

# Article Identifying urban neighborhood names through user-contributed online property listings

Grant McKenzie<sup>1</sup>, Zheng Liu<sup>2</sup>, Yingjie Hu<sup>3</sup>, and Myeong Lee<sup>2</sup>

- <sup>1</sup> McGill University, Montréal, Canada;
- <sup>2</sup> University of Maryland, College Park, USA;
- <sup>3</sup> University at Buffalo, Buffalo, USA
- \* Correspondence: grant.mckenzie@mcgill.ca

Version September 22, 2018 submitted to ISPRS Int. J. Geo-Inf.

- Abstract: Neighborhoods are vaguely defined, localized regions that share similar characteristics.
- <sup>2</sup> They are most often defined, delineated, and named by the citizens that inhabit them rather than
- <sup>3</sup> municipal government or commercial agencies. The names of these neighborhoods play an important
- role as a basis for community and sociodemographic identity, geographic communication, and
- <sup>5</sup> historical context. In this work we take a data-driven approach to identifying neighborhood names
- <sup>6</sup> based on the geospatial properties of user-contributed rental listings. Through a random forest
- ensemble learning model applied to a set of spatial statistics for all n-grams in listing descriptions,
- \* we show that neighborhood names can be uniquely identified within urban settings. We train a
- model based on data from Washington, DC and test it on listings in Seattle, WA and Montréal, QC.
- <sup>10</sup> The results indicate that a model trained on housing data from one city can successfully identify
- neighborhood names in another. In addition, our approach identifies less common neighborhood
- names and suggestions alternative or potentially new names in each city. These findings represent a
- <sup>13</sup> first step in the process of urban neighborhood identification and delineation.
- Keywords: neighborhood; neighborhood name; random forest; spatial statistics; housing; craigslist

# 15 PRE-PRINT

# 16 1. Introduction

In 2014, Google published a neighborhood map of Brooklyn, the most populous borough in 17 New York City, a seemingly harmless step in providing its users with useful geographic boundary 18 information. The backlash was swift. Residents of Brooklyn responded angrily, many stating that a 19 commercial company such as Google had no right to label and define boundaries within their city [1]. 20 This was not a lone incident [2], as many mapping agencies, both government and commercial, have 21 come to realize that regional boundaries and names are a contentious issue. Google and others are 22 frequently placed in the difficult situation of publishing hard boundaries and definitive names for 23 regions that are in dispute or poorly defined [3,4], often applying names to parts of the city that few 24 residents have even heard before [5]. This poses a problem as the names assigned to neighborhoods are 25 important for understanding one's identity and role within an urban setting. Names provide a bond 26 between a citizen and a place [6]. In many cases neighborhood names are much more than just a set of 27 characters, they have a history that is situated in religious beliefs [7], gender identity [8], and/or race [9]. 28 Neighborhood names evolve over time and are given meaning by the neighborhood's inhabitants. 29 Applying a top-down approach to naming neighborhoods, a practice often done by municipalities and 30 commercial agencies, can produces unforeseen, even anger-inducing, results. 31

Historically, neighborhood identification has also been predominantly driven through financial
 incentives. The term *redlining*, which describes the process of raising service prices or denying loans

in selective neighborhood and communities based on demographics such as race, was coined in the 34 1960s [10] and is one of the foundational examples of neighborhood delineation driven by financial 35 interests. In many ways, the neighborhood boundaries of many U.S. cities today are at least a partial 36 result of these practices. Real estate companies still rely on neighborhood boundaries for comparable 37 pricing [11] and being associated with a neighborhood name can significantly impact one's social 38 capital [12] as well as mortgage rate [13]. Today, web-based real estate platforms such as Zillow, Redfin, 39 and Trulia each curates their own neighborhood dataset [14]. These companies realize the immense 40 value of these boundaries and names [15] and actively invest in promoting their brand's datasets.<sup>1</sup> While commercial mapping companies and real estate platforms engage in the complex process 42 of geographically splitting up a city into neighborhoods and labeling those regions, the inhabitants 43 and citizens themselves often have their own understanding of the region in which they live. Their 44 historically-rooted understanding of a neighborhood can sometimes be at odds with the neighborhood 45 identification methods employed by these commercial entities. The urban historian, Lewis Mumford 46 stated that "Neighborhoods...exist wherever human beings congregate, in permanent family dwellings; 47 and many of the functions of the city tend to be distributed naturally-that is, without any theoretical 48 preoccupation or political direction" [16]. That is to say that neighborhoods differ from other regional 49 boundaries (e.g., city, census tract) in that they are constructed from the *bottom-up* by citizens, rather 50 than *top-down* by governments or commercial entities. Any attempt to interfere with this bottom-up 51 approach is met with resentment from residents of the neighborhoods, as evident by Google's Brooklyn 52 neighborhood map. In fact, one of the goals of public participatory GIS has been to enable citizens to 53 construct, identify, and contribute to their communities and neighborhood [17,18], thus defining the 54 regions themselves. 55 Today, information is being generated and publicly disseminated online by everyday citizens at 56 an alarming rate. While governments and industry partners have increased their involvement in public

participatory GIS and engagement platforms,<sup>2</sup> the vast majority of content is being contributed through 58 social media applications, personal websites, and other sharing platforms, many of which include 59 location information. Online classified advertisements are an excellent example of this recent increase 60 in user-generated content. People post advertisements for everything from local services to previously 61 used products, and most notably, rental properties. *Craigslist* is by far the most popular online website 62 for listing and finding rental properties in the United States, Canada, and many other countries<sup>3</sup> 63 and is therefore a rich source of information for understanding regions within a city. As inhabitants, 64 property owners, or local rental agencies post listings for rental properties on such a platform, they 65 geotag the post (either through geographic coordinates or local address), and provide rich descriptive 66 textual content related to the property. Much of this content includes standard information related 67 to the property such as square footage, number of bedrooms, etc., but other information is related to 68 the geographic location of the listing, namely nearby restaurants, public transit, grocery stores, etc. 69 Neighborhood names are also frequently included in rental listing descriptions. Those posting the 70 rental properties realize that by listing the neighborhood name(s) in which the property exists, they 71 are effectively situating their property within a potential renter's existing idea and understanding of 72 the region. While the motivation and biases surrounding which neighborhoods are included in the 73 textual descriptions of a listing are important (will be discussed in Section 6.2), these data offer a novel 74 opportunity to understand how citizens, property owners, and local real estate companies view their 75 urban setting, and label and differentiate the neighborhoods that comprise the city. 76 Given our interest in both identifying and delineating neighborhoods, this work tackles the 77

preliminary, but essential step of extracting and identifying neighborhood names. The specific

78

<sup>&</sup>lt;sup>1</sup> Zillow for example freely offer access to their neighborhood boundaries and real estate APIs.

<sup>&</sup>lt;sup>2</sup> See ArcGIS Hub and Google Maps Contributions, for example.

<sup>&</sup>lt;sup>3</sup> Over 50 billion classified page views per month. Source: http://web.archive.org/web/20161114220514/http://www.craigslist.org/about/factsheet

<sup>79</sup> contributions of this work are outlined in the five research questions (RQ) below. Each builds on the
 <sup>80</sup> findings of the previous question and direct references to these RQs can be found in the manuscript.

RQ1 Can neighborhood names be identified from natural language text within housing rental listings?
 Specifically, can spatially descriptive measures of geo-tagged n-grams be used to separate

neighborhood names from other terms? A set of spatial statistical measures are calculated for all

n-grams<sup>4</sup> in a set of listings and used to identify neighborhoods names.

RQ2 Does an ensemble learning approach based on spatial distribution measures more accurately
 identify neighborhood names than the spatial distribution measures alone? Given spatial
 statistics for each n-gram in a set of listings, we show that combining these in a random forest

model, produces higher accuracy than individual measures alone.

<sup>89</sup> RQ3 Can an identification model trained on a known set of neighborhood names be used to identify

- uncommon neighborhood names or previously unidentified neighborhoods? Training a random
   forest model on spatial statistics of common neighborhood names within a city, we demonstrate
- that lesser known neighborhood names can be identified. In some cases, alternative names or

<sup>93</sup> other descriptive terms are proposed through the use of such a model.

RQ4 Can a neighborhood name identification model trained on data from one city be used to identify
 neighborhood names in a different city? A random forest model constructed from neighborhood
 names in Washington, DC is used in the identification of neighborhood names in Seattle, WA
 and Montréal, QC.

<sup>98</sup> RQ5 What are the biases associated with neighborhood names mentioned in rental property listings?

<sup>99</sup> Lastly, we report on the spatial distribution biases associated with Craigslist rental listings in

100 Washington, DC.

The remainder of this manuscript is organized as follows. Previous research related to this topic is discussed in Section 2 and an overview of the data is provided in Section 3. The spatial statistics and random forest methods are introduced in Section 4 including measures of accuracy. Section 5 presents the results of this work which are then discussed in Section 6. Finally, conclusions, limitations, and future work are the subjects of Section 7.

# 106 2. Related Work

Defining neighborhoods has been the subject of numerous research projects spanning many 107 different domains. Understanding how neighborhoods are defined as well as identifying characteristics 108 that distinguish one neighborhood from another has a long history within geography, regional 1 0 9 science, sociology and social psychology (see [19–21] for an overview). Many previous studies in the 110 social sciences have contrasted citizen-defined neighborhoods to regions defined by government or 111 commercial entities. Coulton et al. [22] provide an example of this type of research, having asked 112 residents of a city to draw boundaries on a map, defining their version of neighborhoods within a city. 113 This process inevitably results in some overlap between neighborhood names and boundaries, but also 114 quite a few significant differences. These citizen-defined boundaries are then often compared to census 115 or other government designated areas [23,24]. An outcome of these works is a clear need to better 116 understand what a neighborhood is and how it can be identified based on the people that inhabit it. 117

From a geographic perspective, a substantial amount of work has aimed at defining geographic *areas of interest*. While many researchers steer clear of the term 'neighborhood,' many of the methods employed, focus on delinated a sub-urban region for its surrounding components based on some characteristic or spatial property. Many of these rely on analyzing user-contributed point data accompanied by a names, categories, or descriptive tags. For instance, Schockaert and De Cock [25]

<sup>&</sup>lt;sup>4</sup> An *n-gram* is a sequence of *n* items (often words) identified in text. For example 'kitchen' is a uni-gram, 'small kitchen' is a bi-gram, etc.

identified the spatial footprints of neighborhoods from geotagged content while a number of
studies [26,27] identified areas of interest based on user-contributed photograph tags. Tags have
been used in the identification of vaguely defined regions as well. For instance social media tags and
text were use to differentiate Southern California from Northern California [28].

Recent work has focused on extracting *functional regions* based on human activities and category-labeled places of interest [29] while other work has identified *thematic regions* such as the bar district or shopping regions of the city based on the popularity of social media check-ins [30]. Though not explicitly labeled as neighborhoods, the characteristics and activities afforded by these regions often result in them being referred to colloquially as neighborhoods. The livehoods project [31] aimed to identify regions based on the similarities of geosocial check-in patterns in various cities around the United States. This project, however, did not involve naming the small livehood regions.

From a data source perspective, existing work has used housing posts to better understand, 1 34 explore, and in some cases, define neighborhoods [32,33]. Chisholm and Cohen [34] developed The 1 35 Neighborhood Project, a web map based on combining geocoded craigslist posts with neighborhood 136 names extracted from text in the posts. The neighborhood names themselves, however, were 1 37 determined by experts and user-contributed knowledge of the region. Hu et al. [35] used housing 1 38 advertisements as training data for a natural language processing (NLP) and geospatial clustering 139 framework that harvests local and colloquial place names in order to enrich existing gazetteers. Our 140 work further enhances this approach, combining measures from a range of statistical techniques to 141 specifically extract sub-urban regional names. Zhu et al. [36] explored the use of spatial statistics to 142 differentiate geographic feature types and disambiguate place names [37]. In these works they showed 143 that different feature types and place names exhibit different spatial patterns and it is through these 144 individual patterns that geographic features can be compared (e.g., mountain tops to beaches). 145

While a considerable amount of previous work has focused on neighborhood boundary 146 identification and delineation, far less work has focused on the extraction of neighborhood names. 147 Brindley et al. [38,39] took a data-driven approach to mapping urban neighborhoods, using postal 148 addresses to extract neighborhood names and boundaries. Specifically, the authors extracted commonly 149 found sub-district names from within postal addresses, and used a kernel density function to estimate 150 the geographic boundary. While similar to our work in their usage of publicly available geo-tagged 1 5 1 content, their approach did not combine various spatial statistics with natural language text for the 152 extraction of neighborhood names, nor did it produce a prediction model that could be learned from 153 one city and applied to another. 154

Place name extraction has been an important topic within geographic information retrieval 155 community for some time. Jones et al. [40] focused on the extraction of place names and vague 156 regions from natural language on websites while others were able to extract spatial relations from 157 natural language in web documents [41]. In that same thread, additional research has looked at the 158 identification of place names based on their context within descriptive documents [42]. Further work 159 has focused on disambiguation of terms within a geographic context. For example, Buscaldi and 160 Rosso [43] used term similarity and context to disambiguate place names from one another. The rise 161 of social media content has lead to new sources of geotagged content that has been used for named 162 geographic entity extraction [44,45]. Co-occurrence of place names and other geographic locations 163 within natural language text has been shown to correspond with close spatial proximity [46]. Still 1 64 other research has proposed machine learning approaches to identify and disambiguate places within 165 a document based on contextual terms [47,48]. The work presented in this manuscript continues with 166 this *leitmotif*, proposing a novel approach to identifying neighborhood names based on the spatial 167 168 distribution and content of rental property listings.

# 169 3. Data

Two sources of data are used in this work, namely rental property listings and curated lists of neighborhood names. Both sets of data were collected for three cities in North America. Further details on these data are described below.

# 173 3.1. Rental Property Listings

Rental property listings were accessed from the online classified advertisements platform 1 74 Craigslist.<sup>5</sup> Specifically, postings in the *apts/housing for rent* section of the subdomains for three 175 cities, Washington, DC; Seattle, WA; and Montréal, QC were accessed over a 6-month period starting in 176 September of 2017. These three cities were chosen based on the availability of content and geographic 177 locations (two different coasts of the United States and one bilingual city in Canada). The content 178 collected for each city consists of rental housing property listings such as the one shown in Figure 1. 179 At a minimum each listing contains geographic coordinates, a title and unstructured textual content 1 80 describing the rental property. 1 81



Figure 1. An example Craigslist rental listing in Washington, DC.

Table 1 presents an overview of the data collected for each of the cities. The first column, *Listings*, 1 82 reflects the total number of rental housing listings collected in and around each city over the course of 183 6 months. The Unique Locations column lists the number of unique rental housing listings for each city 1 84 after data cleaning. Cleaning involved removing duplicate entries and restricting posts to only those 185 listed with a unique pair of geographic coordinates. This had to be done due to the fact that many 186 posts were repeated for the exact same listing location but with slightly different titles and content 187 (presumably an advertising tactic to saturate the market). Those listings with no textual content were 188 removed. 189

<sup>&</sup>lt;sup>5</sup> http://craigslist.org

**Table 1.** Number of craigslist housing listings, unique housing locations, unique number of n-grams across all city listings, and cleaned unique n-grams.

City	Listings	Unique Locations	Unique n-grams	Cleaned n-grams
Washington, DC	60,167	13,307	1,294,747	3,612
Seattle, WA	68,058	17,795	1,053,297	5,554
Montréal, QC	10,425	4,836	571,223	2,914

#### 190 3.1.1. N-grams

All the textual content, including titles, for each listing in a city were combined into a corpus and the Natural Language Toolkit [49] was employed to tokenize words in the corpus and extract 1 92 all possible n-grams (to a maximum of 3 words). The total number of unique n-grams per city are 193 shown in Table 1. The frequency of occurrence within the corpus was calculated for each n-gram 1 94 and those with frequency values above 4 standard deviations from the mean were removed as well 195 as all n-grams that occurred less than 50 times within each city. Furthermore, all n-grams consisting 196 of less than 3 characters were removed. The removal of the exceptionally high frequency n-grams was done to reduce computation given that it is highly unlikely that the most frequent words are 198 neighborhood names. For example, the top five most frequent, greater than 2 character words in each 199 of the cities are *and*, *the*, *with*. Similarly, the removal of n-grams occurring less than 50 times was done 200 to ensure robustness in our neighborhood identification model and elicit legitimized neighborhood 201 names. Given the long tail distribution of n-gram frequencies, this latter step removed most of the 202 n-grams including single occurrence phrases such as included and storage, throughout painted, and for 203 rent around. 204

#### 205 3.1.2. Geotagged N-grams

Provided the reduced set of n-grams for each city, the original geo-tagged listings were revisited 206 and any n-grams found in the textual content of the listings were extracted and assigned the geographic 207 coordinates of the listing. This resulted in a large set of <latitude, longitude, n-gram> triples 208 for each city. These geo-tagged n-grams were intersected with the 1km buffered boundaries for each 209 city to remove all listings that were made outside of the city. The buffers were added to account for 210 listings that described neighborhoods on city borders (e.g., Takoma Park on the District of Columbia – 211 Maryland border). Figure 2 shows two maps of geo-tagged n-grams in Washington, DC, (2a) depicts 212 the clustering behavior of neighborhood names (three examples shown in this case) and (2b) shows a 213 sample of three generic housing-related terms. 214

### 215 3.2. Neighborhood Names & Boundaries

Since neighborhoods in the United States and Canada are neither federally nor 216 state/province-designated geographical units, there is no standard, agreed upon set of neighborhood 217 names and boundaries for each city. In many cases, neighborhood boundaries are arbitrarily defined 218 and there is little agreement between neighborhood data sources. Zillow, for example, provides a freely 219 available neighborhood boundaries dataset<sup>6</sup> for large urban areas in the United States that is heavily 220 based on property values. Platforms such as Google Maps also contain neighborhood boundaries for 221 most cities in the United States. However, Google considers this data proprietary and does not make it 222 available for use in third-party applications. There are numerous other sources of neighborhood or 223 2 24 functional region boundaries available for specific cities but few of these sources offer boundaries for more than a single urban location. Table 2 lists four sources of neighborhood names and boundaries 225 along with the number of neighborhood polygons available for each city. Notably, the number of 226

<sup>&</sup>lt;sup>6</sup> https://www.zillow.com/howto/api/neighborhood-boundaries.htm



(a) N-grams of three neighborhood names

(b) N-grams of three non-neighborhood names

**Figure 2.** N-grams mapped from rental property listings in Washington DC. (a) shows the clustering behavior of three neighborhood names while (b) visually depicts the lack of clustering for a sample of generic housing terms.

neighborhood names and polygons range substantially between data sources. Washington, DC, for
example, consists of 182 neighborhood boundaries according to *Zetashapes* compared to 46 listed on *DC.gov*.

**Table 2.** Neighborhood names and boundary sources including polygon counts for each city. The \* indicates that this source assigns many neighborhood names (comma delimited) to larger than average neighborhood regions. Note that Zillow and Zetashapes do not provide neighborhood names outside of the United States.

Source	Washington, DC	Seattle, WA	Montréal, QC
Wikipedia	129	134	73
Zillow	137	115	N/A
Zetashapes / Who's On First	182	124	N/A
City Government / AirBnB	46*	106	23
Common Neighborhoods	95	79	23

To build a training set for our machine learning model, we attempted to match each of the neighborhood names in each of the sources and exported those names that occurred in the majority of the sources. We label these our *Common Neighborhoods* and use them as the foundation on which to build the identification model.

# 234 4. Methodology

In this section we first give an overview of the various spatial statistics used to spatially describe the n-grams. This is followed by assessing the prediction power of each spatial statistic *predictor* in identifying neighborhood names and finally describing how the predictors are combined in a random forest ensemble learning model. Figure 3 depicts a flow chart of the process, with example data, from data-cleaning to random forest model.



**Figure 3.** A flow chart showing the process and example data for the methodology in this work. Note that the data is simplified/rounded for example purposes.

#### 240 4.1. Spatial Statistics

The fundamental assumption in this work is that different categories of words can be described by an array of statistics associated with the locations of their use. We hypothesize that neighborhood names exhibit unique spatial statistical patterns which can be used to specifically identify and extract these neighborhood names from other terms. With this goal in mind, we identified a few foundational spatial statistics that can be applied to representing point data in space. In total, 24 different spatial statistics measures, roughly grouped in to three categories, are used in describing each of the n-grams in our dataset. To be clear, we do not claim that this list of spatial statistics is exhaustive, but rather intend to show what is possible with a select set of measures.

#### 249 4.1.1. Spatial Dispersion

Nine measures of spatial dispersion were calculated for each n-gram in our datasets. *Standard Distance*, a single measure representing the dispersion of points around a mean centroid, was calculated along with average nearest neighbor and pairwise distance. We hypothesize that neighborhood names will be identified by this measure as neighborhood n-grams are likely to display a unique spatial dispersion pattern, different from most other non-geographic terms. Standard distance is shown in Equation 1 where *x* and *y* are individual point coordinates,  $\bar{X}$  and  $\bar{Y}$  are the mean centroid coordinates and *n* is the total number of geographic coordinates associated with the n-gram.

$$StandardDistance = \sqrt{\frac{\sum_{i=1}^{n} (x_i - \bar{X})^2}{n} + \frac{\sum_{i=1}^{n} (y_i - \bar{Y})^2}{n}}$$
(1)

<sup>257</sup> Within the category of average nearest neighbor (ANN), we calculated the mean and median for <sup>258</sup> each point's closest n-gram neighbor (NN1), second nearest (NN2), and third nearest (NN3) resulting <sup>259</sup> in six unique measures. Finally we computed the mean and median pairwise distance, or distance <sup>260</sup> between all pairs of points assigned to a single n-gram. ANN and Pairwise calculations were done <sup>261</sup> using the spatstat package in *R* [50]. Similarly to *Standard Distance*, we hypothesize that the average <sup>262</sup> spatial distance between the closest (2nd closest, and 3rd closest) n-grams that describe the same <sup>263</sup> neighborhood will be unique for neighborhoods, thus allowing us to include this measure in our <sup>264</sup> approach to neighborhood name identification.

#### <sup>265</sup> 4.1.2. Spatial Homogeneity

The spatial homogeneity of each geo-tagged n-gram was calculated through a binned approach to Ripley's L, or variance stabilized Ripley's K [51,52]. Ripley's L measures the degree of clustering across different spatial scales. Specifically, our approach split the resulting Ripley's L clustering function into ten 500m segments and averaged the range of clustering values for each n-gram within each segment. Figure 4 shows the binned Ripley's L approach for two n-grams in Washington, D.C., one a neighborhood name (Columbia Heights) and the other what should be an a-spatial term (Wood
Flooring). From a conceptual perspective, one might expect that most neighborhood names will show
a higher than expected degree of clustering around a certain distance mark. Higher than expected
clustering at a small distance might identify landmarks, while clustering at a large distance might be
useful for the identification of metro stations. Ripley's L allows us to assess clustering vs expected
clustering across these different distances. This approach of binning spatial homogeneity functions has
been employed successfully in differentiating point of interest types (e.g., Bars vs. Police Stations) [53].



**Figure 4.** Ripley's L function over 5 km for two n-grams, *Columbia Heights* (a neighborhood name) and *Wood Flooring*. The points show the averaged 'binned' values over 500 m.

In addition to the ten binned relative clustering values, the *kurtosis* and *skewness* measures for each Ripley's L function over 5km was recorded for each n-gram. The kurtosis and skewness provide overall measures of the Ripley's L function instead of a single measure based on binned distance.

#### 281 4.1.3. Convex Hull

The convex hull [54] is the smallest convex set (of listings in this case) that contain all listings. 282 Using the chull R package<sup>7</sup>, we computed the area of the convex hull for all geo-tagged n-grams in 283 our dataset as well as the density of the convex hull based on the number of listings in the set divided by 284 the area. These two measures offer a very different description of the property listings as they represent 285 the maximum area covered by all listings. Convex hull area simply assigns a numerical value for the 286 region covered by all listings. This measure is heavily impacted by outliers (e.g., random mention 287 of a neighborhood across town) as one occurrence can drastically alter the area of the convex hull. 288 Conceptually, density of points within the convex hull is useful for comparing n-grams as we would 289 expect to find a higher than average density of points within a region identified as a neighborhood, 290 compared to an a-spatial term such as wood flooring. 291

292 4.1.4. Spatial Autocorrelation

As part of our initial exploratory analysis for this project, spatial autocorrelation was investigated as a meaningful spatial feature due to its potential relatedness to neighborhood names. This form of measurement, however, is substantially different from many of the other measures mentioned here as

<sup>7</sup> https://stat.ethz.ch/R-manual/R-devel/library/grDevices/html/chull.html

there is really no way to report spatial autocorrelation through a single value per n-gram. As with
other measures of correlation, this inference statistic involves interpreting the results through a number
of values, not least of which are *P*-values and *Z*-scores. Running Moran's I across our set of geo-tagged
n-grams we found the results inconclusive overall. At least half of the values for global Moran's I were
not of a high enough significance including what many would consider 'prototypical' neighborhood
names in Washington, DC, such as *Georgetown* and *Capital Hill*. For these reasons, we elected to leave
Moran's I after the exploratory phase of analysis and not use it in the final random forest models.

## 303 4.2. Data Setup

In setting up the data for input to a prediction and identification model, we calculated the above 304 statistics for each n-gram in our dataset. These values were then combined into a single data table, one 305 for each city with rows as n-grams and columns as statistical measures. From this point on we will 306 refer to the spatial statistic values as predictor variables. The n-grams in the common neighborhood names 307 dataset (see Section 3.2) were programmatically compared against all n-grams in the merged dataset 308 and matches were recorded. While in an ideal world this would have resulted in clean matches, a 309 number of matches were not made due to slight misspellings, abbreviations (e.g., Columbia Hts. for 310 Columbia Heights), and variations of n-grams that include the neighborhood names (e.g., to Dupont or 311 *Logan Circle from*). These neighborhoods were identified and matched manually by two researchers 312 and disagreements were resolved by a third person. Again, manual matching was only based on 313 the common neighborhood names, not all potential neighborhood names. As a result of this process, all 314 n-grams were given a binary value, either identified as neighborhood matches or not. 315

# 316 4.3. Individual Predictors

Having calculated spatial statistics values for each of our n-grams based on the methods described 317 in the previous sections, we now turn to **RQ1**, examining how well each individual statistic performs 318 at identifying a neighborhood from within the set of descriptive n-grams. All predictor variables were 319 normalized to between 0 and 1 based on their minimum and maximum values to allow for simpler 320 comparison. The Pearson's correlation matrix of all predictors and neighborhood matches is shown in 321 Appendix A1. A single star (\*) indicates p < 0.1, three stars (\*\*\*) indicates no significance, and all other 322 values are significant to p < 0.01. Notably there tends to be a negative correlation between the mean 323 and median nearest neighbor values and Neighborhood Match and a positive correlation to all binned 324 Ripley's L variables. 325

Each of the individual variables was then used to predict which of the n-grams was a neighborhood name and the accuracy of each prediction was recorded. The  $F_{score}$  (Equation 2), harmonic mean of precision and recall, was used to assess prediction power. Accuracy measures were recorded at 0.05 threshold increments, meaning the first time the model was run, any predictor variable value above (and including) 0.05 was considered a match and everything below was not. The threshold value that produced the best  $F_{score}$  for each predictor variable was identified. The best scores are reported in Section 5.1.

$$F_{score} = 2 \cdot \frac{\text{precision} \cdot \text{recall}}{\text{precision} + \text{recall}}$$
(2)

#### 333 4.4. Random Forest Ensemble Learning

In addressing **RQ2**, we now combine our predictor variables and take a machine learning approach using a *Random Forest* [55,56] ensemble learning method to identify neighborhood names within our n-gram dataset. Random forest models have proven quite successful in numerous other classification and regression tasks involving geographic components [57–59]. Random forest models are touted as being better suited to lower dimensional problems (as opposed to support vector machines for example) and correct for over-fitting which tends to happen when training decision trees. Random forest is a supervised learning algorithm that approaches the problem of learning through a combination of decision trees. Multiple decision trees are built and merged together to get a more accurate and stable prediction. This follows the idea that *en masse*, these individual learners will combine to form a learner better than its individual parts. In this work we used the R *randomForest* package ported from the Breiman et al.'s original Fortran code.<sup>8</sup>

345 4.4.1. Training & Testing

The first random forest model was trained with a randomly selected 50% of the n-grams in the 346 Washington, DC dataset (both neighborhood matches and non-neighborhood matches) and tested for 347 accuracy against the remaining 50% of the data. This combination of training and testing was done 348 100 times in order to produce robust measures for results, each time training on a different randomly 349 selected 50% of the Washington, DC data. When each model was trained, it was applied to the testing data in order to predict which n-grams were neighborhoods and which were not. This was done using 351 a probability method of prediction with the resulting probability for each n-gram bounded between 0 352 (not a neighborhood) and 1 (definitely a neighborhood). The F-scores (Equation 2), were recorded at 353 0.05 probability increments every time the prediction model was run and the probability threshold 354 that produced the best mean F-score was identified. By way of comparison, we also randomly selected n-grams as neighborhood matches in our dataset and trained a separate random forest model on these 356 data. The purpose of this was to provide a baseline on which our true neighborhood matching model 357 could be compared. The same number of n-grams were chosen, at random, so as to provide comparable 358 results 359

The purpose of this research is not only to show that spatial features can be used to identify 360 existing neighborhoods, but also can be used in the identification of less common and previously 361 unidentified neighborhoods (**RQ3**). To this end, the model probability threshold should be adjusted 362 to alter the precision as we want to identify those false positives that may actually be neighborhood 363 names but may either not have been found in our dataset, were not matched to a *common* neighborhood 364 name, or are more *colloquial*, or unofficial, neighborhood names. After computing the optimal threshold 365 value (based on F-score), we manually examined the false positives, those that the model identified as 366 neighborhoods, but were not considered matches to our *common neighborhood* list. Through manual 367 inspection we discovered a number of interesting false positives. Many were neighborhood names 368 that appeared in one or more of the curated neighborhood lists, but did not appear in enough sources 369 to be considered part of the *common neighborhood set*. Provided these new neighborhood matches, we 370 build a subsequent random forest model this time with the addition of those newly identified false 371 positive n-grams that are in fact neighborhood names. The resulting accuracy of both of these models 372 are reported in Section 5. 373

#### 374 4.4.2. Variable Importance

The random forest models described in the previous section were constructed with 500 tries 375 and 4 variables tried at each split. As a result of these splits, the model produced a ranking of the 376 predictor variables based on their contribution to the overall accuracy of the model. Figure 5 shows 377 the importance of these variables by way of the mean decrease in Gini index of node purity. What this 378 demonstrates is that some variables are more useful than others at splitting mixed label nodes in to 379 pure single class nodes. The higher the value, in this case, the more important the predictor variable 380 is to the model's accuracy. We see here that larger bin distances of Ripley's L are substantially more 381 important to the success of the model than the mean nearest neighbor measures, for example. To some 382 extent, this mirrors the ranking of individual predictor accuracy that is reported in Section 5.1. 383

<sup>&</sup>lt;sup>8</sup> https://cran.r-project.org/package=randomForest

	0	5 1	0	15
	L	Τ	Υ	Ì
Kurtosis	•			
Mean NN1	•••••			
Mean NN3	•••••			
Ripley 2	•			
Mean NN2	•			
Median NN1	•••••			
Skewness	••••			
Ripley 3	•			
Convex Area				
Ripley 4	· · · · · · · · · · · · · · · · · · ·			
Convex Density		•		
Ripley 1		•		
Median Pair		•		
Median NN3		•		
Median NN2		•••••		
Ripley 5		•		
Ripley 6	· · · · · · · · · · · · · · · · · · ·	•		
Ripley 7		: : :		
Count		•		
Standard Distance			• • • • • • • • • • • • • • • • • • • •	
Ripley 10			•••••	
Ripley 8			•	
Mean Pair			•	
Ripley 9				

Mean Decrease in Gini Index of Node Purity

**Figure 5.** Mean decrease in Gini Index of node purity due to shuffling of values within predictive variable.

#### 384 4.5. Evaluation

Having trained two random forest models based on housing rental n-grams from Washington, 385 DC, we next turn our attention to RQ4, namely evaluating the accuracy of such a model using data 386 from two other North American cities, Seattle, WA and Montréal, QC. As described in Section 4.2, 387 the predictor variables for each n-gram were merged into city-specific datasets and matched against 388 existing common neighborhood lists for their respective cities. Manual inspection and matching was 389 done as before, and those n-grams that matched neighborhood names were marked as matches while 390 all others were not. The random forest model trained on the Washington, DC data was then tested 391 against the geo-tagged Seattle and Montréal n-grams independently using the highest performing 392 probability threshold value from the Washington DC testing results. 393

# 394 5. Results

In this section we present the results of the methods presented in the previous section. Specifically
 we focus on the accuracy values of the individual predictors as well as the combined random forest
 model.

#### 398 5.1. Individual Predictors

The maximum  $F_{score}$  accuracy values for the individual predictor variables are shown in Table 3. On average, the accuracy of each predictor variable independently is not high. However, the binned Ripley's L variables produced the best predictive results with the 4500m bin (Ripley 9) producing the best  $F_{score}$  of 0.633 with a recall and precision of 0.724 and 0.562 respectively. These results demonstrate measure can perform reasonably well at differentiating neighborhoods from non-neighborhoods.

<sup>405</sup> Notably however, not all spatial statistics are useful for this endeavor independently. Next, we explore

combining these individual predictor variables with the purpose of improving neighborhood nameidentification.

**Table 3.** Max F-scores for individual predictor variables trained and tested on data from Washington,DC.

Measure	Max F-Score
Standard Distance	0.047
Count	0.083
Mean NN1	0.047
Mean NN2	0.047
Mean NN3	0.047
Med. NN1	0.050
Med. NN2	0.050
Med. NN3	0.050
Mean Pair	0.047
Med. Pair	0.047
Ripley 1	0.099
Ripley 2	0.279
Ripley 3	0.405
Ripley 4	0.500
Ripley 5	0.548
Ripley 6	0.570
Ripley 7	0.587
Ripley 8	0.624
Ripley 9	0.633
Ripley 10	0.624
Kurtosis	0.101
Skewness	0.047
Convex Area	0.047
Convex Density	0.144

#### 408 5.2. Ensemble Learning

The first step in matching common neighborhoods to n-grams (both programmatically and manually) resulted in 59 neighborhood names, out of 95, being identified in the 3,612 unique n-grams in Washington, D.C. Of these, 30 were direct matches, with 29 indirect, manually identified matches. There are a number of reasons why not all common neighborhood names were found in our dataset which will be discussed in Section 6.

The first random forest model was trained on the predictor variables of n-grams tagged as either 414 *common neighborhoods* or not. The resulting averaged  $F_{score}$  is shown in Table 4. This value is based on a 415 prediction probability threshold of 0.35. This is a high F-score given the noisiness of the user-generated 416 content on which the model was constructed. The recall value indicates how well the model did at 417 identifying known neighborhoods whereas the precision tells us how well the model did at identifying 418 neighborhoods n-grams as neighborhoods and non-neighborhood n-grams as such. As mentioned in 419 Section 4, these results allowed us to re-examine our dataset and uncover neighborhood names that 420 were not previously identified, i.e., those that did not appear in our common set but rather one of the 421 individual neighborhood sources such as Wikipedia. Through manual inspection, we increased the 422 number of neighborhood / n-gram matches in our dataset and trained a new random forest model on 423 the data. The results of this second random forest model are shown in the second row of Table 4. The 4 2 4  $F_{score}$  has improved as have both the precision and the recall with the largest increase occurring in the 425 recall value. 426

**Table 4.** F-score, precision, and recall values for two random forest models trained and tested on listings from Washington, D.C. Accuracy values for a model built on random assignments is also shown for comparison.

Model	F-Score	Precision	Recall
Common matched neighborhoods	0.807	0.845	0.777
Common + secondary matches	0.872	0.863	0.882
Randomly assigned matches	0.047	0.063	0.037

As a base-line we also included the  $F_{score}$  results of a random forest trained on randomly assigned matches (not necessarily neighborhood names). As expected, the results are considerably lower than the previous two models with an accuracy of roughly 0.05.

430 5.3. Identifying Neighborhoods in Other Cities

Equipped with the best performing random forest model trained and tested on the Washington, DC n-grams, we then tested it against our two other North American cities, as outlined in **RQ4**.

## 433 5.3.1. Predicting Seattle Neighborhoods

The first row of Table 5 shows the results of the random forest model trained on Washington, DC n-grams. This first model used the *common* Seattle neighborhoods as matches. As was reported in the previous section, the results of the first RF model prediction lead to an investigation of the precision of the model resulting in the identification of a number of neighborhoods that were not previously identified as such. This was rectified and the model was run again producing the values show in the second row of table.

**Table 5.** F-score, precision, and recall values for two random forest models trained on listings from Washington, DC and tested on listings from Seattle, WA (the first two rows). The last row shows the results of a model trained and tested on listings from Seattle, WA.

	F-Score	Precision	Recall
Common matched neighborhoods	0.671	0.625	0.724
Common + secondary matches	0.733	0.702	0.767
Trained on Seattle (common)	0.786	0.782	0.791

The third row of Table 5 presents the results of a random forest model trained on half of the Seattle data rather than the Washington, DC n-grams, and tested on the other half of the Seattle data. These results indicate that while the DC-trained RF models do perform well at predicting Seattle neighborhoods, a model trained on local data, still performs better.

5.3.2. Predicting Montréal Neighborhoods

In many ways, Seattle, WA is very similar to Washington, DC. Both are major metropolitan, 445 predominantly English speaking cities. Both host populations of roughly 700,000 and have similar 446 population densities, median age, and median income. To test the robustness of the DC-based random 447 forest model, we chose to test it against a very different city, namely Montréal, Quebec in Canada. 448 Montréal is a bilingual French/English speaking island city, boasting French as it's official language. 449 Montréal has a population of roughly 2 million (on island) residents. Craigslist rental housing listings 450 in Montréal are written in either French or English and often both. In addition to all of this, the city 4 5 1 has a historically unique rental market with the majority of leases beginning and ending on July 1 [60]. 452 Given the data collection dates, far fewer rental postings were accessed for the city compared to both 453 Washington, DC and Seattle, WA. These factors combined, this city offers a unique dataset on which to 4 5 4 test our model. 455

**Table 6.** F-score, precision, and recall values for two random forest models trained on listings from Washington, DC and tested on listings from Montréal, QC. The last row shows the results of a model trained and tested on listings from Montréal, QC.

	F-Score	Precision	Recall
Common matched neighborhoods	0.397	0.353	0.453
Common + secondary matches	0.483	0.412	0.583
Trained on Montréal (common)	0.655	0.559	0.792

As shown in Table 6, the first random forest model built from the DC n-grams produces an  $F_{score}$ of roughly 0.4. Upon examining the results of this model, additional non-common neighborhoods were identified and a second model was run resulting in a slightly higher F-score. While clearly not as high as the Seattle results, these values are still substantially higher than a model built on randomly matched n-grams. As was the case with Seattle, a model built on local Montréal data produced the best results with an F-score of 0.655 and notably a recall inline with that of Seattle's. A set of n-grams identified as neighborhoods by this model is presented in Appendix A2.

#### 463 6. Discussion

The results presented in this work offer evidence as to how neighborhoods can be identified by the 4 64 spatial distribution of rental housing advertisements. These findings demonstrate that identification of 465 a sample of common neighborhood names with spatial distribution patterns can be used to accurately 466 predict additional, less common neighborhood names within a given city. Furthermore, we find that an 467 array of spatial distribution measures from neighborhoods identified in one part of North America can 468 be used to train a machine learning model that can then be used to accurately identify neighborhoods 469 on another part of the continent. While rental housing data from local listings produces a more accurate 470 model, we find that this model can also span linguistic barriers, admittedly producing less accurate, 471 but quite significant, results. In this section we further delve into the nuanced results of using such a 472 machine learning approach and identify unique aspects and biases within the dataset. 473

#### 474 6.1. False Positives

The F-score values presented in the Tables 4-6 depict an overall view of the accuracy of the model, 475 but omit the nuances of the actual on-the-ground data and neighborhoods. Specifically some regions of 476 the city are better represented by the dataset than others and this is reflected in the analysis results. The 47 size, dominance, and popularity of a neighborhood all impact the probability of a neighborhood being 478 identified in the n-gram datasets. For example, many of the historic neighborhoods in Washington, 479 D.C. (e.g., Georgetown, Capital Hill, Brightwood) were clearly represented in the original data thus 480 resulting in high accuracy results. These prevalent neighborhoods then had a much larger impact 481 in contributing to the construction of the neighborhood identification model. This often meant that 48 smaller and less dominant neighborhoods, e.g., Tenleytown, were less likely to be identified through 483 the machine learning process and other, non-neighborhood regions were more likely to be identified. 4 84

**Table 7.** Examples of n-grams falsely identified as neighborhood names split by city (columns) and category (rows).

Category	Washington, DC	Seattle, WA	Montréal, QC
Landmarks	Capitol Building	Space Needle	Place Jacques-Cartier
Academic Inst.	Catholic University	University of Washington	McGill University
Streets	Wisconsin Ave.	Summit Ave.	Cavendish Blvd.
<b>Broader Regions</b>	National Mall	Waterfront	Saint-Laurent River
Transit Stations	Union Train Station	King Street Station	Jolicoeur Station
Companies	Yes, Organic	Amazon	Atwater Market
Misc.	blvd	concierge	du vieux

While the model performed well provided training data from within the city, there were an 485 expected set of false positives (see Table 7 for examples). Further examination of these false positives 486 allow us to categorize them into 6 relatively distinct groupings. Landmarks such as the Capitol Building or the White House were falsely identified as neighborhoods given the importance of these 488 landmarks within Washington, DC. Many housing rental listing specifically mentioned a proximity 489 to these landmarks thus resulting in spatial distribution measures similar to those of neighborhoods. 490 Similarity, some important streets, academic institutions and popular transit stations were labeled as 4 91 neighborhoods given their dominance within a region of the city. This reiterates the argument from the introduction of this paper that neighborhoods are simply regions with distinct characteristics that 4 93 are given a descriptive name by inhabitants and visitors. It therefore follows that many neighborhood 1 01 names come from important streets (e.g., George Ave), Transit stations (Union Station) and Universities 495 (Howard). While many of these n-grams identified as neighborhoods by our model were labeled as 496 false positives, there is an argument to be made that the n-grams do exist as neighborhood names. 497

Though many of these *false positives* can be explained given knowledge of the region, spatial dominance of a certain term, or prevalence of the geographic feature, a small portion of the false positives appeared to be non-spatially related. For example, terms such as *concierge* and *du vieux* appear to not be related to any geographic feature or place within a city and rather are n-grams within the data that happen to demonstrate spatial distribution patterns similar to neighborhoods. In addition to these, a number of real-estate company names were falsely identified as neighborhoods in our initial models given that many real estate companies are focused specifically on one region of a city. These *real estate company* related n-grams were removed early in the data cleaning process.

#### 506 6.1.1. Washington, DC

Washington, DC is a particularly interesting city, arguably representative of many east coast 507 U.S. cities, namely in the way that many populated regions run into one another. Washington, DC 508 itself is part of the larger Metro DC area which includes cities in the neighboring states of Virginia 5.09 and Maryland. Since rental housing listings were clipped to the buffered boundary of Washington, 510 DC, this meant that some neighborhoods were identified by the model that do not appear in the 511 common DC neighborhood set as they technically exist outside the district boundary. Examples of 512 such neighborhoods identified by our model are Alexandria and Arlington in Virginia and Silver Spring 513 and *College Park* in Maryland. 514

Within the district boundaries a number of neighborhoods were identified through the machine 515 learning model that did not originally exist in the common neighborhoods set for the district such 516 as Cleveland Park and University Heights, both labeled as neighborhoods on Wikipedia. Moreover, alternative or secondary names for neighborhoods were identified in the results, such as Georgia Ave, 518 a secondary name for *Petworth*, and *Howard*, the name of a University that has taken on a colloquial 519 reference to a sub-region within or overlapping the Shaw neighborhood. While many of the false 520 positives were smaller than a typical neighborhood area (e.g., Capitol Building), the ensemble learning 521 model also identified a number of larger regions, such as the National Mall, an important tourist attraction within Washington, DC, and the broader Northeast region of the district. Notably, Washington 523 addresses are divided into quadrants based on intercardinal directions. As stated previously, a few 524 major streets were identified, namely Wisconsin Ave., Connecticut Ave., and Rhode Island Ave., all 525 major thoroughfares leading from outside of the district to the city center. As demonstrated with 526 Georgia Avenue, many street names have taken on neighborhood-like statuses being used to describe 527 regions of similar socioeconomic status, demographics, or other characteristics.

#### 529 6.1.2. Seattle, WA

Further qualitative discussion of the n-gram neighborhood identification results in Seattle expose
some unique aspects of the city. As was the case in Washington, DC, investigation of false positives
exposed a number of neighborhood names that did exist as neighborhoods in one of the neighborhood

datasets (e.g., Wikipedia) but not in the common neighborhood set. Examples of these are Columbia 533 City, Lake Union, and Wallingford. Neighborhoods outside the Seattle city boundary such as Bothell 5 34 or Mercer Island were also identified as were neighborhoods such as Lincoln Park, a large park which has given rise to a new neighborhood name, and Alki Beach, a sub-neighborhood within West 536 Seattle along the waterfront. While popular streets, e.g., Summit and Anderson, were labeled as 537 neighborhoods, the biggest difference in false positives compared to Washington, DC, is an increase in 538 company/foundation names identified as neighborhoods. Amazon.com Inc, The Bill and Melinda 539 Gates Foundation, and Microsoft (outside of Seattle) were each clearly identified as neighborhoods and the first Starbucks location (in Pike Place Market) was initially identified as a neighborhood when 541 the model was built on local training data. 542

## 543 6.1.3. Montréal, QC

Examination of the n-gram results in Montréal produced some interesting insight into how a machine learning model such as this is actually quite language-independent, at least as it relates to English and French rental listings. Importantly, though a single rental listing may contain both French text and English translation, the neighborhood names in Montréal are either in French *or* in English, not both, at least according to the reference datasets we employed. This means that each neighborhood does not have two names (one in each language) and implies that a model does not have to be adjusted for sparsity in the labels, but rather can be run as is.

As in the previous two cities, non-common neighborhoods were identified through the model such as Mile End and Quartier Latin as well as academic institutions such as Loyola college/high school. Colloquial references to existing neighborhoods such as *NDG* for *Notre-Dame-de-Grâce* were also identified as were many important street names in Montréal such as Crescent or Ste.-Catherine. Interestingly since these street names were referenced either in French or English, the n-gram which includes the generic type, e.g., Street or Rue (in French), is often not identified as a neighborhood, only the specific name. This is notably different than the other two English-language-dominant cities.

#### 558 6.2. Listing Regional Bias & False Negatives

In the previous section we discuss a number of the false positives and examine some possible explanations. Here we investigate instances where our model did not correctly identify common neighborhoods as well as some of the potential reasons for this. Data from Washington, DC in particular is the subject of further examination and Figure 6 presents a good starting point for this discussion.

The regions represented in purple in this figure are neighborhoods in our *common neighborhood* set that were correctly identified in the initial RF model. The regions shown in orange are those 5 64 neighborhoods that did not appear in the *common neighborhood* set but did appear in at least one of the 565 source-specific neighborhood datasets (Government defined neighborhoods in this case). These are the 566 neighborhoods that were successfully identified by the first iteration of the RF model that were then 567 properly tagged as neighborhoods for input into the second RF model (for use in training a model for other cities). Green regions of the map depict those neighborhoods that were never identified 569 (false negatives), or did not exist, in the n-grams from the Craigslist data. Dark gray regions can be 570 ignored as they represent uninhabitable space such as the Potomac and Anacostia rivers, Rock Creek 571 Park, Observatory Circle, and Joint Base Anacostia-Bolling (military controlled). In observing Figure 6, 572 there is a clear geographic bias between the true positives (blue and orange) and unmentioned or false negatives (green). The green regions are predominantly in the east-southeast region of Washington, 5 74 DC, east of the Anacostia river in what is municipally defined as Wards 7 and 8.9 In referencing the 575 2015 American Community Survey data, we find that Wards 7 and 8 contain the largest number of 576 residence in the district living below the federal poverty line. In addition, the neighborhoods in Wards 577

<sup>&</sup>lt;sup>9</sup> Washington, DC's planning department splits the District into 8 Wards.



Figure 6. Identified and unidentified neighborhoods in Washington, DC.

<sup>578</sup> 7 and 8 contain a mean of 232.3 (median 290) public housing units.<sup>10</sup> By comparison, neighborhoods
<sup>579</sup> in all other Wards list a mean of 173.2 (median 13) public housing units.

Further investigation into the neighborhood names in Wards 7 and 8 show that none of the names 580 or reasonable partial matches of the names occur in the rental listing-based craigslist dataset. Either 581 listings did not occur in those neighborhoods, were too few and thus removed from the dataset during 582 cleaning, or the neighborhood names themselves were not stated in the listings. The mean number 583 of listings per square kilometer or neighborhoods in Wards 7 and 8 is 0.0063 (median 0.0054, SD 5 84 0.0035) whereas for the rest of the neighborhoods showed a mean of 0.0526 (median 0.0352, SD 0.0539) 585 suggesting that the lack of n-gram neighborhood identification was due to the lack of listings, not 586 necessarily missing names in the text or false negatives. This bias in rental listings related to poverty 587 supports existing research in this area [61]. 588

#### 589 7. Conclusions & Future Work

Neighborhoods are an odd concept related to human mobility and habitation. They are difficult
to quantify, and within the domain of geographical sciences, have been historically ill defined.
Neighborhoods are given meaning by the people that inhabit a region based on a set of common or
shared characteristics. Part of the problem, is that a top-down approach to defining a neighborhood is

<sup>&</sup>lt;sup>10</sup> Housing provided for residents with low incomes and subsidized though public funds. Data: http://opendata.dc.gov/datasets/public-housing-areas

fraught with problems and the resulting names and boundaries are often at odds with the citizens that 5 94 live and work within them. In this work, we take a bottom-up and data-driven approach to identifying 5 95 neighborhood names within urban settings. Using geotagged rental property listings from the popular classifieds platform, Craigslist, we demonstrate that neighborhood names can be identified from 597 natural language text within housing rental listings (RQ1). Using an ensemble learning approach based 5 98 on spatial descriptive statistics we demonstrate that it is possible to differentiate neighborhood names 599 from other descriptive natural language terms and phrases (RQ2). Three unique cities within North 600 America are used as focal study sites with listings from one (Washington, DC) being used to train a model that is tested on the other two (Seattle, WA and Montréal, QC). The results of this approach 602 demonstrate that neighborhood names can successfully be identified within the trained city and across 603 different cities (RQ4). In some cases, new, alternative, or previously unidentified neighborhood names 604 are proposed based on this approach (RQ3). Finally, the biases associated with these data are further 605 exposed through this method (RQ5) and are discussed in further detail. 606

As mentioned when discussing the biases associated with this approach, these data really represent the property listers' views of the city. In most cases, the people listing these properties represent a small subset of the city's population, either property owners or real estate agents, both of which tend to exist within a narrow socio-economic group. The neighborhood names identified in the results are therefore heavily influenced by this group. While the methods presented are agnostic to the underlying source of the data, it is important to understand that the neighborhood results depicted in this work are reliant on data contributed to a single online platform.

Similarly, the three example cities used in this research are all within North America. Future work 614 should examine how the results and accuracy values are affected by a change in location. European 615 Cities such as Berlin, for example, could be vastly different given the unique historical context through 616 which the city is understood. Additional work will focus on increasing the diversity of the data 617 sources, languages of the rental property listings, and inclusion of additional structured content (e.g., 618 number of bedrooms, price, etc.). From a statistical perspective, further research will attempt to reduce 619 the dimensionality of this approach by further investigating the correlations between the various 620 spatial statistical measures. Furthermore, a deeper investigation into the role of spatial-autocorrelation, 621 specifically the lack of significance in the results of the Moran's I analysis, will be conducted as this 622 lack of significance is quite interesting and surprising to the researchers. Finally, this work presents the 623 first step of identifying neighborhood names. Our next step is to identify the boundaries associated 624 with these neighborhood names with the goal of developing local listing-based neighborhood datasets. 625

#### 626 References

627 628

629

630 631

632 633 634

1.	Riesz, M. Borders disputed! Brooklynites take issue with Google's neighborhood maps, 2014. https://www.brooklynpaper.com/stories/37/18/all-google-maps-neighborhoods-2014-04-25-bk_37_18.html.
2.	Folven, E. Residents Voice Anger of Redistricting Maps, 2012. http://beverlypress.com/2012/02/ residents-voice-anger-of-redistricting-maps/.
3.	Usborne, S. Disputed territories: where Google Maps draws the line. <i>The Guardian</i> <b>2016</b> . Accessed 20-07-2018.
4.	Sutter, J. Google Maps border becomes part of international dispute. CNN 2010. Accessed 20-07-2018.
E	NI A C I M D NI III I D I D I I I D I NI VAL THE SOLO A

- 5. Nicas, J. As Google Maps Renames Neighborhoods, Residents Fume. *The New York Times* 2018. Accessed
  20-07-2018.
- 6. Taylor, R.B.; Gottfredson, S.D.; Brower, S. Neighborhood naming as an index of attachment to place.
  638 Population and Environment 1984, 7, 103–125.
- Mitrany, M.; Mazumdar, S. Neighborhood design and religion: Modern Orthodox Jews. *Journal of Architectural and Planning Research* 2009, pp. 44–69.

641 642	8.	Knopp, L. Gentrification and gay neighborhood formation in New Orleans. <i>Homo economics: Capitalism, community, and leshian and gay life</i> <b>1997</b> , pp. 45–59.
643	9.	Alderman, D.H. A street fit for a King: Naming places and commemoration in the American South. <i>The</i>
644		Professional Geographer <b>2000</b> , 52, 672–684.
645 646	10.	Hernandez, J. Redlining revisited: mortgage lending patterns in Sacramento 1930–2004. <i>International Journal of Urban and Regional Research</i> <b>2009</b> , <i>33</i> , 291–313.
647	11.	Northcraft, G.B.; Neale, M.A. Experts, amateurs, and real estate: An anchoring-and-adjustment perspective
648		on property pricing decisions. Organizational behavior and human decision processes <b>1987</b> , 39, 84–97.
649	12.	Altschuler, A.; Somkin, C.P.; Adler, N.E. Local services and amenities, neighborhood social capital, and
650		health. Social Science & Medicine 2004, 59, 1219–1229.
651	13.	Calem, P.S.; Gillen, K.; Wachter, S. The neighborhood distribution of subprime mortgage lending. The
652		<i>Journal of Real Estate Finance and Economics</i> <b>2004</b> , 29, 393–410.
653	14.	Romero, M. How real estate websites define Fishtown's boundaries, 2016. https://philly.curbed.com/
654		2016/10/31/13458206/fishtown-neighborhood-boundaries-map.
655	15.	Grether, D.M.; Mieszkowski, P. Determinants of real estate values. Journal of Urban Economics 1974,
656		1, 127–145.
657	16.	Mumford, L. The neighborhood and the neighborhood unit. <i>Town Planning Review</i> <b>1954</b> , 24, 256.
658	17.	Talen, E. Constructing neighborhoods from the bottom up: the case for resident-generated GIS. <i>Environment</i>
659		and Planning B: Planning and Design <b>1999</b> , 26, 533–554.
660	18.	Sieber, R. Public participation geographic information systems: A literature review and framework. <i>Annals</i>
661	10	of the association of American Geographers <b>2006</b> , 96, 491–507.
662	19.	United States. Dept. of Housing and Urban Development. Office of Policy Development and Research.
663	20	The Behavioral Foundations Of Neighborhood Change; University of Michigan Library, 1979.
664	20.	Keller, S.I. The urban neighborhood: A sociological perspective; Vol. 33, Kandom House, 1968.
665	21. 22	Hoyt, H. The structure and growth of residential neighborhoods in American cities; Washington, U.S. Govt., 1959.
666	22.	methodological pote. American journal of community neuchology 2001, 20, 271, 282
667	23	Lee B A : Reardon S E: Eirobaugh C : Farrell C R : Matthews S A : O'Sullivan D Beyond the concus
668	23.	tract: Patterns and determinants of racial segregation at multiple geographic scales. American Sociological
670		Review 2008, 73, 766–791.
671	24.	Sampson, R.J.; Morenoff, J.D.; Gannon-Rowley, T. Assessing "neighborhood effects": Social processes and
672		new directions in research. <i>Annual review of sociology</i> <b>2002</b> , <i>28</i> , 443–478.
673	25.	Schockaert, S.; De Cock, M. Neighborhood restrictions in geographic IR. Proceedings of the 30th annual
674		international ACM SIGIR conference on Research and development in information retrieval. ACM, 2007,
675		pp. 167–174.
676	26.	Hollenstein, L.; Purves, R. Exploring place through user-generated content: Using Flickr tags to describe
677		city cores. <i>Journal of Spatial Information Science</i> <b>2010</b> , 2010, 21–48.
678	27.	Hu, Y.; Gao, S.; Janowicz, K.; Yu, B.; Li, W.; Prasad, S. Extracting and understanding urban areas of interest
679		using geotagged photos. Computers, Environment and Urban Systems 2015, 54, 240–254.
680	28.	Gao, S.; Janowicz, K.; Montello, D.R.; Hu, Y.; Yang, J.A.; McKenzie, G.; Ju, Y.; Gong, L.; Adams, B.; Yan, B.
681		A data-synthesis-driven method for detecting and extracting vague cognitive regions. <i>International Journal</i>
682		of Geographical Information Science <b>2017</b> , 31, 1245–1271.
683	29.	Gao, S.; Janowicz, K.; Couclelis, H. Extracting urban functional regions from points of interest and human
684	•	activities on location-based social networks. <i>Transactions in GIS</i> <b>2017</b> , 21, 446–467.
685	30.	McKenzie, G.; Adams, B. Juxtaposing Thematic Regions Derived from Spatial and Platial User-Generated
686		Content. 15th International Conference on Spatial Information Theory (COSIT 2017); Clementini, E.;
687		Domieny, IVI.; Tuan, IVI.; Kray, C.; Fognaroni, F.; Danatore, A., Eds.; Schloss Dagstuni–Leibniz-Zentrum fuer
688		20.1_20.14 doi:10.4230/LIPIcs COSIT 2017 20
089 600	31	Cranshaw I: Schwartz R: Hong II: Sadeh N. The Livehoods Project: Utilizing Social Modia to
691	01.	Understand the Dynamics of a City. The Sixth International AAAI Conference on Weblogs and Social
692		Media. AAAL 2012.

<sup>693</sup> 32. Wahl, B.; Wilde, E. Mapping the World... One Neighborhood at a Time. *Directions Magazine* **2008**.

- 696 34. Chisholm, M.; Cohen, R. The neighborhood project, 2005. https://hood.theory.org/.
- Hu, Y.; Mao, H.; McKenzie, G. A natural language processing and geospatial clustering framework for
   harvesting local place names from geotagged housing advertisements. *International Journal of Geographical Information Science* 2018, pp. 1–25.
- Zhu, R.; Hu, Y.; Janowicz, K.; McKenzie, G. Spatial signatures for geographic feature types: Examining
   gazetteer ontologies using spatial statistics. *Transactions in GIS* 2016, 20, 333–355.
- Zhu, R.; Janowicz, K.; Yan, B.; Hu, Y. Which kobani? a case study on the role of spatial statistics and
   semantics for coreference resolution across gazetteers. International Conference on Geographic Information
   Science, 2016.
- Brindley, P.; Goulding, J.; Wilson, M.L. A data driven approach to mapping urban neighbourhoods.
   Proceedings of the 22nd ACM SIGSPATIAL International Conference on Advances in Geographic
   Information Systems. ACM, 2014, pp. 437–440.
- Brindley, P.; Goulding, J.; Wilson, M.L. Generating vague neighbourhoods through data mining of passive
   web data. *International Journal of Geographical Information Science* 2018, *32*, 498–523.
- 40. Jones, C.B.; Purves, R.S.; Clough, P.D.; Joho, H. Modelling vague places with knowledge from the Web.
   International Journal of Geographical Information Science 2008, 22, 1045–1065.
- 41. Derungs, C.; Purves, R.S. Mining nearness relations from an n-grams web corpus in geographical space.
   *Spatial Cognition & Computation* 2016, *16*, 301–322.
- 42. Vasardani, M.; Winter, S.; Richter, K.F. Locating place names from place descriptions. *International Journal* of *Geographical Information Science* 2013, 27, 2509–2532.
- 43. Buscaldi, D.; Rosso, P. A conceptual density-based approach for the disambiguation of toponyms.
   *International Journal of Geographical Information Science* 2008, 22, 301–313.
- 44. Gelernter, J.; Mushegian, N. Geo-parsing messages from microtext. *Transactions in GIS* 2011, 15, 753–773.
- Inkpen, D.; Liu, J.; Farzindar, A.; Kazemi, F.; Ghazi, D. Location detection and disambiguation from Twitter
   messages. *Journal of Intelligent Information Systems* 2017, *49*, 237–253.
- 46. Liu, Y.; Wang, F.; Kang, C.; Gao, Y.; Lu, Y. Analyzing Relatedness by Toponym Co-O ccurrences on Web
   Pages. *Transactions in GIS* 2014, *18*, 89–107.
- 47. Santos, J.; Anastácio, I.; Martins, B. Using machine learning methods for disambiguating place references
   in textual documents. *GeoJournal* 2015, *80*, 375–392.
- 48. Melo, F.; Martins, B. Automated geocoding of textual documents: A survey of current approaches.
   *Transactions in GIS* 2017, 21, 3–38.
- <sup>727</sup> 49. Bird, S.; Klein, E.; Loper, E. Natural language processing with Python: analyzing text with the natural language toolkit; "O'Reilly Media, Inc.", 2009.
- <sup>729</sup> 50. Baddeley, A.; Rubak, E.; Turner, R. *Spatial Point Patterns: Methodology and Applications with R*; Chapman
   <sup>730</sup> and Hall/CRC Press: London, 2015.
- Ripley, B.D. The second-order analysis of stationary point processes. *Journal of applied probability* 1976, 13, 255–266.
- 733 52. Besag, J.E. Comment on 'Modelling spatial patterns' by BD Ripley. JR Stat. Soc. B 1977, 39, 193–195.
- McKenzie, G.; Janowicz, K.; Gao, S.; Yang, J.A.; Hu, Y. POI pulse: A multi-granular, semantic signature-based information observatory for the interactive visualization of big geosocial data.
   *Cartographica: The International Journal for Geographic Information and Geovisualization* 2015, *50*, 71–85.
- Graham, R.L. An efficient algorithm for determining the convex hull of a finite planar set. *Information processing letters* 1972, 1, 132–133.
- <sup>739</sup> 55. Ho, T.K. Random decision forests. Document analysis and recognition, 1995., proceedings of the third
  <sup>r40</sup> international conference on. IEEE, 1995, Vol. 1, pp. 278–282.
- <sup>741</sup> 56. Breiman, L. Random forests. *Machine learning* **2001**, *45*, 5–32.
- <sup>742</sup> 57. Chesnokova, O.; Nowak, M.; Purves, R.S. A crowdsourced model of landscape preference. LIPIcs-Leibniz
   <sup>743</sup> International Proceedings in Informatics. Schloss Dagstuhl-Leibniz-Zentrum fuer Informatik, 2017, Vol. 86.
- 58. Oliveira, S.; Oehler, F.; San-Miguel-Ayanz, J.; Camia, A.; Pereira, J.M. Modeling spatial patterns of fire
- occurrence in Mediterranean Europe using Multiple Regression and Random Forest. *Forest Ecology and Management* 2012, 275, 117–129.

- <sup>747</sup> 59. Hayes, M.M.; Miller, S.N.; Murphy, M.A. High-resolution landcover classification using Random Forest.
   <sup>748</sup> *Remote sensing letters* 2014, *5*, 112–121.
- <sup>749</sup> 60. George-Cosh, D. July 1 Is Day for Mass, Messy Moves in Montreal. *The Wall Street Journal* **2013**.
- <sup>750</sup> 61. Boeing, G.; Waddell, P. New insights into rental housing markets across the united states: web scraping
- and analyzing craigslist rental listings. *Journal of Planning Education and Research* 2017, 37, 457–476.

752 Appendix 1

	N Match	SD	Count	Mean NN1	Mean NN2	Mean NN3	Med. NN1	Med. NN2	Med. NN3	Mean Pair	Aed. Pair I	Vipley 1 R	ipley 2 Rip	ley 3 Riple	ey 4 Ripley	5 Ripley (	Ripley 7	Ripley 8	Ripley 9	Ripley 10	Kurtosis	Skewness	Convex Area	Convex Density
N Match	1.000	-0.363	0.029*	-0.125	-0.131	-0.136	-0.089	-0.108	-0.114	-0.371 -	0.337 0	.109 0.	289 0.38	1 0.427	0.447	0.443	0.435	0.429	0.422	0.402	0.0212***	-0.121	-0.201	-0.001***
SD	-0.363	1.000	-0.041*	0.179	0.201	0.210	0.137	0.213	0.237	0.981 0		0.166 -0	.362 -0.5	15 -0.61	1 -0.683	-0.738	-0.774	-0.802	-0.823	-0.844	0.056	0.017***	0.419	-0.309
Count	0.029*	$-0.041^{*}$	1.000	-0.793	-0.820	-0.831	-0.428	-0.510	-0.570	0.010*** -		0.296 -0	.219 -0.1	80.0- 68	6 -0.047	-0.010***	0.022***	0.044	0.053	0.065	-0.174	-0.029*	0.676	0.857
Mean NN1	-0.125	0.179	-0.793	1.000	0.961	0.948	0.779	0.807	0.824	0.169 0	.165 -	0.138 -0	.204 -0.2	37 -0.23	4 -0.233	-0.244	-0.259	-0.255	-0.234	-0.221	-0.083	0.101	-0.257	-0.921
Mean NN2	-0.131	0.201	-0.820	0.961	1.000	0.984	0.739	0.824	0.842	0.188 0	- 184 -	0.117 -0	.198 -0.2	11 -0.24	3 -0.245	-0.261	-0.280	-0.278	-0.259	-0.246	-0.095	0.126	-0.276	-0.943
Mean NN3	-0.136	0.210	-0.831	0.948	0.984	1.000	0.732	0.813	0.852	0.194 0	- 186	0.104 -6	.186 -0.2	34 -0.24	0 -0.245	-0.263	-0.281	-0.280	-0.263	-0.252	-0.081	0.123	-0.291	-0.944
Med. NN1	-0.089	0.137	-0.428	0.779	0.739	0.732	1.000	0.906	0.871	0.178 0	-126 -	0.533 -0	.539 -0.4	5 -0.41	7 -0.360	-0.337	-0.321	-0.292	-0.249	-0.218	-0.335	0.143	0.011***	-0.607
Med. NN2	-0.108	0.213	-0.510	0.807	0.824	0.813	0.906	1.000	0.964	0.254 0	- 209 -	0.446 -0	.505 -0.5	J5 -0.45	6 -0.417	-0.401	-0.393	-0.366	-0.323	-0.295	-0.317	0.171	-0.021***	-0.703
Med. NN3	-0.114	0.237	-0.570	0.824	0.842	0.852	0.871	0.964	1.000	0.273 0	.235 -	0.369 -0	.445 -0.4	58 -0.43	7 -0.412	-0.406	-0.404	-0.382	-0.342	-0.316	-0.262	0.162	-0.077	-0.742
Mean Pair	-0.371	0.981	0.010***	0.169	0.188	0.194	0.178	0.254	0.273	1.000 0	- 955 -	0.242 -0	.445 -0.5	-0.69	4 -0.762	-0.810	-0.838	-0.857	-0.870	-0.881	-0.001	0.047	0.465	-0.276
Med. Pair	-0.337	0.932	-0.039*	0.165	0.184	0.186	0.126	0.209	0.235	0.955 1	- 000.	0.147 -0	.352 -0.5	10 -0.62	7 -0.710	-0.768	-0.811	-0.847	-0.877	-0.903	0.072***	0.004***	0.364	-0.285
Ripley 1	0.109	-0.166	-0.296	-0.138	-0.117	-0.104	-0.533	-0.446	-0.369	-0.242 -	0.147 1	.000	906 0.74	9 0.589	0.478	0.414	0.370	0.319	0.268	0.230	0.578	-0.189	-0.555	-0.042*
Ripley 2	0.289	-0.362	-0.219	-0.204	-0.198	-0.186	-0.539	-0.505	-0.445	-0.445 -	0.352 0	.906 1.	000 0.93	0 0.808	0.708	0.641	0.588	0.532	0.479	0.440	0.639	-0.337	-0.581	0.071
Ripley 3	0.381	-0.505	-0.139	-0.237	-0.241	-0.234	-0.495	-0.505	-0.468	-0.591 -	0.510 0	.749 0.	930 1.00	0 0.948	0.875	0.809	0.754	0.698	0.646	0.607	0.569	-0.393	-0.544	0.156
Ripley 4	0.427	-0.611	-0.086	-0.234	-0.243	-0.240	-0.417	-0.456	-0.437	- 0.694 -	0.627 0	.589 0.	808 0.94	8 1.000	0.969	0.916	0.862	0.812	0.767	0.731	0.429	-0.371	-0.503	0.204
Ripley 5	0.447	-0.683	-0.047	-0.233	-0.245	-0.245	-0.360	-0.417	-0.412	-0.762 -	0.710 0	.478 0.	708 0.87	5 0.969	1.000	0.974	0.930	0.888	0.850	0.817	0.298	-0.310	-0.469	0.238
Ripley 6	0.443	-0.738	-0.010***	-0.244	-0.261	-0.263	-0.337	-0.401	-0.406	-0.810 -	0.768 0	.414 0.	641 0.80	9 0.916	0.974	1.000	0.979	0.946	0.912	0.878	0.192	-0.236	-0.445	0.271
Ripley 7	0.435	-0.774	0.022***	-0.259	-0.280	-0.281	-0.321	-0.393	-0.404	-0.838 -	0.811 0	.370 0.	588 0.75	4 0.862	0.930	0.979	1.000	0.984	0.954	0.920	0.112	-0.160	-0.425	0.300
Ripley 8	0.429	-0.802	0.044	-0.255	-0.278	-0.280	-0.292	-0.366	-0.382	-0.857 -	0.847 0	.319 0.	532 0.69	8 0.812	0.888	0.946	0.984	1.000	0.986	0.957	0.039*	-0.090	-0.401	0.314
Ripley 9	0.422	-0.823	0.053	-0.234	-0.259	-0.263	-0.249	-0.323	-0.342	- 0.870	0.877 0	.268 0.	479 0.64	6 0.767	0.850	0.912	0.954	0.986	1.000	0.986	-0.012***	-0.050	-0.380	0.312
Ripley 10	0.402	-0.844	0.065	-0.221	-0.246	-0.252	-0.218	-0.295	-0.316	-0.881 -	0.903 0	.230 0.	440 0.60	7 0.731	0.817	0.878	0.920	0.957	0.986	1.000	-0.037*	-0.035*	-0.360	0.312
Kurtosis	0.0212***	0.056	-0.174	-0.083	-0.095	-0.081	-0.335	-0.317	-0.262	-0.001*** 0	072 0	.578 0.	639 0.56	9 0.429	0.298	0.192	0.112	0.039*	-0.012***	-0.037*	1.000	-0.764	-0.277	-0.042*
Skewness	-0.121	0.017***	-0.029*	0.101	0.126	0.123	0.143	0.171	0.162	0.047 0	- ***	0.189 -0	.337 -0.3	33 -0.37	1 -0.310	-0.236	-0.160	-0.090	-0.050	-0.035*	-0.764	1.000	0.060	-0.077
Convex Area	-0.201	0.419	0.676	-0.257	-0.276	-0.291	0.011***	-0.021***	-0.077	0.465 0	.364 -	0.555 -0	.581 -0.5	44 -0.50	3 -0.469	-0.445	-0.425	-0.401	-0.380	-0.360	-0.277	0.060	1.000	0.294
<b>Convex Density</b>	-0.001***	-0.309	0.857	-0.921	-0.943	-0.944	-0.607	-0.703	-0.742	-0.276 -	0.285 -	0.042* 0.	071 0.15	6 0.204	0.238	0.271	0.300	0.314	0.312	0.312	-0.042*	-0.077	0.294	1.000

*** indicates no significance, and all other values are	
< 0.1, *	
* indicates p -	
A1. Pearson's correlation matrix for all predictive spatial statistics measures.	cant to $p < 0.01$ .
Table	signifi

Neighborhood names (both true and false positives) as identified by the random forest ensemble
 learning model.

756

### 757 Washington, DC

adams, adams morgan, alexandria va, american, and downtown, apartments in alexandria, arlington, 758 arlington va, bloomingdale, branch, brookland, capitol, capitol hill, cathedral, chase, chevy, chevy 759 chase, chinatown, circle, circle and, cleveland, cleveland park, columbia, columbia heights, crystal, 760 crystal city, downtown, downtown bethesda, downtown silver, downtown silver spring, dupont, 761 dupont circle, foggy, foggy bottom, forest, fort, friendship, friendship heights, from downtown, george, 762 georgetown, georgetown and, georgetown university, georgia, glover, glover park, green, heights, 763 howard, in alexandria, in arlington, kalorama, logan, logan circle, morgan, navy, navy yard, noma, of 764 old town, old town, old town alexandria, petworth, pleasant, potomac, shaw, silver spring, silver 765 spring md, spring, spring md, stadium, takoma, takoma park, to downtown, to dupont, to dupont 766 circle, to georgetown, to silver, to union, to union station, town alexandria, triangle, u corridor, union, union station, university, vernon 768

769 770

#### 771 Seattle, WA

admiral, alki, alki beach, and redmond, anne, ballard, ballard and, beacon, beacon hill, belltown, bothell, bothell wa, broadway, by windermere, capitol hill, columbia, columbia city, corridor, eastlake, 773 first hill, fremont, green lake, greenlake, greenwood, heart of capitol, heart of downtown, interbay, 774 international district, junction, lake city, lake union and, lincoln, lower queen, lower queen anne, 775 madison, magnolia, mercer, ne seattle, ne seattle wa, north seattle, northgate, northgate mall, of 776 ballard, of capitol, of capitol hill, of lake union, of queen, of queen anne, phinney, phinney ridge, pike, pike pine, pike place, pike place market, pine, pine corridor, pioneer, pioneer square, queen anne, 778 ravenna, roosevelt, seattle center, seattle central, seattle downtown, seattle university, shoreline, south 779 lake, south lake union, stevens, the junction, the university district, to green, to green lake, u district, 780 union and, university district, university village, uptown, uw campus, wallingford, west seattle, 781 westlake, windermere, woodland, woodland park

783 784

# 785 Montréal, QC

and downtown, canal lachine, cote des, cote des neiges, dame, dame de, des neiges, downtown, 786 downtown and, downtown montreal, du mont royal, du plateau, from downtown, griffintown, heart 787 of downtown, henri, in downtown, in ndg, lachine, lasalle, laurent, le plateau, loyola, marie, mile end, 788 minutes to downtown, monk, monkland, monkland village, mont royal, mont royal et, mount, mount 789 royal, neiges, nord, notre, notre dame de, of downtown, of downtown montreal, of the plateau, old 790 montreal, old port, outremont, plateau, plateau mont, plateau mont royal, rosemont, royal, saint, saint 791 laurent, snowdon, st henri, te des neiges, the lachine, the plateau, to downtown, tro mont, tro mont 792 royal, verdun, villa maria, village, ville, ville marie, villeray, westmount 793

794

795

<sup>796</sup> © 2018 by the authors. Submitted to *ISPRS Int. J. Geo-Inf.* for possible open access
 <sup>797</sup> publication under the terms and conditions of the Creative Commons Attribution (CC BY) license
 <sup>798</sup> (http://creativecommons.org/licenses/by/4.0/).