

Visualizing demographic evolution using geographically inconsistent census data

Fabio Dias and Daniel Silver

Abstract—We propose a visual analytics system that enables the exploration evolutionary patterns in geographically inconsistent data, removing the need to harmonize it into the same geographical regions, a time consuming and error-prone process that is currently used in virtually all longitudinal analysis of geographical data. This work also includes incremental developments in the representation, clustering, and visual exploration of geographical data. We illustrate it considering census data, where it allows an easier understanding of the demographic groups present in a city and their evolution over time. We present the feedback of experts in urban sciences and sociology, along with illustrative scenarios in the USA and Canada on the decennial censuses between 1970 and 2010.

1 INTRODUCTION

URBAN sciences are blooming thanks to a renewed interest in understanding and improving the urban environment. Visual analytics is following this trend, fueled by new public datasets that encompass progressively more of our daily lives [1]. There is no shortage of methods to explore mobility patterns [2], social media [1], traffic [3], and so on, providing experts, planners, policy makers, and the general population with deeper insights about their cities.

These new datasets usually contain GPS coordinates for the records, leading to *point-based* data. Combined with the corresponding timestamps, this data is easily suitable for longitudinal analysis. But most demographic datasets are *region-based*, where the measurements are associated with pre-defined regions, not only for an additional level of privacy protection, but because some measurements only make sense over a defined area. Census data is a classic example of this format, with datasets available from 1790 onwards for the US [4]. Despite this unmatched temporal availability, longitudinal analyses of census data are often restricted in time, especially when smaller tabulation areas are considered, such as census tracts (CT) or dissemination areas, which evolve to reflect changes in population density, leading to geographic inconsistencies across time, and the traditional time-series based approach is no longer viable.

However, these analyses are necessary to understand the urban environment. Indeed, two different regions can have similar average income for a given year, while one is experiencing a process of economic improvement and the other one impoverishment. Quite obviously, a single snapshot cannot be used to identify gentrification, migration, education changes, or any of the relevant processes that happen over time.

To overcome these inconsistencies, the traditional approach is the *geographical harmonization* of the data, the interpolation of the measurements into a common set of

regions [5], [6], [7], so that each variable can be represented using time-series. This is laborious work that inevitably introduces some amount of error [8], even when additional data is provided [9]. Nevertheless, this step is considered mandatory in the current literature: “(...) *tract-by-tract comparison is not possible unless data from 2000 is interpolated to 2010 boundaries (...)*” [10].

The main contribution of this application paper is an visualization-based alternative to the geographical harmonization, a combination of established graph based processing and information visualization techniques allowing tract-by-tract comparison, the identification and visualization of patterns of demographic evolution without geographic harmonization, effectively removing one of the most challenging problems in longitudinal demographic analysis. We also include illustrative scenarios and our prototype is available at <http://uoft.me/piccard>, including more than fourty regions in the US and Canada.

2 RELATED WORK

Since our problem encompasses several fields, we divided this section into specific subproblems: *longitudinal demographic studies*, describing the traditional approach to perform longitudinal studies; *data representation*, exploring how evolving geographic data can be represented for processing; *Data clustering*, briefly reviewing existing clustering methods; and *cluster characterization*, exploring how the clusters can be visually summarized.

2.1 Longitudinal demographic studies

Census data is used not only to discover demographic patterns [11], but to correlate demographic characteristics to other measurements [12]. However, longitudinal studies are rare: “(...) One of the most challenging and fascinating areas in spatial statistics is the synthesis of spatial data collected at different spatial scales(...)” [13].

While CT level data is readily available for the US since 1910 [4], most studies consider the period between 1970 and 2010, using pre-harmonized data [4], [5]. Despite the inherent errors [6], [8], this dataset became the standard

• F. Dias is with the Department of Mechanical & Industrial Engineering, University of Toronto, Toronto, ON, M5S 3G8.
E-mail: fabio.dias@utoronto.ca

• D. Silver is with the Department of Sociology, University of Toronto, Toronto, ON, M5S 2J4.
E-mail: dsilver@utsc.utoronto.ca

Manuscript received MMMM DD, YYYY; revised MMMM DD, YYYY.

75 source for longitudinal demographic data, with similar efforts appearing in other countries [7], [14], [15]. This result
 76 was significant for the field, but it also restricts the useable
 77 data, since new datasets need to be similarly processed.

78 Another option considers the use of grid data [10], [16],
 79 where small rectangular areas are used, in an approach
 80 similar to satellite imagery. Beyond the increased spatial
 81 accuracy, this approach does not require complex harmonizations
 82 when new data is considered. However, demographic
 83 data is usually not available in this format, especially from
 84 older sources, and the conversion from tabulation areas can
 85 introduce significant errors.

86 In the proposed methodology, we avoid the harmonization
 87 by considering each measurement using its actual
 88 geographic region. It does not require the regions to be
 89 consistent across time because they are already represented
 90 as different entities.

92 **2.2 Data representation**

93 Most data is represented in tabular form, where the rows
 94 and columns have coherent definitions. For example, consider
 95 a table with rain measurements over time, with the
 96 rows representing different locations and the columns differ-
 97 ent times. This representation can also be interpreted as
 98 a collection of time-series, one for each location. Geographic
 99 data followed this format, only including an additional field
 100 that describes the associated geographic area. Following the
 101 example, the data would now represent the amount of rain
 102 for a given region and time. As long each region remains
 103 the same, the data is coherent and can be interpreted again
 104 as a collection of time-series.

105 In the proposed method, we remove the requirement
 106 for consistency in the measurement regions by leveraging a
 107 graph-based representation, where each region in time cor-
 108 responds to a different node. Instead of a collection of time-
 109 series, the data is represented as a dynamic graph. Graph
 110 based representation of geographic information is fairly well
 111 explored in the literature, as a basis for topological methods
 112 for event detection [17], leveraging signal processing on
 113 graphs [18], [19] to find patterns and outliers [20], [21],
 114 [22]. Graphs are well suited to represent trajectories as
 115 well [2], [3], [23], allowing the use of graph visualization
 116 methods [24], [25].

117 Graphs were used to represent census data for clustering
 118 purposes before [21], [26], but these works did not explore
 119 temporal evolution, where graphs are particular powerful as
 120 they allow a natural representation of inconsistent regions,
 121 with both spatial and temporal connections. Note that there
 122 are other possible representations that have similar proper-
 123 ties, but we adopted graphs to allow the use of the existing
 124 literature and methods.

125 **2.3 Data clustering**

126 Data clustering is one of the elementary processes for data
 127 analysis, simplifying the data into a smaller number of
 128 homogeneous sets that can be interpreted in the same way.
 129 While there is no shortage of contributions for this prob-
 130 lem [27], most applications still rely on k-means [28], [29]
 131 and, to a lesser extent, Self Organizing Maps [30], [31].

132 However, a method for geographic data analysis should
 133 not ignore the geographic component of the data. One
 134 straightforward option, for agglomerative methods [32], is
 135 to consider only nearby clusters for merging [33], which
 136 can also be done for k-means [34]. Alternatively, the spatial
 137 distance could be directly added to the inter-cluster met-
 138 ric [33] via a mixing parameter, which adds flexibility to the
 139 method, but introduces the problem of finding the correct
 140 application-dependent values.

141 Indeed, one crucial step in most clustering algorithms
 142 is the definition of the number of clusters. We sidestep
 143 this problem by considering hierarchical methods [35],
 144 where the result is not a partition of the data, but a tree
 145 of partitions. This approach is interesting for interactive
 146 methods, because it allows the user to change the num-
 147 ber of displayed clusters with minimal processing. Since
 148 our data is represented as a graph, one option would be
 149 the watershed cuts algorithm [36], inspired by the well
 150 known image processing segmentation and equally prone
 151 to oversegmentation. Considering that the processing time
 152 is also a relevant factor, we opted for an heuristic variation
 153 of the maximum weighted matching algorithm called *sorted*
 154 *maximal matching* [37], which merges clusters based on the
 155 weights of the edges between pairs of clusters.

156 **2.4 Cluster characterization**

157 Visually representing evolving spatial data is a challenging
 158 old problem [38], [39], [40], [41]. Most geographic data
 159 is naturally bidimensional and maps work well in this
 160 case [41], [42], but the temporal dimension cannot be so nat-
 161 urally represented. One straightforward option is to lever-
 162 age tridimensional plots [43], [44], but this can lead to visual
 163 obstructions or scaling problems unless a tridimensional
 164 display device is used. Animation can also be explored in
 165 some specific cases [45], but it is not a general approach.
 166 Glyphs can also be used [46], [47], but this may lead to
 167 cluttering when many small regions are present. A simpler,
 168 well adopted, option is to display a map that corresponds
 169 to a subset of the temporal information, allowing the user
 170 to change the time with an associated control [1], [17], [20],
 171 [22]. Small multiples can be used [2], but only when there
 172 are few temporal snapshots. However, none of these options
 173 is suitable to represent many variables at the same time.

174 Using data clustering, we can represent the region's
 175 cluster instead of all the its variables [2], [20], [22]. While
 176 this simplifies the geographic portion of the visualization, it
 177 introduces the problem of how to summarize the contents of
 178 each cluster. One traditional approach is to use parallel coor-
 179 dinates plot [48], [49], [50], [51], but these they can get clut-
 180 tered representing similar clusters over several variables.
 181 Further, for demographic applications, the clusters are usu-
 182 ally strongly characterized by a small subset of values [29],
 183 [30]. Therefore, in the proposed method, we identify the
 184 variables that are most relevant to the characterization of
 185 each cluster. The distribution of values on that variable is
 186 then represented using a boxplot, a well known statistical
 187 plot displaying basic properties of the distributions.

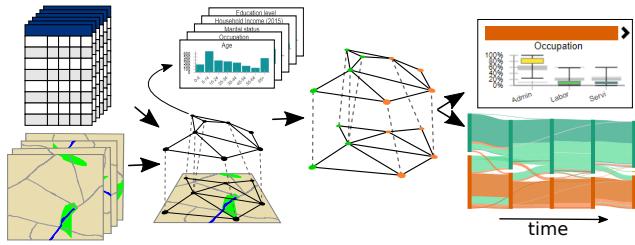


Fig. 1. Overview of the proposed method. A graph is generated combining the original census data, encoding the changing geographical information. The graph is partitioned into an hierarchy [37]. The characteristics and evolution of the clusters are then visually represented.

188 3 VISUALIZING THE DEMOGRAPHIC SPATIO- 189 TEMPORAL EVOLUTION

190 Beyond the objective of allowing the study of inconsistent
191 data, our method includes incremental developments in
192 most steps of the analysis, from data representation to the
193 visualization method for the clusters. Figure 1 presents an
194 overview of the processing steps of the proposed method.

195 3.1 Census methodology and data representation

196 Census data is disseminated in a tabulated form for aggre-
197 gation areas: whole country, state/province, metropolitan
198 region, and so on. To provide as much detail as possible, we
199 focus on the smallest region with available data: *census tracts*
200 (CT). They are usually defined to maintain the anonymity
201 of the population, leading to a population count in the
202 order of thousands in densely populated areas. Physical
203 barriers are usually adopted as borders, so these regions can
204 change because of new roads, construction or removal of
205 high density buildings, and so on. Some census entities also
206 consider demographic characteristics, aiming to establish
207 the CTs as a cohesive unit. Therefore, CTs are the least
208 geographically stable tabulation area.

209 Each CT is associated with a series of variables, with
210 counts derived from the census questionnaires, covering
211 several aspects of the demographic characteristics of the
212 population. Some questions allow for multiple choices
213 or open answers, that are then tabulated into the most
214 frequent categories. Since the census is often used to
215 direct government initiatives, which variables are mea-
216 sured/disseminated is dependent on administrative inter-
217 ests, the general understanding of the population, and cur-
218 rent customs. For instance, income is disseminated with a
219 finer tabulation in the lower portion than on the higher.

220 To match these variables over time and allow for direct
221 comparison across different census years, we aggregated
222 similar ones (e.g. White, Black, Asian, Other) into *aspects*
223 (e.g. Race), encoding the distribution of that facet of the
224 population. In this convention, we refer to the composing parts
225 of an aspect as a *part* or the traditional *variable*. Internally,
226 the aspects are represented using normalized histograms.
227 This normalized representation is crucial for the comparison
228 between inconsistent regions.

229 In our graph based representation, each CT of each
230 census year is represented as a node, and edges are placed
231 between nodes if the corresponding CTs share geographic

borders in the same year. Further, edges are placed between nodes if the corresponding CTs belong to sequential years and there is geographical overlap between them. This approach leads to a single graph representing the whole spatio-temporal space of the data. Our objective then becomes to identify partitions of this graph such that the nodes of each partition are more similar between themselves than to the other nodes. This representation is not the only option, nor unique, but it allows the use of existing graph-based methods for the other steps.

242 3.2 Geographic content clustering

243 To partition the graph we must first establish a distance
244 function between the nodes, measuring the data similarity.
245 This similarity is then associated with the edges, leading
246 to a weighted dynamic graph. Every node has a collection
247 of histograms, each representing the distribution of certain
248 aspect in the population.

249 Let $G = (V, E)$ be a graph, where $V = \{v_1, v_2, \dots, v_n\}$
250 is the set of nodes and $E = \{(v_i, v_j), i \neq j \text{ and } i, j \in [1, n]\}$
251 is the set of edges. A function H associates each node to a
252 set of K histograms. We define the distance D between two
253 nodes v_i and v_j as:

$$254 D(v_i, v_j) = \sum_{k \in [1, K]} w_k d(H_k(v_i), H_k(v_j)) \quad (1)$$

255 where d is a distance metric between histograms and w is
256 a sequence of non-negative weights associated with each
257 aspect, $\sum_{k \in [1, K]} w_k = 1$. While any histogram metric can
258 be used, we adopted a euclidean distance between the
259 vectors, because it led to reasonable results with reduced
260 computational cost. Therefore the distance between two
261 nodes is defined as the weighted average distance between
262 its associated histograms, where the weights can be adjusted
263 by the user.

264 Once the distances are associated to the edges, we use
265 watershed cuts [36] to create an initial clustering, which
266 is then refined into a hierarchy using the Sorted Maximal
267 Matching (SMM) [37] with median linkage. The initial
268 watershed step is performed to create an initial clustering and
269 reduce the running time of the SMM. For completeness, we
270 briefly review this method, but we refer the reader to the
271 original paper [37] for more details, including a complete
272 performance evaluation using several metrics.

272 We included two application-specific parameters: the
273 maximum number of clusters to be shown and a distance
274 threshold. Contrarily to the original SMM, which merges all
275 clusters in all steps, we only merge two clusters where the
276 distance is above the threshold after we reach the maximum
277 number of displayed clusters. Without this restriction, signif-
278 icantly different clusters would be merged early, leading
279 to increased intra cluster variance and the disappearance of
280 small outlier regions. Further, after the maximum number
281 of clusters is reached, we create one step of the hierarchy
282 for each merge, leading to a binary partition tree. In this
283 structure we can directly access a result with an arbitrary
284 number of clusters.

285 Each resulting cluster is contiguous in the graph. This
286 means that two similar, but non-contiguous, sets of CTs will

be classified into two different clusters, which can be counterintuitive. To overcome this issue, we *augment* the graph with two new edges per node from a nearest neighbors graph [52] using only the distances between the histograms. These edges connect nodes with similar content, if they are not already connected, providing a path for the algorithm to group similar nodes. Theoretically, adding more of these content based edges could be used to decrease the impact of the spatio-temporal edges, controlling the balance between content and topology in the result. In practice, the effect is dependent on the data itself, and the results are not consistent, or predictable, across different cities. We fixed it at two edges because it was the lowest number that empirically led to consistent clusters, but we believe that this idea warrants further investigation, as an alternative to mixing parameters in the distance metric [33].

3.3 Cluster characterization and variable relevance

The composition of each cluster is determined by simple statistical measures, considering each aspect separately. We compute the minimum, maximum, median, 25%, and 75% quantiles for each part of each aspect for all clusters in the hierarchy. While interpreting these values is more complex than interpreting just the average, they provide far more information about the underlying distribution.

We also use these statistical measurements to discover what characterizes each cluster, that is, what makes it different from the others. We define the *relevance* of a part of an aspect based on the distance between the interquartile ranges (IQR) of the clusters in the same hierarchical level. If the IQRs overlap for all clusters, that variable is not relevant to the characterization of the cluster, but if the IQRs are distant, it means that this specific range of values is something that only occurs in this cluster.

3.4 Clusters and trajectories

While the partition of the data into different clusters helps the user to understand what demographic groups exist and where they are, we are also interested in the changes in those groups. To this end, we introduce the concept of *trajectories*, which are composed by regions that are classified into the same sequence of clusters over the considered period. This enables direct access to regions that evolved in the same manner. While the interface provides access to data by census tract, the trajectories are the main unit of exploration in this work.

3.5 Colors

As illustrated by Figure 3 and further explored in the next subsection, our proposed interface heavily relies on color to express cluster-related information. We adopted this convention because colors can be used in all our visual tools in a coherent manner. However, this also introduced significant challenges. The first is the limit on the number of clusters that can be visually represented. We limited the number of clusters to eight because this was the largest number of colors that we could reliably use, derived from the 8-class Dark2 set from ColorBrewer [53].

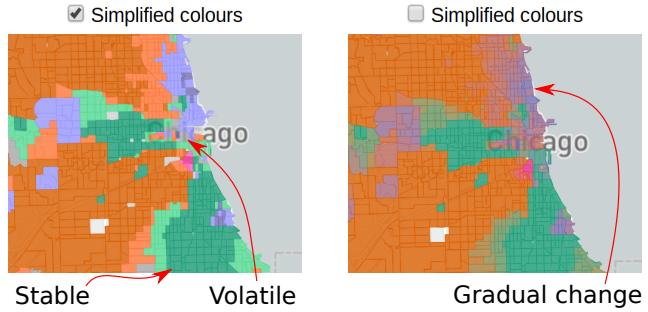


Fig. 2. Different color schemes for Chicago with four clusters. Left: simplified, right: average color.

While we can reasonably limit the number of clusters, there are far more possible trajectories. And the color associated with each trajectory should bear some resemblance to the clusters included in it. Therefore, we were left with a conundrum: *Should we associate each trajectory with a unique color, which the user probably cannot distinguish, or should we use a reduced set of colors and associate the same color to different trajectories?* Since there are advantages and disadvantages for each of those options, we adopted both. The user can control which color policy is used via a checkbox in the configuration panel, on the top left of the interface.

By default, the interface adopts a simplified color scheme, where a trajectory is painted in the same color of a cluster if the regions were associated to that cluster for *all* times; in a slightly less saturated version of that color if the regions were associated to that cluster for *the simple majority* of the time, and gray otherwise. In this mode, the colors will mostly represent stability, immediately identifying the regions that were consistently associated with each cluster. It also easily identifies volatile regions, painted gray.

When this simplified color scheme is disabled, each trajectory will be painted using the average of the colors of the involved clusters, in the LAB color space. In this mode, the map becomes more similar to a heatmap, where stronger presence of a color indicates more temporal affinity to the cluster. Volatile regions will also tend to be displayed in gray, as the average of three or more colors.

While both approaches will use more than eight colors, in practice this is not as significant because most cities can be explained using less than eight clusters. In fact, articles in the literature usually employ from two to five, which are fairly stable across time. For the more dynamic scenarios, user interaction can be used to alleviate the shortcomings of both approaches.

3.6 User interface

The initial interface is illustrated in Figure 3. Since demographic data can be nuanced, with intricate interconnections, we decided against presenting the interface using synthetic data, considering instead data from the Chicago region between 1970 and 2010, using previous published studies as corroboration. This region is known for its entrenched racial divide and the emergence of a '*young urban*' population with a higher education level [29], [30]. More details about this scenario are presented in Section 4.

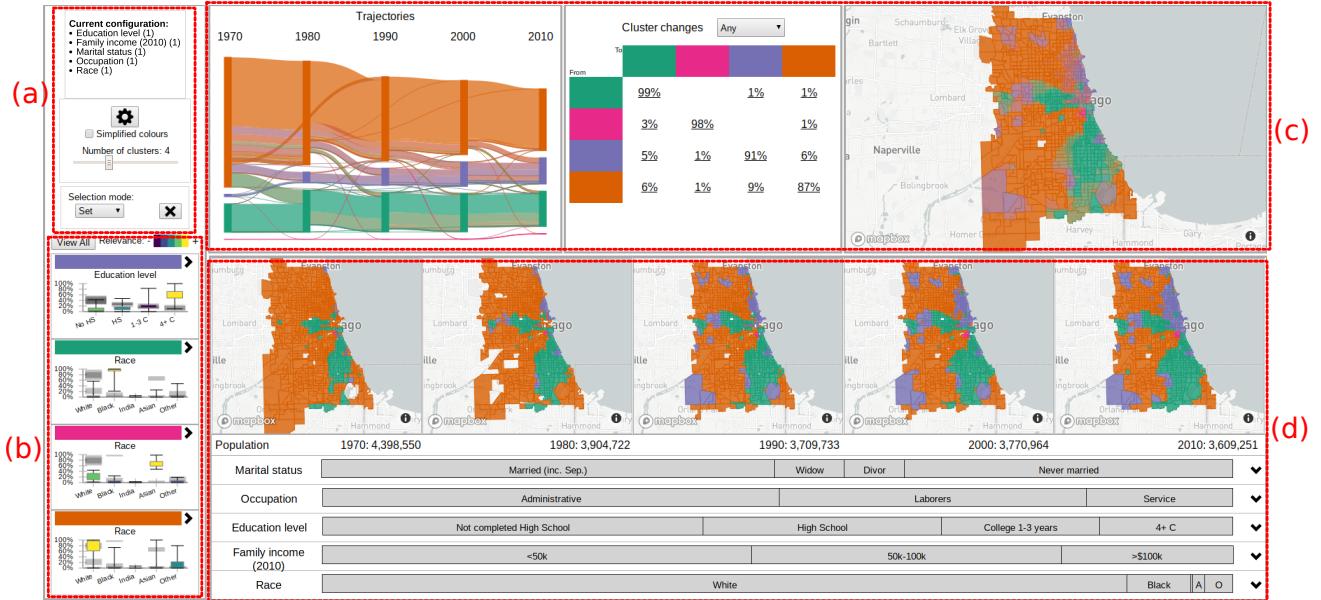


Fig. 3. Initial interface of our method showing the demographic evolution of Chicago. (a): Configuration panel with the current clustering parameters and controls. (b): Cluster overview illustrating the most relevant aspect for each cluster. (c): Trajectories overview and the general evolution of the population, geographical information, and how it changed. (d): Details of the selected trajectories, including precise geographic locations, population numbers, and the composition of the aspects.

The configuration panel, on top left in Figure 3, displays which aspects were considered and their weights (following Equation 1). It also includes other configuration options that can be altered without re-processing the data, such as the number of clusters and the color option. The gear button allows access to the other configuration options that do require further processing, such as changing location, aspects, and weights. This panel also includes the configuration of the selection mode for the trajectories, which allows the user to set, add, or remove the next selected trajectories to the current selection. This feature enables the analysis of complex sets of trajectories.

The cluster overview panel, on the bottom left in Figure 3, displays a brief summary of each cluster, based on the distance between the IQRs, as detailed in Section 3.3. The *View all* button opens a new panel where all aspects are represented, while the chevron at the side of the color lets the user expand each cluster separately. While the standard approach to represent cluster characteristics is to use parallel coordinates [50], [51], this representation occupies screen space proportional to the number of variables and can get cluttered with a higher number of clusters, or when the clusters are not well defined for multiple variables. To save space and leverage the familiarity scientists have with statistical tools, we opted to use boxplots to properly convey the distribution of each variable in the current cluster. However, a simple boxplot would not include information about the other clusters, forcing the user to mentally compare them to find what is relevant.

We adopt an *enhanced* version of the traditional boxplot, which includes the minimum, maximum, 25% and 75% quantiles for the current cluster, but also the IQRs for the other clusters, in slightly larger and faded black rectangles. We also color the current IQR according to its relevance. While there might be some degree of similarity between the

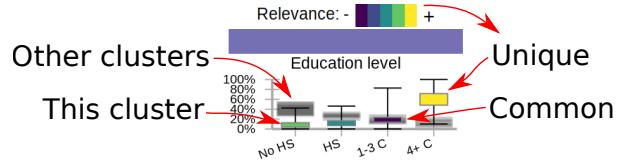


Fig. 4. Enhanced boxplot of the clusters' characteristics allows a quick comparison to the other clusters.

color schema for relevance and for trajectory identification, none of the experts consulted reported confusion. These simple changes allow the user to easily understand the composition of the cluster and how it relates to the others. Violin plots [54] could also be used, but with increased potential for cluttering the representation of the other clusters, hindering the comparisons.

For instance, the boxplot that summarizes the purple cluster, detailed on Figure 4, illustrates that this cluster is best defined by the proportion of the population with four or more years of college. The user can quickly see that this is relevant because the corresponding IQR is colored with the highest relevance present in the legend. It is also clear that, while this cluster includes CTs that have between 10% to 90% of people in this variable, approximately, half of them have about 60% of the population with four or more years of college. Since all the other IQRs are well separated, this is a defining characteristic of this cluster. By clicking on the colored bar above the boxplot, the user can select all trajectories that contain this cluster at any point in time.

The other clusters identified on Figure 3 correspond to higher concentration of people that identify as Black in the green cluster, people that identify as "Asian, Hawaiian, other pacific islander" in the pink cluster, and people that identify as White in the orange cluster. From these plots,

446 it is clear that the city is indeed racially divided [29], with
 447 several CTs that are almost exclusively occupied by people
 448 of the same racial category.

449 The trajectories overview aims to convey basic information
 450 about the trajectories present in the data, where they
 451 are, and what changes are involved. This is done using
 452 three sub panels. The first, on the left, contains a Sankey
 453 diagram illustrating the evolution of the clusters over time.
 454 The widths are proportional to the population involved, the
 455 colors follow a policy detailed in Section 3.5. A stacked
 456 graph could also be used to represent the proportions of
 457 each cluster [20] with less clutter, since the transitions be-
 458 tween clusters would not be represented. However, this is
 459 only viable if more temporal steps are available, making
 460 the plot smoother. Another option to remove clutter is to
 461 remove portions of this plot when trajectories are selected,
 462 but this would change the layout and compromise the user’s
 463 mental map.

464 In our example in Figure 3, we can see that the total
 465 population of Chicago is decreasing. Additionally, the or-
 466 ange and green clusters contain most of the population and
 467 are fairly stable over time. The pink cluster is small and
 468 mostly stable. The purple cluster is increasing, mostly by
 469 incorporating areas that were previously orange. Since the
 470 purple group corresponds to the emergent ‘young urban’
 471 group, this corroborates the findings of Delmelle [29], [30].
 472 This diagram can also be used to select specific trajectories,
 473 by clicking on the bands, or all trajectories that contain
 474 a specific cluster at a specific time, by clicking on the
 475 rectangles.

476 The next sub panel, in the top middle of Figure 3, is
 477 a transition matrix between the clusters. It indicates the
 478 percentage of the population whose area changed in each
 479 pattern. This kind of table can be found in the related
 480 literature [29], so it is familiar to the advanced users, and
 481 it not only informs the proportional changes, but allows
 482 the selection of the corresponding trajectories for further
 483 analysis.

484 Contrary to the trajectories plot, this representation is
 485 more Markovian, where only the current and next state are
 486 considered. This panel also enables easier access to trajec-
 487 tories with specific changes, by clicking on the corresponding
 488 percentage values. The combo box allows the user to refine
 489 the transitions, from ‘Any’, which includes all transitions be-
 490 tween years, to specific transitions, to changes from the first
 491 year to the last year. In this example, approximately 99%
 492 of the population in areas classified as green were also in
 493 areas classified as green in the next year, while 1% changed
 494 to purple at some point and another 1% to orange. The total
 495 is over 100% due to rounding errors. Regions changed from
 496 the orange to the green cluster for 6% of its population, 1%
 497 to pink, and 9% to purple. This further corroborates the fact
 498 that most of the growth of the purple cluster came from the
 499 orange cluster. Additionally, the lack of transitions is also
 500 relevant, for instance, no CT changed from majority of Black
 501 population (green) to Asian population (pink), and no CT
 502 with significant Asian population had significant increase
 503 in education levels (purple).

504 The last sub panel, in the top right of Figure 3, is a map of
 505 the region under analysis, summarizing the evolution over
 506 time. The objective of this panel is to provide an overview

507 of the general position of the trajectories and clusters. The
 508 colors are derived from the clusters involved in each trajec-
 509 tory as detailed in Section 3.5, which are consistent across
 510 the linked views.

511 The bottom part of the interface contains the details for
 512 the selected trajectories. If no trajectory is selected, then the
 513 details for the whole city are displayed, as illustrated in
 514 Figure 3. This panel contains two main regions: the small
 515 multiple maps, depicting the clusters at each year, and the
 516 stacked bar plots that summarize the overall composition
 517 of these regions. Some finer localization information is lost
 518 using small multiples, such as small border changes, but
 519 that information is available at the larger map. All the maps
 520 are linked with synchronized navigation, and the use of small
 521 multiples allows the exploration of each temporal census
 522 individually, and its comparison to the others, with minimal
 523 interaction.

524 In this example, the maps show the transition from
 525 orange to green and purple in several regions over time.
 526 Clicking on a region in these maps will bring up a new
 527 panel with the original census data of this specific region.
 528 The actual population numbers are below the maps, and
 529 they confirm that the total population is indeed decreasing.

530 Each aspect is represented by a stacked bar plot, where
 531 the width of each rectangle corresponds to the average
 532 percentage of that variable. We chose stacked bar plots to
 533 represent the composition of the regions because they can
 534 accurately and succinctly inform the proportions of each
 535 aspect, without any interaction. In this case, about half of the
 536 people in Chicago in the considered period are married, and
 537 the percentage that are Widowers or Divorced is roughly
 538 similar. About half of the population work in Adminis-
 539 trative jobs, a third never completed high-school, approxi-
 540 mately half have gross family income below 50,000USD
 541 per year. The vast majority identify as white. Placing the
 542 mouse over one of the bars will open a small panel with
 543 the temporal evolution of that specific variable, and clicking
 544 on the chevron on the right side expands the corresponding
 545 aspect, showing details of the temporal evolution of each
 546 variable and also the corresponding IQRs for the whole city.

4 ILLUSTRATIVE SCENARIOS

547 We used decennial census data from the United States [4]
 548 and Canada¹, tabulated by CTs, from 1970 to 2010. The
 549 prototype allows access to 40 regions, 28 in the US and 12
 550 in Canada. Due to the high number of CTs, New York City
 551 was split into its boroughs.

552 We used five aspects for the USA: Education level, Fam-
 553 ily income, Marital status, Occupation, and Race; and seven
 554 for Canada: Age, Education level, Home language, House-
 555 hold Income, Marital status, Occupation, Place of birth, and
 556 Religion. While our method does not require geographic
 557 harmonization, it requires matching variables over time.
 558 The supplementary material contains the details of which
 559 census columns were used for each aspect. Income is slightly
 560 inaccurate, even though we did correct for inflation. We
 561 grouped the original ranges into three larger ranges, but
 562 they do not match precisely.

563 1. <http://datacentre.chass.utoronto.ca/census/>

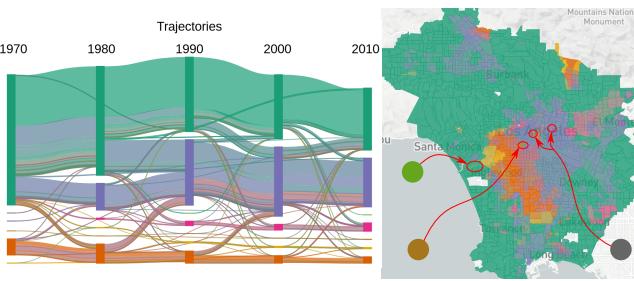


Fig. 5. Result for Los Angeles with 8 clusters, including three small and ephemeral clusters. Cluster characterization is displayed in Figure 6.

564 4.1 Los Angeles

565 We selected a region around the metropolitan area of Los
 566 Angeles (LA). This region follows urban density, not adminis-
 567 trative boundaries. The summary of the results using all
 568 aspects with equal weights and eight clusters is illustrated
 569 in Figure 5, while the statistical description of the clusters
 570 is illustrated in Figure 6, where the most relevant aspect
 571 of each cluster is highlighted. From the trajectories plot, we
 572 can see that there is a large but shrinking cluster, depicted in
 573 green, one increasing cluster in purple, an almost constant
 574 orange cluster, a smaller but increasing pink cluster, and
 575 three other small clusters. The corresponding map illus-
 576 trates where these clusters are located, and that they are somewhat
 577 geographically stable, with some movement between the
 578 green, orange, and purple clusters.

579 From Figure 6, we can see that the green cluster is
 580 characterized by a high percentage of White population,
 581 low percentage of population in the lowest income bracket,
 582 mostly administrative occupations, and about 30% of the
 583 population with four or more years of college. The orange
 584 cluster is characterized by a high percentage of Black pop-
 585 ulation, with few people in the highest range of income and
 586 education. The purple cluster corresponds to a high concen-
 587 tration of "Other" in race, which includes Hispanic for this
 588 dataset, high concentration of Laborers, and low education
 589 and income. The pink cluster contains a high percentage of
 590 Asian population and about 30% of the population with four
 591 or more years of college. The light green cluster contains
 592 very few people in the lower income bracket, mostly White
 593 population, with the highest percentage of population with
 594 four or more years of college, working administrative jobs,
 595 and a high concentration of singles. The yellow cluster
 596 represents Black population, with higher level of education
 597 and income, mostly working administrative jobs. The brown
 598 cluster represent a majority of single population, working
 599 administrative jobs with mostly low income. The dark gray
 600 cluster is characterized by all its population in the lowest
 601 income bracket, low education level, with a majority of
 602 White population.

603 Further, the larger clusters, green, orange, and purple,
 604 present a significant intra-cluster variance in most variables,
 605 as indicated by the maximum and minimum whiskers of the
 606 boxplots. For instance, while fifty percent of the CTs in the
 607 green cluster have between 20% and 40% of people in the
 608 lowest income bracket, that cluster also includes CTs where
 609 none and all the population belongs to that bracket. This
 610 might indicate that this cluster could be further divided, and

611 currently represents different groups of people that are not
 612 different enough to be separated into clusters in this level of
 613 the hierarchy.

614 While the larger clusters are expected, the light green,
 615 brown, and dark gray clusters are more surprising. These
 616 clusters are small and ephemeral, including only a few
 617 CTs, and are peculiarly different enough to be separated
 618 into their own clusters, at this level of the hierarchy. While
 619 these small clusters might be interesting enough to warrant
 620 further study, they might suppress larger, but more subtle,
 621 clusters from appearing in this interface, due to the limited
 622 number of clusters.

623 The orange area in the map in Figure 5 presents move-
 624 ment, indicated by the presence of green and purple tones
 625 mixed with the orange, which may warrant further explo-
 626 ration. By clicking on the orange bar above the boxplot,
 627 we select all trajectories that contain the orange cluster. The
 628 corresponding details are illustrated in Figure 7. This shows
 629 a location change, where the orange cluster is progressively
 630 replaced by the purple cluster on its east side, and in
 631 turn expanding to the west. Interestingly, the population
 632 increased, the racial profile changed, but the distribution of
 633 income was reasonably stable, with a higher amount of the
 634 population in the lowest income range and very few people
 635 in the highest income range. Indeed, the income difference
 636 is significant when compared to the city-wide distribution.

637 A portion of this region is classified into the green cluster
 638 in 2010, indicating a majority white population. To further
 639 understand that change, we clear the current selection, and
 640 select all regions that changed from orange in 1970 to green
 641 in 2010, using the transition matrix. A portion of the result-
 642 ing region, near the Florence-Graham region, is depicted in
 643 Figure 8, along with the temporal evolution of Race. Despite
 644 this difference, the other aspects are similar to the ones
 645 from the region in Figure 7, with slightly lower income
 646 and education profiles. While the racial aspect changed
 647 considerably, the economic and educational aspects stayed
 648 the same.

649 4.2 Toronto

650 We considered a region that is approximately the adminis-
 651 trative border of the current city of Toronto and all seven
 652 available aspects. The results are summarized in Figure 9,
 653 considering eight clusters.

654 The population with low percentage of University de-
 655 grees is represented in orange, mostly anglophone popula-
 656 tion in green, Asian immigrants in yellow, high percentage
 657 of income in the highest bracket in purple, high percentage
 658 of Jewish people in light green and brown, high percentage
 659 of Eastern Non-Christian religion in pink, and high concen-
 660 tration of single people in dark gray.

661 From the trajectories plot, we can see that Toronto is
 662 more dynamic than both Chicago and LA, with one cluster
 663 constantly shrinking, and others growing more prominent.
 664 In the 1970s, the city was divided into four clusters: low
 665 number of university degrees, Jewish population, majority
 666 anglophones, and high income, with the latter being rather
 667 small. There is some variance in the other characteristics of
 668 the orange cluster, but none as significant as the education
 669 level, especially when compared to the other years. While

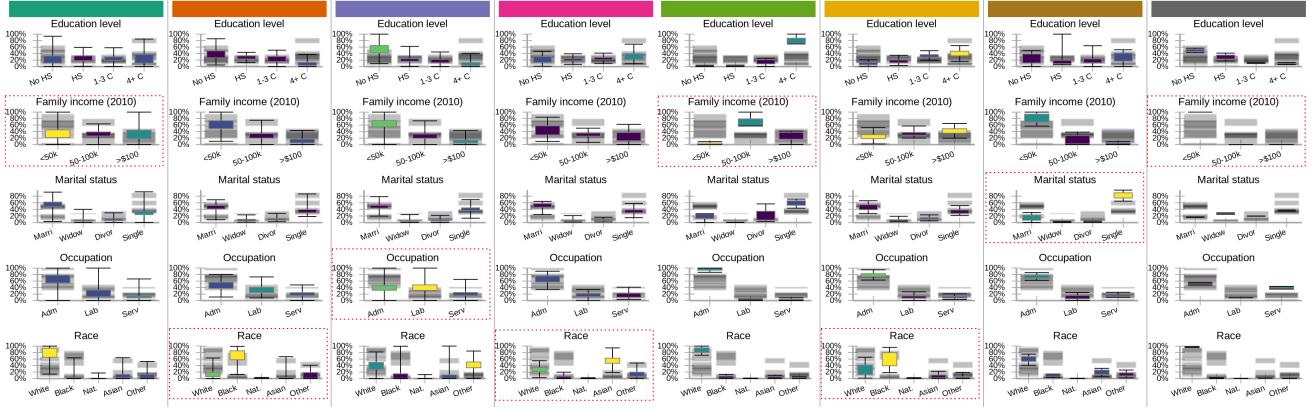


Fig. 6. Full characterization of the eight clusters found for LA. The red rectangles indicate the most relevant aspects for each cluster.

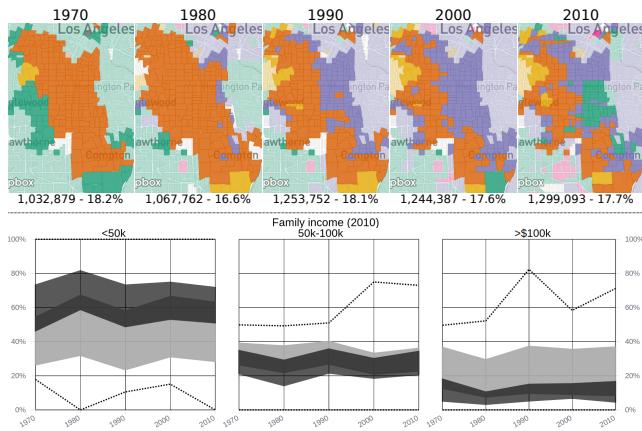


Fig. 7. Top: Geographic changes in the majority Black population cluster (orange) and Laborers cluster (green). Bottom: Income evolution for this region (black) and the whole city (gray).

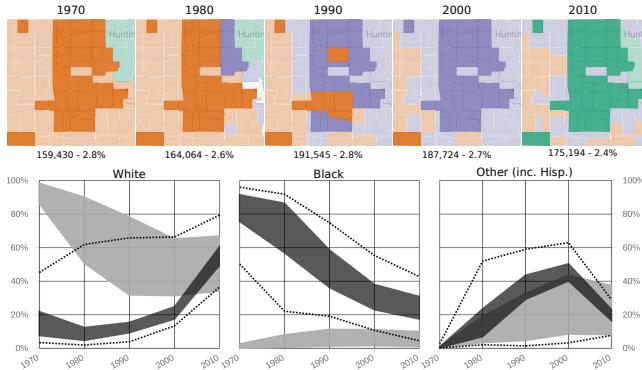


Fig. 8. Details for a volatile region contained in the area of Figure 7. This region went from Black to Hispanic to White.

the orange cluster persists until 2011, the other clusters present have similar education profiles to one another, leading the clustering method to use other aspects to characterize them.

Toronto's dynamism also appears in that we see the influx of Asian population starting in the 1980s and building thereafter, as represented by the yellow and pink clusters. While both include a high percentage of people born in Asia, the pink region is more defined by religion, with low

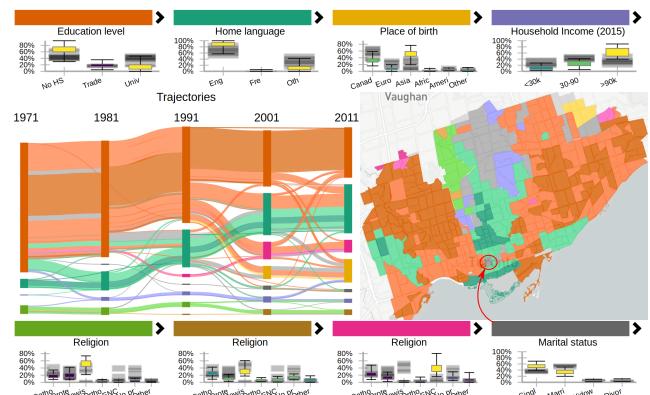


Fig. 9. Results for Toronto, with eight clusters.

percentage of university degrees, and contains the lowest percentage of people in the highest income bracket for these clusters; the yellow is less defined by religion, and has higher education and income, geographically corresponding to the Markham region, known for its Chinese population. A similar division also happens for the two Jewish clusters, where the light green cluster has lower education and income levels than the brown cluster. The purple cluster of high income is somewhat stable. This cluster includes until 2011 the Bridle Path neighborhood, known for its wealthy population. In 2011, however, it was classified into the yellow cluster that includes Asian immigration, since about 35% of the population for this CT were born in Asia. It maintained its income profile, with 85% with 90k CAD or more.

The grey regions of the map are another indication of Toronto's volatility over this period. While there is a cluster that is painted dark grey, annotated by a red circle in Figure 9, the light gray regions on the north part of the city indicate volatility. Here the same region was classified into three or more clusters for this five census period. Using the 'Add' mode for the trajectory selection, we select these volatile regions. A subset of the details is illustrated in Figure 10.

These regions account for about 5% of Toronto's population. The whole region was classified into the orange cluster in 1971 (low level of university degrees). By 1991,

670
680
681
682
683
684
685
686
687
688
689
690
691
692
693
694
695
696
697
698
699
700
701
702
703
704
705

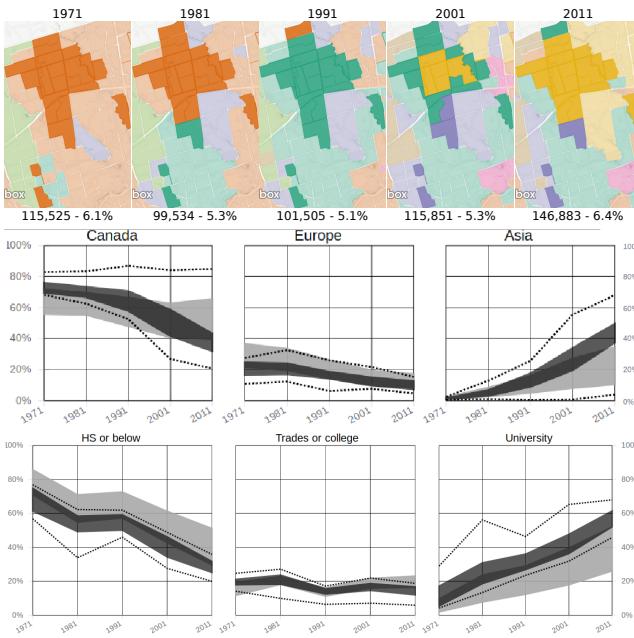


Fig. 10. Details for some regions of Toronto that were classified into 3 or more clusters over time.

most of the region was classified into the green cluster, representing anglophone population, mostly Canadian born, with a higher level of education. As the corresponding plot indicates, this trend in increasing education is city-wide, but this region has people with better education than most.

In 2001, the purple cluster of high income annexes neighboring parts of the volatile region, and the Asian born population increases sharply, as illustrated by the appearance of the yellow cluster. This cluster, once again, indicates well educated, higher income, and about 30%-50% Asian born population. By 2011, the yellow cluster increased considerably, annexing parts of the high income purple cluster, including the neighboring Bridle Path area.

While Toronto is more dynamic than Los Angeles, possibly due to size differences, the volatile regions shown in Figure 10 did not change as quickly or dramatically as the ones shown in Figure 8, which involved twice as many people. We found this trend to be related to the countries themselves, Canadian cities have larger areas undergoing slow, gradual changes, whereas American cities have more general stability, but quicker changes in smaller scales. The supplementary material contains brief summaries of all the regions accessible in this prototype.

5 EXPERT FEEDBACK

To assess the proposed method, we contacted academic and industry experts in sociology and urban sciences. The supplementary material contains all our communication with five experts, identified from A to E.

Their response was positive, mentioning that the prototype allows them to analyze census data without the additional work of obtaining and cleaning the data (A, B, E), and it allows the inclusion of geographic visual analysis tools in their research process (D). It enables the users to tell different stories about neighborhoods/cities and their

changes (A), visualize the relationship between key urban variables over time (D), offering a quick way to identify particular neighborhoods that one may be interested in studying more in depth around a particular issue or efficiently understanding the context of an area (E). Indeed, the experts identified gentrification processes in Manhattan (B) and Dallas (E), reinforced a hypothesis for occupational clustering (D), and highlighted how the method can be used to compare neighborhoods and cities (A). These remarks attest that we satisfied the design requirements regarding the existing demographic groups and their changes.

While the interface was "easy to navigate" (B), it was also considered "overwhelming" (A), "intimidating" (E), and "tricky to interpret" (C), possible side-effects of our effort to increase representational accuracy, where we avoided using simplified labels. Specifically, while identifying clusters by their most relevant variables was welcome, the overlap of information from different clusters in the boxplot was "a bit confusing" (C). While the map of trajectories was mentioned as a "good summary map", how it related to the clustering method was unclear (C). While we provide different options on how the colors are used on this map, both are sub optimal, reliably representing several distinct entities using colors is still an open problem in visualization.

The experts also mentioned the poor responsiveness of the method when changes in the clustering parameters required server-side processing (B,D). Indeed, the current implementation can take a few minutes to cluster regions with high number of CTs, like Los Angeles or Brooklyn.

Most of the experts demonstrated interest in using our method in their research (A, B, D, E), aiming to use the census data as a backdrop for other datasets, providing demographic context. They also mentioned the need to export subsets of data, plots, and maps to be used in reports and publications (C, D, E).

More importantly, while these experts were aware that our method does not perform geographic normalization, none of them mention it. We did not specifically ask if this difference led to unexpected results, but rather inquired if they found interesting insights, leading them to consider the interface, its options, and results as a whole. Therefore the method was neutrally compared against other, traditionally normalized results. We interpret the fact that most of them were interested in the next steps as an affirmation of the accuracy of the method.

6 DISCUSSION AND LIMITATIONS

The objective of this work was to leverage a graph based data representation with visual tools to allow for the exploration of geographically inconsistent census data. While we successfully replicated and corroborated results from the literature, this method still has significant limitations.

Removing the geographical normalization/interpolation step greatly reduces the amount of work necessary, but the method still requires consistent variables across the years. Matching the fields of the public census can be trivial for some aspects, like Age, but challenging for others, like Income. The divulged income ranges vary over time, the actual value changes due to inflation and other factors, and so on. Moreover, some fields were not considered in earlier

799 censuses, such as Race in Canada, or Hispanic population
 800 in the USA, hampering its use when they are available. We
 801 matched some aspects, but a deeper demographic analysis
 802 would greatly benefit from all available information.

803 Another limitation is the lack of control on how much the
 804 geographical information will impact the clustering result.
 805 While the adopted method met our needs for this work,
 806 a configurable control would add another dimension to
 807 the exploration, allowing for more intra-cluster variance
 808 to obtain more 'compact' clusters. We explored changing
 809 the number of content based augmented edges, but this
 810 proved to be unreliable and hard to interpret. The *ClustGeo*
 811 method [33] can be a viable option for this, allowing a graph
 812 based input and a hierarchical output, combined using a
 813 single mixing parameter. Alternatively, one could cluster the
 814 changes [55] instead of the stable states.

815 There are also technological limitations, such as memory
 816 use on the visualization client. To allow for changes on the
 817 CTs over the years, we use a geographic file that contains
 818 all possible intersections, which can grow rather large if the
 819 original city was expansive and contained several CTs, like
 820 NYC or LA. However, the most significant technological
 821 limitation relates to parameters that are not immediately
 822 interactive, such as the clustering configuration. Since the
 823 clustering is computationally expensive and performed on
 824 the server, which allows for cached results, some changes
 825 can take a few minutes to be considered, removing any
 826 possibility of a continuous exploration.

827 Indeed, the cognitive load on the user is already sig-
 828 nificant, as we compromised simplicity for accuracy. While
 829 other works labelled the clusters, as 'young urban', 'strug-
 830 gling', and so on [29], [30], we show the statistical char-
 831 acteristics of the clusters, which are harder to interpret, as the
 832 data may have subtle nuances that labels would otherwise
 833 hide. This also led to a crowded interface, mitigated some-
 834 what the use of pop-up panels and collapsible sections. For
 835 some cities, especially if they are small and stable, the panels
 836 can appear redundant, but each provide a different way to
 837 interact with the information that can ease the exploration
 838 process for larger and dynamic cities.

839 7 CONCLUSION

840 Our objective was to allow for the exploration of census data
 841 without geographical harmonization, an approach that is
 842 usually not considered. Our method was able to corroborate
 843 previous findings from the specialized literature, with an
 844 increased level of detail due to our data representation
 845 and visualization choices. The feedback from experts was
 846 remarkably positive, most of them were able to extract
 847 insight from the prototype and demonstrated interest in
 848 using it on their research efforts.

849 While we focused on census data in this work, the
 850 proposed method can be used with most similarly area-
 851 based datasets. Further, while there is interest in the use
 852 of census data, it is mostly in combination with other, larger
 853 and more dynamic, datasets. We postulate that, while census
 854 data is interesting, the amount of effort required to obtain
 855 and prepare it often discourages its use as complementary
 856 information, reinforcing our hypothesis that there is need

for methods that are able to deal with geographically inconsis-
 857 tent data.

858 "... tract-by-tract comparison is not possible unless data from
 859 2000 is interpolated to 2010 boundaries (...)" [10]

861 ACKNOWLEDGMENTS

862 This research was supported by a University of Toronto
 863 Connaught Global Challenge grant and is part of the Urban
 864 Genome Project.

865 REFERENCES

- [1] W. Chen, Z. Huang, F. Wu, M. Zhu, H. Guan, and R. Maciejewski, "Vaud: A visual analysis approach for exploring spatio-temporal urban data," *IEEE Transactions on Visualization & Computer Graphics*, no. 1, pp. 1–1, 2017.
- [2] T. Von Landesberger, F. Brodkorb, P. Roskosch, N. Andrienko, G. Andrienko, and A. Kerren, "Mobilitygraphs: Visual analysis of mass mobility dynamics via spatio-temporal graphs and clustering," *IEEE transactions on visualization and computer graphics*, vol. 22, no. 1, pp. 11–20, 2016.
- [3] W. Chen, F. Guo, and F.-Y. Wang, "A survey of traffic data visualization," *IEEE Transactions on Intelligent Transportation Systems*, vol. 16, no. 6, pp. 2970–2984, 2015.
- [4] S. Manson, J. Schroeder, D. V. Riper, and S. Ruggles. (2017) Ipums national historical geographic information system: Version 12.0 [database]. Minneapolis: University of Minnesota. [Online]. Available: <http://doi.org/10.18128/D050.V12.0>
- [5] J. R. Logan, Z. Xu, and B. J. Stults, "Interpolating us decennial census tract data from as early as 1970 to 2010: A longitudinal tract database," *The Professional Geographer*, vol. 66, no. 3, pp. 412–420, 2014.
- [6] E. Hallisey, E. Tai, A. Berens, G. Wilt, L. Peipins, B. Lewis, S. Graham, B. Flanagan, and N. B. Lunsford, "Transforming geographic scale: a comparison of combined population and areal weighting to other interpolation methods," *International Journal of Health Geographics*, vol. 16, no. 1, p. 29, Aug 2017.
- [7] J. Allen and Z. Taylor, "A new tool for neighbourhood change research: The canadian longitudinal census tract database, 19712016," *The Canadian Geographer / Le Géographe canadien*, vol. 0, no. 0, 2018. [Online]. Available: <https://onlinelibrary.wiley.com/doi/abs/10.1111/cag.12467>
- [8] J. R. Logan, B. J. Stults, and Z. Xu, "Validating population estimates for harmonized census tract data, 2000–2010," *Annals of the American Association of Geographers*, vol. 106, no. 5, pp. 1013–1029, 2016.
- [9] C. L. Eicher and C. A. Brewer, "Dasymetric mapping and areal interpolation: Implementation and evaluation," *Cartography and Geographic Information Science*, vol. 28, no. 2, pp. 125–138, 2001.
- [10] A. Dmowska, T. F. Stepinski, and P. Netzel, "Comprehensive framework for visualizing and analyzing spatio-temporal dynamics of racial diversity in the entire united states," *PLOS ONE*, vol. 12, no. 3, pp. 1–20, 03 2017.
- [11] G. Firebaugh and C. R. Farrell, "Still large, but narrowing: The sizable decline in racial neighborhood inequality in metropolitan america, 1980–2010," *Demography*, vol. 53, no. 1, pp. 139–164, Feb 2016. [Online]. Available: <https://doi.org/10.1007/s13524-015-0447-5>
- [12] A. V. Diez-Roux, F. J. Nieto, C. Muntaner, H. A. Tyroler, G. W. Comstock, E. Shahar, L. S. Cooper, R. L. Watson, and M. Szklo, "Neighborhood environments and coronary heart disease: a multilevel analysis," *American journal of epidemiology*, vol. 146, no. 1, pp. 48–63, 1997.
- [13] C. A. Gotway and L. J. Young, "Combining incompatible spatial data," *Journal of the American Statistical Association*, vol. 97, no. 458, pp. 632–648, 2002.
- [14] X. Liu, Y. Song, K. Wu, J. Wang, D. Li, and Y. Long, "Understanding urban china with open data," *Cities*, vol. 47, pp. 53 – 61, 2015, current Research on Cities (CRoC).
- [15] A. C.-D. Lee and C. Rinner, "Visualizing urban social change with self-organizing maps: Toronto neighbourhoods, 1996–2006," *Habitat International*, vol. 45, pp. 92–98, 2015.

- [16] A. Dmowska and T. F. Stepinski, "Spatial approach to analyzing dynamics of racial diversity in large u.s. cities: 199020002010," *Computers, Environment and Urban Systems*, vol. 68, pp. 89 – 96, 2018. [Online]. Available: <http://www.sciencedirect.com/science/article/pii/S019897151730371X>
- [17] H. Doraiswamy, N. Ferreira, T. Damoulas, J. Freire, and C. T. Silva, "Using topological analysis to support event-guided exploration in urban data," *IEEE Transactions on Visualization and Computer Graphics*, vol. 20, no. 12, pp. 2634–2643, Dec 2014.
- [18] D. I. Shuman, S. K. Narang, P. Frossard, A. Ortega, and P. Vandergheynst, "The emerging field of signal processing on graphs: Extending high-dimensional data analysis to networks and other irregular domains," *IEEE Signal Processing Magazine*, vol. 30, no. 3, pp. 83–98, 2013.
- [19] A. Sandryhaila and J. M. Moura, "Discrete signal processing on graphs," *IEEE transactions on signal processing*, vol. 61, no. 7, pp. 1644–1656, 2013.
- [20] P. Valdivia, F. Dias, F. Petronetto, C. T. Silva, and L. G. Nonato, "Wavelet-based visualization of time-varying data on graphs," in *2015 IEEE Conference on Visual Analytics Science and Technology (VAST)*, Oct 2015, pp. 1–8.
- [21] F. Dias and L. G. Nonato, "Some operators from mathematical morphology for the visual analysis of georeferenced data," in *Workshop on Visual Analytics, Information Visualization and Scientific Visualization - SIBGRAPI*, 2015.
- [22] A. Dal Col, P. Valdivia, F. Petronetto, F. Dias, C. T. Silva, and L. G. Nonato, "Wavelet-based visual analysis of dynamic networks," *IEEE Transactions on Visualization and Computer Graphics*, vol. PP, no. 99, pp. 1–1, 2018.
- [23] X. Huang, Y. Zhao, C. Ma, J. Yang, X. Ye, and C. Zhang, "Trajgraph: A graph-based visual analytics approach to studying urban network centralities using taxi trajectory data," *IEEE Transactions on Visualization and Computer Graphics*, vol. 22, no. 1, pp. 160–169, Jan 2016.
- [24] C. Vehlow, F. Beck, and D. Weiskopf, "The State of the Art in Visualizing Group Structures in Graphs," in *Eurographics Conference on Visualization (EuroVis) - STARS*, R. Borgo, F. Ganovelli, and I. Viola, Eds. The Eurographics Association, 2015.
- [25] F. Beck, M. Burch, S. Diehl, and D. Weiskopf, "The State of the Art in Visualizing Dynamic Graphs," in *EuroVis - STARS*, R. Borgo, R. Maciejewski, and I. Viola, Eds. The Eurographics Association, 2014.
- [26] T. Setiadi, A. Pranolo, M. Aziz, S. Mardiyantri, B. Hendrajaya, and Munir, "A model of geographic information system using graph clustering methods," in *2017 3rd International Conference on Science in Information Technology (ICSI Tech)*, Oct 2017, pp. 727–731.
- [27] A. Fahad, N. Alshatri, Z. Tari, A. Alamri, I. Khalil, A. Y. Zomaya, S. Foufou, and A. Bouras, "A survey of clustering algorithms for big data: Taxonomy and empirical analysis," *IEEE Transactions on Emerging Topics in Computing*, vol. 2, no. 3, pp. 267–279, sep 2014. [Online]. Available: <https://doi.org/10.1109/tetc.2014.2330519>
- [28] A. K. Jain, "Data clustering: 50 years beyond k-means," *Pattern recognition letters*, vol. 31, no. 8, pp. 651–666, 2010.
- [29] E. C. Delmelle, "Mapping the dna of urban neighborhoods: Clustering longitudinal sequences of neighborhood socioeconomic change," *Annals of the American Association of Geographers*, vol. 106, no. 1, pp. 36–56, 2016.
- [30] ——, "Differentiating pathways of neighborhood change in 50 u.s. metropolitan areas," *Environment and Planning A: Economy and Space*, vol. 49, no. 10, pp. 2402–2424, 2017.
- [31] C. Ling and E. C. Delmelle, "Classifying multidimensional trajectories of neighbourhood change: a self-organizing map and k-means approach," *Annals of GIS*, vol. 22, no. 3, pp. 173–186, 2016.
- [32] J. Han, M. Kamber, and A. K. Tung, "Spatial clustering methods in data mining," *Geographic data mining and knowledge discovery*, pp. 188–217, 2001.
- [33] M. Chavent, V. Kuentz-Simonet, A. Labenne, and J. Saracco, "Clustgeo: an r package for hierarchical clustering with spatial constraints," *Computational Statistics*, pp. 1–24, 2017.
- [34] S. Soor, A. Challa, S. Danda, B. D. Sagar, and L. Najman, "Extending k-means to preserve spatial connectivity," 2018.
- [35] P. Soille and L. Najman, "On morphological hierarchical representations for image processing and spatial data clustering," in *Applications of Discrete Geometry and Mathematical Morphology*. Springer, 2012, pp. 43–67.
- [36] J. Cousty, G. Bertrand, L. Najman, and M. Couperie, "Watershed cuts: Minimum spanning forests and the drop of water principle," *IEEE Transactions on Pattern Analysis and Machine Intelligence*, vol. 31, no. 8, pp. 1362–1374, Aug 2009.
- [37] M. D. Dias, M. R. Mansour, F. Dias, F. Petronetto, C. T. Silva, and L. G. Nonato, "A hierarchical network simplification via non-negative matrix factorization," in *2017 30th SIBGRAPI Conference on Graphics, Patterns and Images (SIBGRAPI)*, Oct 2017, pp. 119–126.
- [38] M. Monmonier, "Strategies for the visualization of geographic time-series data," *Cartographica: The International Journal for Geographic Information and Geovisualization*, vol. 27, no. 1, pp. 30–45, 1990.
- [39] N. Andrienko, G. Andrienko, and P. Gatalsky, "Exploratory spatio-temporal visualization: an analytical review," *Journal of Visual Languages & Computing*, vol. 14, no. 6, pp. 503–541, 2003.
- [40] N. Ferreira, "Visual analytics techniques for exploration of spatiotemporal data," Ph.D. dissertation, Polytechnic Institute of New York University, 2015.
- [41] Y. Zheng, W. Wu, Y. Chen, H. Qu, and L. M. Ni, "Visual analytics in urban computing: An overview," *IEEE Transactions on Big Data*, vol. 2, no. 3, pp. 276–296, Sept 2016.
- [42] M. O. Ward, G. Grinstein, and D. Keim, *Interactive data visualization: foundations, techniques, and applications*. AK Peters/CRC Press, 2015.
- [43] G. Andrienko, N. Andrienko, H. Schumann, and C. Tominski, "Visualization of trajectory attributes in space-time cube and trajectory wall," in *Cartography from Pole to Pole*. Springer, 2014, pp. 157–163.
- [44] C. Tominski and H.-J. Schulz, "The Great Wall of Space-Time," in *Vision, Modeling and Visualization*, M. Goesele, T. Grosch, H. Theisel, K. Toennies, and B. Preim, Eds. The Eurographics Association, 2012.
- [45] S. Buschmann, M. Trapp, and J. Döllner, "Real-time animated visualization of massive air-traffic trajectories," in *Cyberworlds (CW), 2014 International Conference on*. IEEE, 2014, pp. 174–181.
- [46] D. Seebacher, J. Häufner, M. Hundt, M. Stein, H. Müller, U. Engelke, and D. Keim, "Visual analysis of spatio-temporal event predictions: Investigating the spread dynamics of invasive species," in *2017 IEEE Visualization Conference (VIS)*, 2017.
- [47] G. Andrienko, N. Andrienko, G. Fuchs, and J. Wood, "Revealing patterns and trends of mass mobility through spatial and temporal abstraction of origin-destination movement data," *IEEE Transactions on Visualization and Computer Graphics*, vol. 23, no. 9, pp. 2120–2136, Sept 2017.
- [48] N. Ferreira, M. Lage, H. Doraiswamy, H. Vo, L. Wilson, H. Werner, M. Park, and C. Silva, "Urbane: A 3d framework to support data driven decision making in urban development," in *Visual Analytics Science and Technology (VAST), 2015 IEEE Conference on*. IEEE, 2015, pp. 97–104.
- [49] M. Li, Z. Bao, T. Sellis, S. Yan, and R. Zhang, "Homeseeker: A visual analytics system of real estate data," *Journal of Visual Languages & Computing*, vol. 45, pp. 1 – 16, 2018.
- [50] J. Johansson, P. Ljung, M. Jern, and M. Cooper, "Revealing structure within clustered parallel coordinates displays," in *Information Visualization, 2005. INFOVIS 2005. IEEE Symposium on*. IEEE, 2005, pp. 125–132.
- [51] D. Guo, J. Chen, A. M. MacEachren, and K. Liao, "A visualization system for space-time and multivariate patterns (vis-stamp)," *IEEE transactions on visualization and computer graphics*, vol. 12, no. 6, pp. 1461–1474, 2006.
- [52] F. Pedregosa, G. Varoquaux, A. Gramfort, V. Michel, B. Thirion, O. Grisel, M. Blondel, P. Prettenhofer, R. Weiss, V. Dubourg, J. Vanderplas, A. Passos, D. Cournapeau, M. Brucher, M. Perrot, and E. Duchesnay, "Scikit-learn: Machine learning in Python," *Journal of Machine Learning Research*, vol. 12, pp. 2825–2830, 2011.
- [53] M. Harrower and C. A. Brewer, "Colorbrewer.org: An online tool for selecting colour schemes for maps," *The Cartographic Journal*, vol. 40, no. 1, pp. 27–37, 2003.
- [54] J. L. Hintze and R. D. Nelson, "Violin plots: a box plot-density trace synergism," *The American Statistician*, vol. 52, no. 2, pp. 181–184, 1998.
- [55] J. Bian, D. Tian, Y. Tang, and D. Tao, "A survey on trajectory clustering analysis," *arXiv preprint arXiv:1802.06971*, 2018.