follows:**

- 94 ■ We propose adopting the technique of term frequency and inverse document frequency in
 95 geographic contexts to identify local terms used in POIs in different metropolitan areas.
- 96 ■ We find an uneven usage of local terms in the names of POIs across POI types, and such
 97 an uneven usage is highly consistent across the seven studied metropolitan areas.
- 98 ■ We test two types of models, count-based vector and word2vec, for understanding and
 99 capturing the distance decay effect of the similarity of POI names.

100 The remainder of this paper is structured as follows. Section 2 reviews related work
 101 on place names and toponym disambiguation. Section 3 describes the dataset used in this
 102 study and an exploratory data analysis. Section 4 presents methods and experiments for
 103 identifying local terms from POI names, examining their usages across POI types, and
 104 modeling the distance decay effect of POI name similarity. Section 5 summarizes this work
 105 and discusses future directions.

106 2 Related Work

107 Place names have attracted the interest of many researchers in geography. For decades,
 108 geographers have been collecting and categorizing place names, studying their origins, and
 109 understanding their meanings [50, 52, 35]. It has been argued that the act of assigning a
 110 name to *space* plays a key role in producing the social construct of *place* [40]. As suggested
 111 by Carter [10], place names transform space into knowledge that can be read. The social,
 112 cultural, and political implications of place names have been widely studied [5, 6]. Ex-
 113 amples include the renaming of streets after the establishment of a new regime to memorize
 114 new stories [30, 41], the use of street names to challenge racism [2, 3], and assigning more
 115 marketable names to local businesses and hospitals [39, 24].

116 Digital gazetteers provide systematic organizations of place names (N), place types (T),
 117 and spatial footprints (F) [16, 13]. As valuable knowledge bases, gazetteers provide import-
 118 ant functions for various applications by connecting the three core components. The key
 119 functions of a gazetteer include lookup ($N \rightarrow F$), type-lookup ($N \rightarrow T$), and reverse-lookup
 120 ($F(\times T) \rightarrow N$) [19]. The first case, for example, corresponds to a query for the spatial
 121 footprint of the place name *CMS Auto Care*, the second to the place type, and the third to
 122 the place names given the spatial footprint and a place type (e.g., *Automotive*). Research
 123 was conducted to enrich gazetteers with (vague) place names and their fuzzy spatial foot-
 124 prints. Jones et al. [21], for instance, used a search engine to harvest geographic entities
 125 (e.g., hotels) related to vague place names (e.g., “Mid-Wales”), and utilized the locations of
 126 these harvested entities to construct vague boundaries. Flickr photos present a natural link
 127 between textual tags and locations, and have been used in many studies on identifying the
 128 boundaries of vague places and regions [15, 26, 18, 28]. Twaroch and Jones [46] developed a
 129 Web-based platform, called “People’s Place Names”, which invites local people to contribute
 130 vernacular place names.

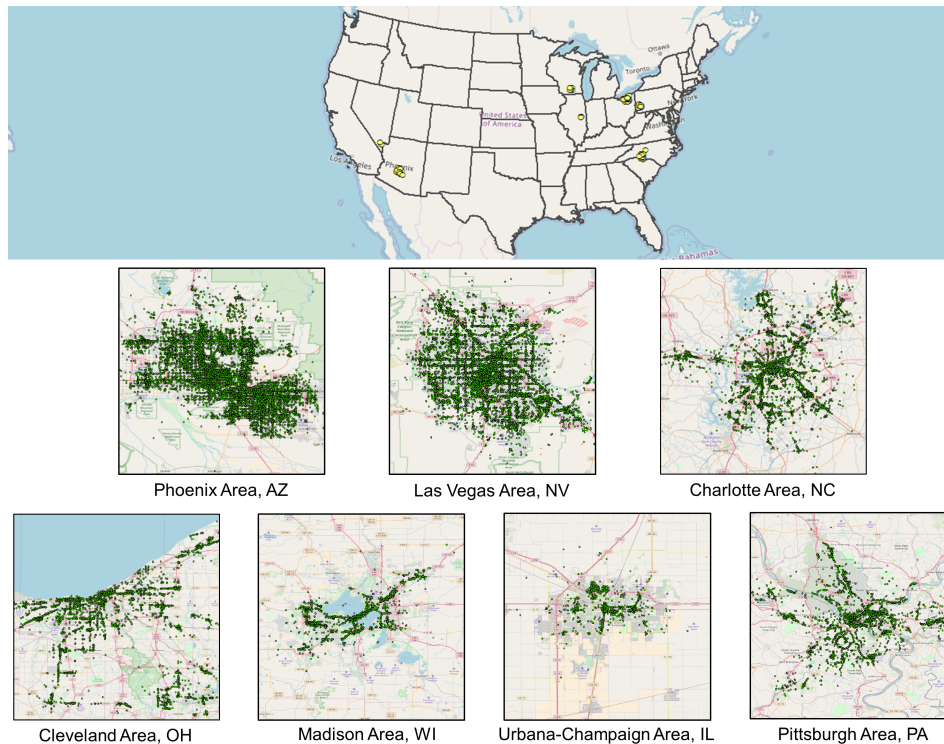
131 In geographic information retrieval [20], place names are frequently discussed in the
 132 context of place name disambiguation. Since different place names can refer to the same
 133 place instance and the same place name can refer to different place instances, it is challenging
 134 to determine which place instance was referred to by a name in text, e.g., the abstract of
 135 a news article [4, 27]. Gazetteers have been used in many ways for supporting place name
 136 disambiguation. Based on the related places in a gazetteer (e.g., higher-level administrative
 137 units), researchers developed methods, such as co-occurrence models [37] and conceptual
 138 density [9], to disambiguate place names. Based on the spatial footprints of place instances,
 139 researchers designed heuristics for place name disambiguation, e.g., place names mentioned
 140 in the same document generally share the same geographic context [29, 43]. The process of
 141 recognizing and resolving place names from texts is called *geoparsing* [12, 23, 14, 49]. Place
 142 names are also examined in studies on toponym matching and geo-data conflation [44].

143 Few existing studies, however, have empirically examined the term usage of place names
 144 and their relations with geographic locations based on large datasets. Longley, Cheshire,
 145 and colleagues [31, 11] investigated the geospatial distributions of surnames based on the
 146 data from the Electoral Register for Great Britain and delineated surname regions. Their
 147 study is related to our work, since family names are included in the names of some local
 148 business. We perform an empirical study based on a large number of POI names in different
 149 US metropolitan areas. Compared with the existing literature, this work is unique in that
 150 it examines the local terms in POI names, explores the term usage patterns, and analyzes
 151 the relations of POI names to geographic locations as well as their decay in this relationship
 152 over distance.

153 **3** Dataset

154 We first describe the data used in this empirical study, which is an open POI dataset from
 155 Yelp (<https://www.yelp.com/dataset>). The original dataset contains POIs from 11 met-
 156 ropolitan areas in four countries: the US, Canada, the UK, and Germany. Considering the
 157 language differences in POI names (e.g., German and English) and the barrier effects of
 158 country borders, we focus on the seven metropolitan areas within the US, which contain
 159 112,071 POIs. Each POI data record has the POI name, city name, state name, latitude-
 160 longitude coordinates, and other information, such as the number of reviews and average
 161 rating. Figure 2 shows the general locations of the seven metropolitan areas and the geo-

graphic distributions of the POIs in each of these areas.



■ **Figure 2** The seven US metropolitan areas and their POIs used for this study.

162

163 We start by performing an exploratory analysis on the term usage frequency in the POI
 164 names. It has been found that Zipf's law exists in the usage of terms in natural language
 165 texts [32], namely the frequency of a term is proportional to the inverse of its frequency
 166 rank among all terms (Equation 1).

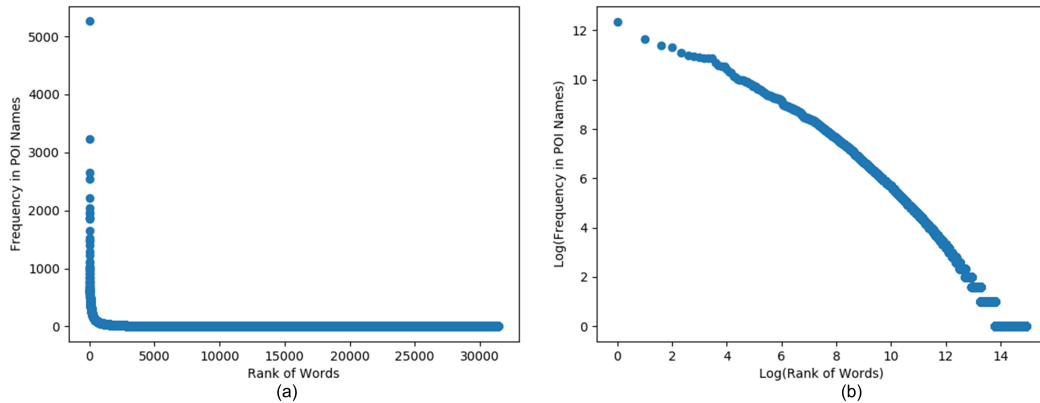
167

$$f \propto \frac{1}{r} \quad (1)$$

168 where f is the frequency of a term and r is the rank of the term among all terms based
 169 on frequency. According to Zipf's law, a small number of terms are used highly frequently
 170 while most others are used only occasionally. The names of POIs are different from nat-
 171 ural language texts in that they are typically not complete sentences but phrases. In this
 172 situation, does Zipf's law still hold in POI names?

173

174 To answer this question, we develop a Python script which reads through the names
 175 of the POIs in the seven metropolitan areas, counts the frequencies of all terms contained
 176 in each name, and ranks the terms based on their frequencies. We then use the ranks as
 177 the horizontal coordinates and term frequencies as the vertical coordinates, and the result
 178 is shown in Figure 3(a). As can be seen, there is a highly skewed distribution of term
 179 frequency with a long tail, which suggests that a small number of terms are used much more
 180 frequently than most other terms. In fact, Figure 3(a) shows almost a right angle fall-off
 181 since the term frequency decreases rapidly with a small increase of the rank. The log-log
 182 plot of the frequencies and ranks is shown in Figure 3(b), and we see almost a straight line.
 183 To quantitatively measure the match of term usage in POI names to Zipf's law, we fit a
 linear regression model with $\log f = A + b \log r$, and obtained an R-squared value of 0.962.



■ **Figure 3** Term frequencies and their ranks in POI names: (a) original values; (b) log-log plot.

184 Based on this exploratory analysis, we conclude that the term usage in POI names also
 185 follow Zipf’s law, even though POI names are usually not complete sentences. The top 10
 186 most frequent terms in POI names in this Yelp dataset are: *the, and, of, center, pizza, grill,*
 187 *spa, bar, auto, restaurant.* These most frequent terms reflect the inherent characteristics of
 188 POI names and POI types. It is worth noting that the most frequent terms in POI names
 189 may change across countries, depending on the corresponding cultures and lifestyles.

190 4 Data Analysis

191 In this section, we perform in-depth analyses on POI names. We organize this section into
 192 three subsections based on the three core components of gazetteers [16]. Thus, the first
 193 subsection focuses on *place names*, and aims to identify the local-specific terms used in
 194 these POI names. The second subsection looks into the interaction between POI names and
 195 *place types*, and examines the usage of local terms in different POI types. Finally, the third
 196 subsection analyzes the change of POI names with geographic distance based on the *spatial*
 197 *footprints* of the POIs.

198 4.1 Identifying local terms from POI names

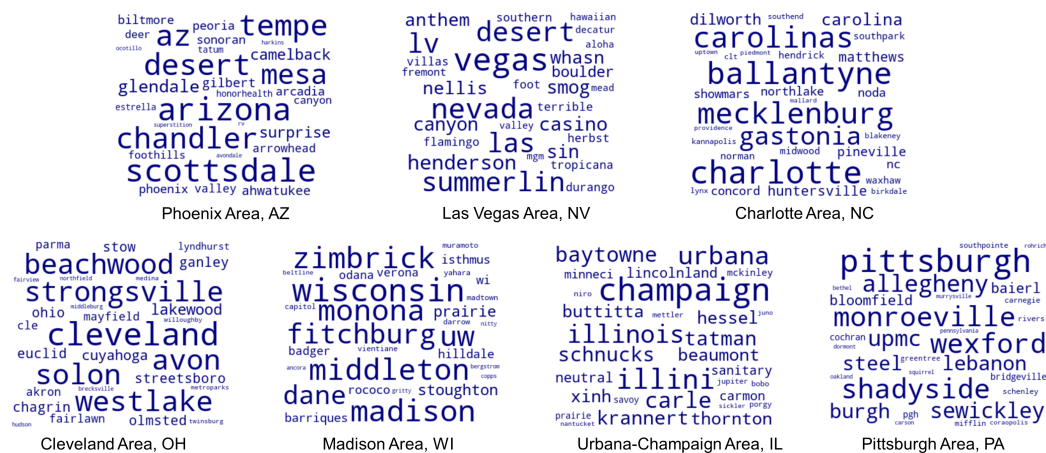
199 In this first analysis, we attempt to answer the question: *what are the local terms used in*
 200 *the names of POIs in a geographic area?* While not every POI name contains local specific
 201 terms, some names are influenced by local factors, such as the “Vol” example discussed in
 202 the Introduction. We consider local terms as those frequently used in a local geographic
 203 area but less likely to be used in other areas. Identifying these local terms can help enhance
 204 computational models for place name disambiguation. We make use of the technique, term
 205 frequency and inverse document frequency (TF-IDF), a method commonly used in inform-
 206 ation retrieval, and adapt it to the context of geography. Equation 2 shows the adapted
 207 version of TF-IDF.

$$208 \quad w_{ij} = tf_{ij} \times \log \frac{|G|}{|G_j|} \quad (2)$$

209 where w_{ij} is the weight of a term j in geographic area i , tf_{ij} is the frequency of term j in area
 210 i , $|G|$ is the total number of geographic areas in a study (which is seven in our case), and
 211 $|G_j|$ is the number of geographic areas that contain the term j . TF-IDF will highlight the

212 terms that are frequently used in a local area, while reducing the weights of those commonly
 213 exist in POI names everywhere. In fact, the weights of the terms that occur in all seven
 214 metropolitan areas will become zero based on Equation 2.

215 Before applying the adapted TF-IDF to the POI names, we perform several data pre-
 216 processing steps. All POI names are converted to lowercase, and punctuations in POI names
 217 are removed. We did not remove typical stop words, such as “the” and “of”, since the term
 218 frequencies in POI names are not the same as other natural language texts, as shown in the
 219 exploratory analysis. Thus, typical stop words may not be stop words in the names of POIs.
 220 We also performed one special step for this analysis by counting the exact same POI names
 221 only once within a metropolitan area. The rationale behind this step is that term frequency
 222 can be increased in two situations: 1) one term is used by many different POIs (e.g., the
 223 term “Vol” is used in the names of many POIs); and 2) one word is used by the same
 224 POI business which simply shows up many times in a metropolitan area (e.g., “walmart”).
 225 We would prefer to keep the terms in the first situation, since those are endorsed by many
 226 different POIs and are more likely to be valid local terms than those in the second situation.
 227 After removing these repeating POI names, we group the names that belong to the same
 228 metropolitan areas using the bag-of-words model. We then use the adapted TF-IDF to
 229 identify local terms. Figure 4 shows the top 30 local terms identified for each of the seven
 metropolitan areas.



230 ■ **Figure 4** Local terms identified based on the POI names in the seven US metropolitan areas.

230

231 We can group the identified local terms into the following categories:

232

232 ■ **City names:** This is the most common type. POI names in all seven metropolitan areas
 233 contain city names, such as *scottsdale*, *las vegas*, *charlotte*, and *cleveland*.

234

234 ■ **State names:** This is similar to city names. State names, such as *arizona* and *wisconsin*,
 235 are used in POI names. There are also name abbreviations, such as *az* and *wi*.

236

236 ■ **Natural features:** Examples include *desert* and *canyon* in Phoenix and Las Vegas
 237 areas, *prairie* in Madison and Urbana-Champaign areas, and *rivers* in Pittsburgh area.

238

238 ■ **Sports teams:** Examples include *badger* in Wisconsin and *illini* in Illinois.

239

239 ■ **Family names:** A notable example is *zimbrick* in Madison, Wisconsin, which is a re-
 240 gional car dealer started by *John Zimbrick* ([http://www.zimbrickbuickgmceast.com/
 241 Zimbrick-History](http://www.zimbrickbuickgmceast.com/Zimbrick-History)).

242

242 ■ **Local cultures:** Terms such as *sin* and *casino* are observed in the POI names in Las
 243 Vegas, while the term *steel* is observed in the POI names in Pittsburgh area.

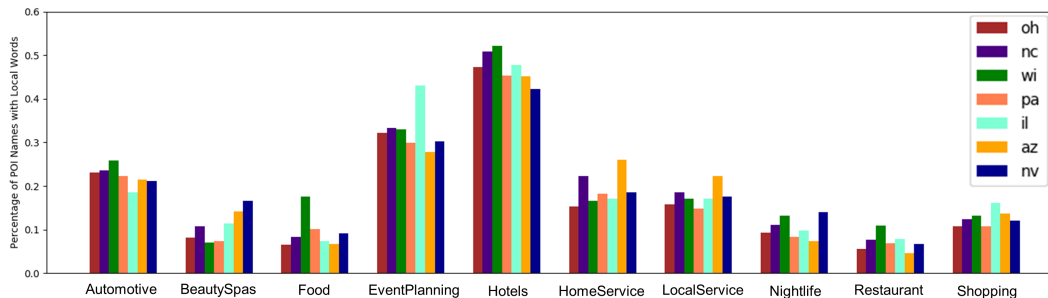
244 **4.2 Examining local term usage in different POI types**

245 The first analysis identified the local terms used in POI names in each geographic area.
 246 However, do POIs in different types have similar probabilities in using local terms as part
 247 of their names? In addition, are there regional differences in using local terms for names
 248 among POI types? In this second analysis, we aim to answer these questions.

249 In order to examine the interaction between POI names and POI types, we need to first
 250 divide the dataset based on POI types. Yelp has grouped their POIs into 23 root categories
 251 which include *Restaurants, Shopping, Food, Hotels & Travel*, and other categories. We make
 252 use of these Yelp POI categories, and the POIs in each metropolitan area are divided into
 253 subsets based on their categories. Yelp allows one POI to belong to multiple categories (e.g.,
 254 one POI can be both *Restaurants* and *Nightlife*), and therefore the same POI is put into
 255 more than one subset when multiple categories exist. Not all metropolitan areas contain
 256 POIs in all 23 categories. In addition, one metropolitan area may contain only a small
 257 number of POIs in a certain category, which can cause a biased result if those POIs are
 258 directly used for analysis. Thus, we only examine the POI types which are shared by all
 259 seven metropolitan areas and have at least one hundred POI instances in each area. Based
 260 on these criteria, we are left with ten categories, which are *Automotive, Beauty & Spas,*
 261 *Food, Event Planning & Services, Hotels & Travel, Home Services, Local Services, Nightlife,*
 262 *Restaurants, and Shopping.* The TF-IDF weights from the first analysis are then re-used,
 263 and we extract the top 100 terms that have the highest TF-IDF weights in each metropolitan
 264 area and use them as the local terms. The percentage of POI names in each POI type that
 265 contain local terms is calculated using Equation 3:

$$266 \quad Pr_{ij} = |LP_{ij}|/|P_{ij}| \quad (3)$$

267 where $|LP_{ij}|$ is the number of POI names that contain any of the local terms in metropolitan
 268 area i in POI type j , $|P_{ij}|$ is the total number of POI names in metropolitan area i in POI
 type j , and Pr_{ij} is the calculated percentage. The result is shown in Figure 5.



269 **Figure 5** The percentages of POI names that contain local terms across POI types and different
 270 metropolitan areas.

270 Two things can be observed in Figure 5. First, there is an uneven usage of local terms
 271 across POI types. Overall, it seems people (business owners) are more likely to include local
 272 terms in the names of hotels, event planning services, and automotive shops. In contrast,
 273 local terms are less likely to be used in the names of restaurants, shopping places, and
 274 beauty spas. This is understandable since we frequently see hotels (especially hotel chains)
 275 include city names as part of their names to indicate locations, such as *holiday inn charlotte*
 276 *center city*. Meanwhile, restaurant names may focus on describing food and cuisine styles
 277 to attract customers. Second, the uneven usage of local terms is highly consistent across the

278 seven metropolitan areas. This result suggests that the identified local term usage patterns
 279 are not specific to a particular region but can be generalized to other geographic areas.

280 To quantify the similarity and difference of local term usage in different POI types
 281 across geographic regions, we employ Jensen-Shannon divergence (JSD), which measures
 282 the similarity between two probability distributions. Equation 4 and 5 show the calculation
 283 of Jensen-Shannon divergence, where $KLD(P||Q)$ is the Kullback–Leibler divergence. The
 284 output of JSD is in $[0, 1]$, with 0 indicating that the two distributions are highly similar and
 285 1 suggesting that the two distributions are largely different.

$$286 \quad JSD(P||Q) = \frac{1}{2}KLD(P||M) + \frac{1}{2}KLD(Q||M) \quad (4)$$

$$287 \quad KLD(P||Q) = \sum_i P(i) \ln \frac{P(i)}{Q(i)} \quad (5)$$

288 JSD requires the input probabilities to sum to 1. To satisfy this criterion, we normalize
 289 the initial percentage values using Equation 6:

$$290 \quad NPr_i = \frac{Pr_i}{\sum_j Pr_j} \quad (6)$$

291 We then iterate through the seven metropolitan areas and calculate the pair-wise JSD, and
 292 finally calculate the average JSD value (there are in total 21 values). The obtained average
 293 JSD is 0.007, suggesting that the local term usage in different POI types are highly similar
 294 across geographic regions. The findings in this subsection can help us select suitable POI
 295 types in future for building computational models. For example, in the task of place name
 296 disambiguation, we may choose to focus on the POI names of certain types, such as *Hotels*
 297 and *Automotive* rather than *Restaurant* and *BeautySpas*, to extract more local terms which
 298 can then be associated with the related place names.

299 4.3 Analyzing POI name change with geographic distance

300 In this third analysis, we examine the change of POI names with geographic distance. Many
 301 phenomena follow Tobler’s First Law and show a distance decay effect. Do POI names,
 302 which reflect many underlying social and cultural processes, also show such an effect? Here,
 303 we look into the *collective similarity* of POI names between metropolitan areas, namely how
 304 the POI names in one area are overall similar or dissimilar to the POI names in another area.
 305 For instance, we may expect the Phoenix metropolitan area to have more similar POI names
 306 compared with the Las Vegas metropolitan area than with the Cleveland metropolitan area.

307 One major challenge for this analysis is how to measure the *collective similarity* of POI
 308 names between metropolitan areas. We propose two approaches to achieve this goal. The
 309 first and a straightforward approach is to group POI names in the same metropolitan area
 310 into a bag of words. This is similar to the TF-IDF approach discussed in our first analysis.
 311 However, we use only term frequency here, since TF-IDF artificially exaggerates the im-
 312 portance of local terms. While such an exaggeration is desired for local term extraction, it
 313 distorts the true frequencies of terms in POI names and therefore is not used in this analysis.
 314 We also do not remove the repeating POIs as we did in the first analysis. In short, we try to
 315 keep the POI names and term frequencies as they are in the real world in order to objectively
 316 model their change with geographic distance. The terms used in the POI names in each
 317 metropolitan area are combined together into a vector. We will refer to this approach as
 318 *count-based vector*. To formally define this approach, let r_1 and r_2 represent two geographic
 319 regions, and each region contains a set of POIs. We derive the vector for a geographic region

23:10 POI Names and Geographic Distance

320 by counting the frequencies of terms in POI names. A common vocabulary V is constructed
 321 based on all the terms of the POI names in a dataset. Thus, the names of POIs in the two
 322 regions, r_1 and r_2 , can be collectively represented as two vectors:

$$323 \quad \langle w_{11}, w_{12}, \dots, w_{1|V}| \rangle \quad (7)$$

$$324 \quad \langle w_{21}, w_{22}, \dots, w_{2|V}| \rangle \quad (8)$$

325 where $|V|$ represents the size of the vocabulary, and w_{ij} represents the count of term j used
 326 in the POI names in geographic region i .

327 While the count-based vector approach is straightforward, it does not capture the se-
 328 mantic similarity between terms. For example, the terms *kiku* and *sakana* are both used
 329 for the names of sushi restaurants in the dataset. The count-based vector will treat the two
 330 terms as completely different with a similarity of zero. However, the fact that these two
 331 terms both co-occur with *sushi* suggests there exists certain degree of similarity between
 332 them. *Word2vec* [34] is a model that has been found to effectively capture the semantic
 333 similarity between terms. It is a neural network model which learns *embeddings* (low di-
 334 mension vectors) for terms. In this work, we use the word2vec model to learn embeddings
 335 for metropolitan areas based on POI names. The embeddings are learned by predicting the
 336 terms used in POI names based on a given region (e.g., what terms are likely to be used for
 337 POI names if the region is *Phoenix, AZ*). The embeddings are condensed vectors, and the
 338 POI names in r_1 and r_2 can be represented as the two vectors below:

$$339 \quad \langle u_{11}, u_{12}, \dots, u_{1|d}| \rangle \quad (9)$$

$$340 \quad \langle u_{21}, u_{22}, \dots, u_{2|d}| \rangle \quad (10)$$

341 where d is the dimensionality of the embeddings, which can be decided empirically. In this
 342 analysis, we set $d = 300$ following the recommendation from the literature [34]. u_{ij} is a
 343 weight value learned from the POI dataset. The word2vec model aims to minimize the
 344 objective function in Equation 11:

$$345 \quad J = -\log\sigma(\mathbf{w}_o^T \mathbf{r}) - \sum_{k=1}^K \log\sigma(-\mathbf{w}_k^T \mathbf{r}) \quad (11)$$

346 where \mathbf{r} is the embedding of one geographic region, \mathbf{w}_o is the embedding of a term that is
 347 used for the POI names in region \mathbf{r} , while \mathbf{w}_k is the embedding of a term not used in region
 348 \mathbf{r} (which serves as negative samples). σ is a sigmoid function: $\sigma(x) = \frac{1}{1+e^{-x}}$.

349 With different geographic regions represented as vectors in the same dimension, cosine
 350 similarity can be employed to measure the similarity of two vectors (Equation 12). $s(r_1, r_2)$
 351 is then used as the collective similarity between regions r_1 and r_2 .

$$352 \quad s(r_1, r_2) = \frac{\sum_{j=1}^d w_{1j} w_{2j}}{\sqrt{\sum_{j=1}^d w_{1j}^2} \sqrt{\sum_{j=1}^d w_{2j}^2}} \quad (12)$$

353 We apply both the count-based approach and word2vec to the Yelp POI dataset to
 354 derive vectors for the seven metropolitan areas. The center point of each metropolitan
 355 area is derived by averaging the location coordinates of the POIs in that area. We then
 356 employ Vincenty's formulae [48], which is based on the assumption of an oblate spheroid,
 357 to calculate the distance between two metropolitan areas. We then perform both Pearson's
 358 and Spearman's correlation to examine the relation between the collective similarity of
 359 POI names and the geographic distance of the corresponding metropolitan areas. Table 1

■ **Table 1** Pearson and Spearman correlation coefficients between the collective similarity of POI names and geographic distance.

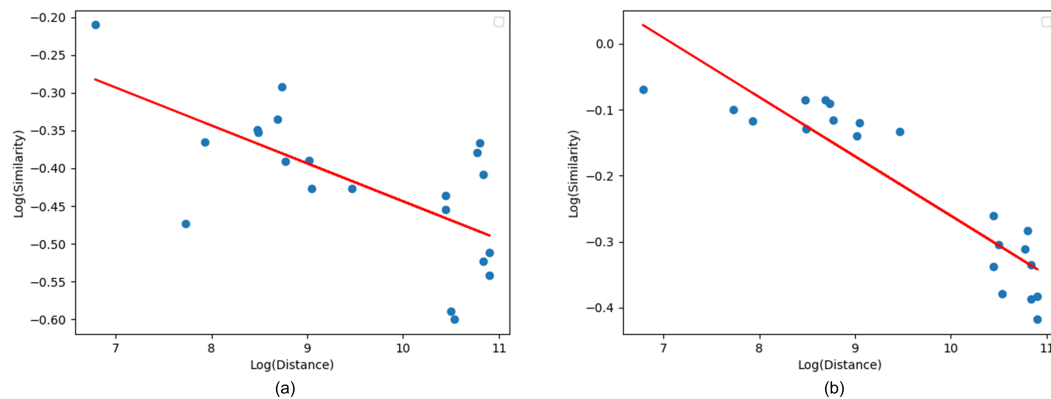
	Count-based vector	word2vec
Pearson	-0.612 (p<0.01)	-0.963 (p<0.001)
Spearman	-0.626 (p<0.01)	-0.917 (p<0.001)

360 shows the correlation results. Overall, the collective similarity of POI names negatively and
 361 significantly correlates with geographic distance based on the four correlation coefficients
 362 in Table 1, which suggests that POI names indeed *gradually* become less similar with the
 363 increase of geographic distance. We emphasize *gradually* here because either no change
 364 or abrupt change can lead to no correlation between POI name similarity and geographic
 365 distance. It is often natural to assume that place names at different locations are of course
 366 different, but our experiment result suggests that place names are not randomly different
 367 but follows a distance decay pattern. The statistical significance of the result is especially
 368 exciting given the fact that we have only 21 data points (21 region pairs from the seven
 369 metropolitan areas) for this correlation analysis. Such a result suggests that there is indeed a
 370 clear negative relation between POI name similarity and distance. In addition, it seems that
 371 word2vec better captures the POI name changes with geographic distance, as demonstrated
 372 by the higher correlation coefficients and stronger significances.

373 To further quantify the distance decay effect, we use a model $s = A * \frac{1}{d^\beta}$ to fit our data.
 374 We first transform it into its logarithmic form:

$$375 \quad \log s = A + \beta * \log d \quad (13)$$

376 where s is the collective similarity of POI names between two metropolitan areas, A is a
 377 constant, β is the slope, and d is the geographic distance between them. We fit a linear
 regression model based on the logged values. Figure 6 shows the result. In the count-

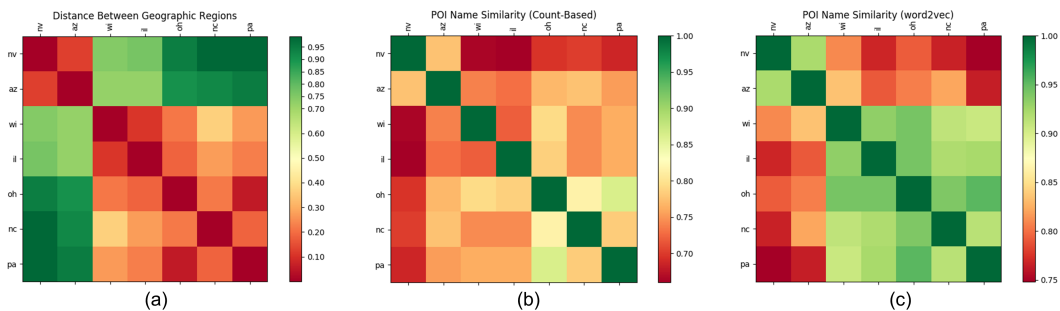


■ **Figure 6** Fitting the collective similarity of POI names with geographic distance: (a) count-based vector; (b) word2vec.

378 based vector approach, we obtained an R-squared value 0.434 and a slope of -0.050 . Using
 379 word2vec, we obtained a R-squared value 0.828 and a slope of -0.090 . More credibility
 380 can be given to the result from word2vec since it better captures the semantic similarity
 381 between terms in POI names. A slope of -0.090 indicates there is a clear distance decay
 382 effect with the increase of geographic distance. Besides, it is interesting to see how the data
 383

384 points clearly fall in two groups in Figure 6(b), which is consistent with their geographic
 385 distributions shown in Figure 2 (a group of city pairs has closer geographic distances, while
 386 the other group of city pairs has farther geographic distances). It would be interesting to
 387 examine the POI names in more metropolitan areas to see if their POI names also follow
 388 the general trend along the red line in Figure 6(b).

389 To further examine the result difference between the count-based vector and word2vec,
 390 Figure 7 shows the matrices of the geographic distances and the collective similarities ob-
 391 tained using the two approaches. It can be seen that the similarity pattern obtained using
 392 word2vec in sub figure (c) is closer to the distance pattern in sub figure (a) compared with
 393 the pattern from the count-based vector in sub figure (b). This result is consistent with the
 394 distance decay pattern observed in Figure 6.



■ **Figure 7** (a) The geographic distances between the seven metropolitan areas; (b) collective similarities based on count-based vector; (c) collective similarities based on word2vec.

395 5 Conclusions and future work

396 Place names are texts given by people to natural or man-made geographic features. The act
 397 of assigning a name to space signifies the important moment of space and human experience
 398 integration, and further enhances the social construct of *place*. Place names, as *city-text*,
 399 reveal a considerable amount of information about the culture, lifestyle, community, and
 400 many other aspects of a city. While place names have long intrigued geographers, existing
 401 research often focuses on case-by-case qualitative descriptions related to the etymology or
 402 taxonomy of place names, or only considers place names as identifiers without analyzing
 403 their term usage and their relations with geographic distances.

404 This paper presents an empirical study on place names and their change with geographic
 405 distance. This study is based on an open dataset from Yelp, and examines more than
 406 110,000 POIs, such as restaurants, hotels, and local services, in seven metropolitan areas
 407 in the United States. We perform an exploratory analysis on the frequencies of terms
 408 used in POI names, and find the term usage follows Zipf's law. We further conduct three
 409 analyses focusing on *place names*, *place types*, and *spatial footprints* respectively. We adapt
 410 the technique of term frequency and inverse document frequency in geographic context to
 411 identify local terms, and examine the term usage in the POI names in different types of
 412 POIs. We find an uneven usage of local terms across POI types (e.g., auto repairs are more
 413 likely to use local terms than restaurants), and such a usage pattern is highly consistent
 414 across different geographic regions. Finally, we test two approaches, count-based vector and
 415 word2vec, to model the collective similarity of POI names in different regions, and find a
 416 distance decay effect in the collective similarity of POI names.

417 This work is only a first step towards quantitatively and systematically examining place
 418 names and their relations with geographic distances. A number of topics can be explored in
 419 the near future. First, all the analyses are conducted based on the seven metropolitan areas
 420 available in the Yelp dataset. While a large number of POI names are examined, it would
 421 be interesting to apply the analyses to more metropolitan areas in other regions (e.g., north
 422 west and mid-south) as well as within local regions to further test the findings from this
 423 work. Second, we have so far used whole terms for the analyses, and it would be interesting
 424 to examine the parts or chunks of a term for measuring the collective similarity of place
 425 names. For example, the place names, *Wauwatosa* in Wisconsin, *Wawatasso* in Minnesota,
 426 and *Wahwahtaysee* in Michigan, share similar chunks, and may have higher similarity values
 427 when a chunk-based approach is used. Third, future research can be conducted on how to
 428 integrate the information extracted from place names with existing computational models
 429 for tasks such as place name disambiguation. While Wikipedia articles and other datasets
 430 have been frequently used for training place-based models, there are situations when we have
 431 only short Wikipedia descriptions or no description for places. Local information extracted
 432 from place names can serve as additional resources to improve existing models.

433 ——— References ———

- 434 1 Derek H Alderman. A street fit for a King: Naming places and commemoration in the
 435 American South. *The Professional Geographer*, 52(4):672–684, 2000.
- 436 2 Derek H Alderman. Street names as memorial arenas: The reputational politics of com-
 437 memorating Martin Luther King in a Georgia county. *Historical Geography*, 30:99–120,
 438 2002.
- 439 3 Derek H Alderman. Place, naming and the interpretation of cultural landscapes. *Heritage*
 440 *and Identity*, edited by Brian Graham and Peter Howard, pages 195–213, 2016.
- 441 4 Einat Amitay, Nadav Har’El, Ron Sivan, and Aya Soffer. Web-a-where: geotagging web
 442 content. In *Proceedings of the 27th annual international ACM SIGIR conference on Re-*
 443 *search and development in information retrieval*, pages 273–280. ACM, 2004.
- 444 5 Maoz Azaryahu. Street names and political identity: the case of East Berlin. *Journal of*
 445 *Contemporary History*, 21(4):581–604, 1986.
- 446 6 Maoz Azaryahu. Renaming the past: Changes in "city text" in Germany and Austria,
 447 1945-1947. *History and Memory*, 2(2):32–53, 1990.
- 448 7 Maoz Azaryahu. The power of commemorative street names. *Environment and Planning*
 449 *D: Society and Space*, 14(3):311–330, 1996.
- 450 8 Daniel L Baggio. *The dawn of a new Iraq: the story Americans almost missed*. US Army
 451 War College, 2006.
- 452 9 Davide Buscaldi and Paulo Rosso. A conceptual density-based approach for the disambiguation
 453 of toponyms. *International Journal of Geographical Information Science*, 22(3):301–
 454 313, 2008.
- 455 10 Paul Carter and Lawrie McKenzie. *The road to Botany Bay: an essay in spatial history*.
 456 Faber & Faber London, 1987.
- 457 11 James A Cheshire and Paul A Longley. Identifying spatial concentrations of surnames.
 458 *International Journal of Geographical Information Science*, 26(2):309–325, 2012.
- 459 12 Judith Gelernter and Nikolai Mushegian. Geo-parsing messages from microtext. *Trans-*
 460 *actions in GIS*, 15(6):753–773, 2011.
- 461 13 Michael F Goodchild and Linda L Hill. Introduction to digital gazetteer research. *Inter-*
 462 *national Journal of Geographical Information Science*, 22(10):1039–1044, 2008.
- 463 14 Milan Gritta, Mohammad Taher Pilehvar, Nut Limsopatham, and Nigel Collier. What’s
 464 missing in geographical parsing? *Language Resources and Evaluation*, pages 1–21, 2017.

- 465 15 Christian Grothe and Jochen Schaab. Automated footprint generation from geotags with
 466 kernel density estimation and support vector machines. *Spatial Cognition & Computation*,
 467 9(3):195–211, 2009.
- 468 16 Linda L Hill. Core elements of digital gazetteers: placenames, categories, and footprints.
 469 In *International Conference on Theory and Practice of Digital Libraries*, pages 280–290.
 470 Springer, 2000.
- 471 17 Yingjie Hu, Krzysztof Janowicz, and Sathya Prasad. Improving Wikipedia-based place
 472 name disambiguation in short texts using structured data from DBpedia. In *Proceedings*
 473 *of the 8th workshop on geographic information retrieval*, pages 1–8. ACM, 2014.
- 474 18 Suradej Intagorn and Kristina Lerman. Learning boundaries of vague places from noisy
 475 annotations. In *Proceedings of the 19th ACM SIGSPATIAL international conference on*
 476 *advances in geographic information systems*, pages 425–428. ACM, 2011.
- 477 19 Krzysztof Janowicz and Carsten Keßler. The role of ontology in improving gazetteer in-
 478 teraction. *International Journal of Geographical Information Science*, 22(10):1129–1157,
 479 2008.
- 480 20 Christopher B Jones and Ross S Purves. Geographical information retrieval. *International*
 481 *Journal of Geographical Information Science*, 22(3):219–228, 2008.
- 482 21 Christopher B Jones, Ross S Purves, Paul D Clough, and Hideo Joho. Modelling vague
 483 places with knowledge from the Web. *International Journal of Geographical Information*
 484 *Science*, 22(10):1045–1065, 2008.
- 485 22 Yiting Ju, Benjamin Adams, Krzysztof Janowicz, Yingjie Hu, Bo Yan, and Grant McK-
 486 enzie. Things and strings: improving place name disambiguation from short texts by
 487 combining entity co-occurrence with topic modeling. In *European Knowledge Acquisition*
 488 *Workshop*, pages 353–367. Springer, 2016.
- 489 23 Morteza Karimzadeh, Wenyi Huang, Siddhartha Banerjee, Jan Oliver Wallgrün, Frank
 490 Hardisty, Scott Pezanowski, Prasenjit Mitra, and Alan M MacEachren. GeoTxt: a web
 491 API to leverage place references in text. In *Proceedings of the 7th workshop on geographic*
 492 *information retrieval*, pages 72–73. ACM, 2013.
- 493 24 Robin A Kearns and J Ross Barnett. To boldly go? Place, metaphor, and the marketing of
 494 Auckland’s Starship Hospital. *Environment and planning D: Society and space*, 17(2):201–
 495 226, 1999.
- 496 25 Robin A Kearns and Lawrence D Berg. Proclaiming place: Towards a geography of place
 497 name pronunciation. *Social & Cultural Geography*, 3(3):283–302, 2002.
- 498 26 Carsten Keßler, Patrick Maué, Jan Heuer, and Thomas Bartoschek. Bottom-up gazetteers:
 499 Learning from the implicit semantics of geotags. *GeoSpatial semantics*, pages 83–102, 2009.
- 500 27 Jochen L Leidner. *Toponym resolution in text: Annotation, evaluation and applications of*
 501 *spatial grounding of place names*. Universal-Publishers, 2008.
- 502 28 Linna Li and Michael F Goodchild. Constructing places from spatial footprints. In *Proceed-*
 503 *ings of the 1st ACM SIGSPATIAL international workshop on crowdsourced and volunteered*
 504 *geographic information*, pages 15–21. ACM, 2012.
- 505 29 Michael D Lieberman, Hanan Samet, and Jagan Sankaranarayanan. Geotagging with
 506 local lexicons to build indexes for textually-specified spatial data. In *2010 IEEE 26th*
 507 *International Conference on Data Engineering (ICDE)*, pages 201–212. IEEE, 2010.
- 508 30 Duncan Light. Street names in bucharest, 1990–1997: exploring the modern historical
 509 geographies of post-socialist change. *Journal of Historical Geography*, 30(1):154–172, 2004.
- 510 31 Paul A Longley, James A Cheshire, and Pablo Mateos. Creating a regional geography of
 511 Britain through the spatial analysis of surnames. *Geoforum*, 42(4):506–516, 2011.
- 512 32 Christopher D Manning and Hinrich Schütze. *Foundations of statistical natural language*
 513 *processing*. MIT press, 1999.

- 514 **33** Grant McKenzie, Krzysztof Janowicz, Song Gao, Jiue-An Yang, and Yingjie Hu. POI pulse:
515 A multi-granular, semantic signature-based information observatory for the interactive
516 visualization of big geosocial data. *Cartographica: The International Journal for Geographic
517 Information and Geovisualization*, 50(2):71–85, 2015.
- 518 **34** Tomas Mikolov, Ilya Sutskever, Kai Chen, Greg S Corrado, and Jeff Dean. Distributed
519 representations of words and phrases and their compositionality. In *Advances in neural
520 information processing systems*, pages 3111–3119, 2013.
- 521 **35** Catherine Nash. Irish placenames: Post-colonial locations. *Transactions of the Institute of
522 British Geographers*, 24(4):457–480, 1999.
- 523 **36** Tessio Novack, Robin Peters, and Alexander Zipf. Graph-based strategies for matching
524 points-of-interests from different vgi sources. In *AGILE 2017*, pages 1–6, 2017.
- 525 **37** Simon Overell and Stefan Ruger. Using co-occurrence models for placename disambiguation.
526 *International Journal of Geographical Information Science*, 22(3):265–287, 2008.
- 527 **38** Kari Palonen. *Reading street names politically*. na, 1993.
- 528 **39** Pauliina Raento and William A Douglass. The naming of gaming. *Names*, 49(1):1–35,
529 2001.
- 530 **40** Reuben Rose-Redwood, Derek Alderman, and Maoz Azaryahu. Geographies of toponymic
531 inscription: new directions in critical place-name studies. *Progress in Human Geography*,
532 34(4):453–470, 2010.
- 533 **41** Reuben S Rose-Redwood. From number to name: symbolic capital, places of memory and
534 the politics of street renaming in New York City. *Social & Cultural Geography*, 9(4):431–
535 452, 2008.
- 536 **42** Reuben S Rose-Redwood. "sixth avenue is now a memory": Regimes of spatial inscription
537 and the performative limits of the official city-text. *Political Geography*, 27(8):875–894,
538 2008.
- 539 **43** Joao Santos, Ivo Anastacio, and Bruno Martins. Using machine learning methods for
540 disambiguating place references in textual documents. *GeoJournal*, 80(3):375–392, 2015.
- 541 **44** Rui Santos, Patricia Murrieta-Flores, Pavel Calado, and Bruno Martins. Toponym match-
542 ing through deep neural networks. *International Journal of Geographical Information Sci-
543 ence*, 32(2):324–348, 2018.
- 544 **45** Yi-Fu Tuan. *Space and place: The perspective of experience*. U of Minnesota Press, 1977.
- 545 **46** Florian A Twaroch and Christopher B Jones. A web platform for the evaluation of ver-
546 nacular place names in automatically constructed gazetteers. In *Proceedings of the 6th
547 Workshop on Geographic Information Retrieval*, page 14. ACM, 2010.
- 548 **47** Maria Vasardani, Stephan Winter, and Kai-Florian Richter. Locating place names from
549 place descriptions. *International Journal of Geographical Information Science*, 27(12):2509–
550 2532, 2013.
- 551 **48** Thaddeus Vincenty. Direct and inverse solutions of geodesics on the ellipsoid with applic-
552 ation of nested equations. *Survey review*, 23(176):88–93, 1975.
- 553 **49** Jan Oliver Wallgrun, Morteza Karimzadeh, Alan M MacEachren, and Scott Pezanowski.
554 GeoCorpora: building a corpus to test and train microblog geoparsers. *International
555 Journal of Geographical Information Science*, 32(1):1–29, 2018.
- 556 **50** John Kirtland Wright. The study of place names recent work and some possibilities. *Geo-
557 graphical Review*, 19(1):140–144, 1929.
- 558 **51** Bo Yan, Krzysztof Janowicz, Gengchen Mai, and Song Gao. From ITDL to Place2Vec–
559 Reasoning About Place Type Similarity and Relatedness by Learning Embeddings From
560 Augmented Spatial Contexts. *Proceedings of 2017 ACM SIGSPATIAL Conference*, 17:7–10,
561 2017.
- 562 **52** Wilbur Zelinsky. Along the frontiers of name geography. *The Professional Geographer*,
563 49(4):465–466, 1997.