

Unharmonised longitudinal data processing using networks

1

June 13, 2019

2

3

Abstract

4

The main contribution of this work is the use of a simple network representation to allow
the processing of spatio-temporal demographic data for the identification of data-driven neigh-
bourhoods and their dynamics, *without the need for geographical harmonisation*, a laborious and
error-prone process that is currently used in virtually all longitudinal analysis of region-based
data.

5

6

7

8

9

To allow for a transparent corroboration of our method, we leverage it as basis for an
interactive and intuitive interface that allows a progressive exploration of the results, from the
characterisation of broad patterns to the identification of individual details and direct access to
the original data, including to non-experts.

10

11

12

13

We validate our method with illustrative scenarios for Chicago, Toronto, and Los Angeles,
with results that match recent literature. The system is publicly available, with data from the
decennial censuses for over forty regions in the USA and Canada between 1970 and 2010, but
the methodology is suitable for any region-based data.

14

15

16

17

1 Introduction

18

Neighbourhoods have increasingly become a central concept in social research and targets for so-
cial policy (Sampson, 2012; Galster, 2019; Stone et al., 2015; Looker, 2015). To be sure, a focus
on neighbourhoods extends to the formative period of the modern social sciences (Abbott, 1997).
Recent interest has at least partly been rekindled through newly available longitudinal demographic
datasets (Logan et al., 2014; Manson et al., 2017), convenient computational tools (Rey et al., 2018),
and new sources of data (Poorthuis, 2018).

19

20

21

22

23

24

25 These advances have made neighbourhood constructs more tractable to quantitative analysis.
26 Yet new challenges have also emerged, especially at the convergence of research on neighbourhood
27 effects and neighbourhood dynamics. Neighbourhood effects research assumes knowledge about the
28 nature and scope of "the neighbourhood" that presumably shapes individual outcomes (Kwan, 2018;
29 Shelton and Poorthuis, 2019). Concurrently, researchers note that neighbourhoods are not neces-
30 sarily fixed containers in which other processes occur, but themselves dynamically evolve (Delmelle,
31 2017; Reades et al., 2019; Li and Xie, 2018). The result is to open up key assumptions about neigh-
32 bourhoods for theoretical and empirical examination: how do we appropriately define and compare
33 neighbourhoods at a given time?; how do we appropriately define and compare the temporal tra-
34 jectories of neighbourhoods?; and can we do both at once, "fully interactionally" (Abbott, 1997):
35 classify neighbourhoods now based on where they came from and where they are going?

36 In principle, much of the recent research is committed to the proposition that neighbourhoods
37 are open and evolving entities. Ironically, its empirical practice tends to rely on methods that require
38 fixed geographical regions. This requirement is difficult to satisfy, as most longitudinal datasets are
39 based on pre-defined tabulation areas that are routinely modified by data collection agencies, usually
40 to follow population changes.

41 The standard approach then is to *geographically harmonise* data. This involves interpolating
42 existing measurements into a common set of regions (Logan et al., 2014; Hallisey et al., 2017;
43 Allen and Taylor, 2018). Recent computational tools have somewhat simplified this process (Rey
44 et al., 2018), but it still involves non-trivial questions: which geometry to use as target, how to
45 apportion the variables, or how to combine data from different sources. Further, these question
46 do not necessarily have optimal answers. Indeed, regardless of how well this process is performed,
47 it still introduces errors (Logan et al., 2016), even when additional data is provided (Eicher and
48 Brewer, 2001). Essentially, harmonisation generates artificial data points that can potentially lead
49 to inaccurate results, even though they are seldom interpreted as such. Nevertheless, because there
50 has been no viable alternative, and the results often appear plausible, these concerns are generally
51 overlooked. The result is that the harmonisation approach is virtually mandatory in the current
52 literature: "(...) *tract-by-tract comparison is not possible unless data from 2000 is interpolated to*
53 *2010 boundaries (...)"* (Dmowska et al., 2017), "(...) *This limits cross-year comparison since data*
54 *are not representative of the same spatial units. (...)"* (Allen and Taylor, 2018).

55 The main contribution of this paper is a method for longitudinal data processing that works with

the original data by leveraging a network based representation. It enables tract-by-tract comparison
and the identification of patterns of demographic evolution *without geographic harmonisation.*

To allow a proper examination of our method and its results, we built an online interactive
system using this representation. It enables users to visualise, interpret, and explore trajectories of
neighbourhood change. This interface helps validate our method, by allowing it to be compared to
existing and future methods. Further, it is a significant contribution to the research community: it
provides a vehicle for quickly and easily grasping complex long-term changes, experimenting with
different parameters to interactively learn from data, and making neighbourhood change research
publicly transparent. The interface thus responds to increasing concerns about reproducibility and
transparency, as well as ongoing attention to the value of visualisation in scientific research and
communication.

We start by presenting an intuitive example of our representation in Section 1, then we review
the relevant literature on longitudinal studies, data representation, clustering, and spatio-temporal
visualisation in Section 3. Our methodology is introduced in detail in Section 4, along with the
included interface. Illustrative scenarios for Chicago, Toronto, and Los Angeles are presented in
Section 5 and the feedback of five field experts are summarised in Section 6. Our prototype system
is available at , including more than forty regions in the US and Canada. The source code is publicly
available at .¹

2 Intuition

While utterly simple, the network model breaks from the deeply rooted traditional tabular paradigm
in a significant way. Instead of requiring the data as a collection of fixed entities which properties
evolve over time, literally rows in a table with temporal values as columns, it represents each measure-
ment as a separate entity and encodes the evolution of these entities over time.

To ease this cognitive transition, we start with an intuitive example of how it works, using a small
portion of a fictitious urban region illustrated on the left part of Figure 1. This example includes
three different times (t_0, t_1, t_2), with different aggregation areas identified as letters from A to H.
For t_0 , the initial time, we have areas A and B, with small houses and a park, respectively. The
park remains stable (B, E, and H), but the houses are partially replaced by larger buildings (C and

¹The editors are considering, at our request, an exception to the double-blind requirement to allow access to the system. We provided them with the URLs of the system, code, and documentation separately.

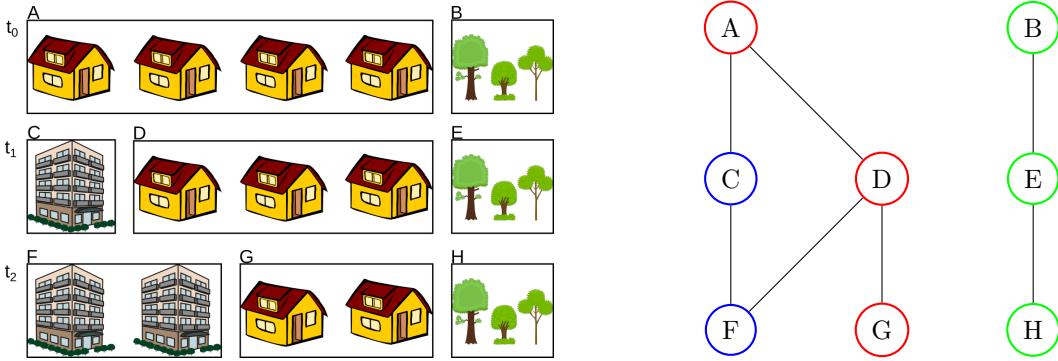


Figure 1: Network based spatio-temporal data representation. **Left:** Three temporal stages of the evolution of a fictitious urban area, with aggregation areas A to H. **Right:** Network representation of the aggregation areas where the colours identify similar regions.

84 F).

85 This example also illustrates some of the challenges of harmonisation. None of these aggregation
 86 areas is clearly suitable as an interpolation target. In fact, adopting any of these areas as a target
 87 would require merging heterogeneous regions and/or dividing homogeneous regions. For instance,
 88 by choosing the regions of t_1 , A would be split to match regions C and D, which appears to be a
 89 rather reasonable approximation in this homogeneous artificial example, albeit the fallacy of division.
 90 Region F would be similarly split to match C, but F would be split and merged with G to match
 91 D, potentially leading to statistical measurements that do not properly represent either region.

92 Real world measurements are seldom as homogeneous and noise-free as this artificial example.
 93 By splitting and merging the data to fit arbitrary borders, that were not necessarily coherent at the
 94 time of the measurement, harmonisation increases the distance between measurement and reality.

95 Instead, we propose a network representation. Intuitively, a *network* (also called a *graph*) is
 96 a collection of entities (nodes) that are related to each other (edges). In this case, each different
 97 aggregation area is represented as a node and we connect nodes that have overlapping geographical
 98 areas, leading to the network illustrated on the right of Figure 1. By partitioning the network into
 99 connected nodes that are similar, we are effectively identifying clusters in the spatio-temporal data,
 100 as illustrated by the colours of the nodes on the right side of Figure 1. Further, the possible evolution
 101 paths can also be obtained by computing temporal sequences of nodes, in this case: (A, C, F), (A,
 102 D, F), (A, D, G), and (B, E, H). This representation is also suited for geographically consistent
 103 regions, as illustrated by the stable park in this example, and is therefore a generalisation of the

traditional paradigm. 104

Note that the edges of this network merely encode that two regions are related. This is a 105
binary information, there is no apportionment, no areal measurements, no population percentages 106
associated with the edge. Indeed, our method also connects regions of the same time that share 107
borders, representing exactly that they are neighbouring areas. 108

3 Related Work 109

Since our problem encompasses several fields, we divide this section into specific sub problems: 110
longitudinal demographic studies, describing the traditional tabular approach to longitudinal studies; 111
data representation, elaborating how evolving geographic data can be represented for processing; *data 112
clustering*, briefly reviewing existing clustering methods; and *cluster characterisation*, articulating 113
how clusters can be visually summarised. 114

3.1 Longitudinal demographic studies 115

Census data is used not only to discover demographic patterns (Firebaugh and Farrell, 2016), but 116
to correlate demographic characteristics to other measurements (Diez-Roux et al., 1997). However, 117
longitudinal studies are rare, because they are difficult : ”(...) *One of the most challenging and 118
fascinating areas in spatial statistics is the synthesis of spatial data collected at different spatial 119
scales(...)*” (Gotway and Young, 2002). While census tract level data is readily available for the US 120
since at least 1910 (Manson et al., 2017), most studies consider the period between 1970 and 2010, 121
using pre-harmonised data from the Longitudinal Tract Data Base (Logan et al., 2014). Despite 122
its inherent errors (Logan et al., 2016; Hallisey et al., 2017), this dataset has become the standard 123
source for longitudinal demographic data at the neighbourhood scale, with similar efforts appearing 124
in other countries (Liu et al., 2015; Lee and Rinner, 2015; Allen and Taylor, 2018). These datasets 125
have been highly significant for the field. Yet they also limit the universe of data that can be used to 126
study neighbourhood change, since any new datasets would need to be similarly processed in order 127
to be rendered compatible with these sources. 128

Another option considers the use of grid data (Dmowska et al., 2017; Dmowska and Stepinski, 129
2018; Stepinski and Dmowska, 2019). Beyond the increased spatial accuracy, this approach does not 130
require complex harmonisation when new data is considered, if the grids are compatible. However, 131

132 demographic data is usually not available in this format, especially from older sources. Additionally,
133 the conversion from tabulation areas can introduce significant errors.

134 Given these challenges, it is worth considering new alternatives. In this work, we propose a
135 novel methodology that entirely avoids the problems of geographical harmonisation, considering
136 each measurement using its actual geographic region. It does not require regions to be consistent
137 across time because they are naturally represented as different entities.

138 3.2 Data representation

139 Network based representation of geographic information is fairly well explored in the literature,
140 as a basis for topological methods for event detection (Doraiswamy et al., 2014), leveraging signal
141 processing on graphs (Shuman et al., 2013; Sandryhaila and Moura, 2013) to find patterns and
142 outliers (Valdivia et al., 2015; Dias and Nonato, 2015; Dal Col et al., 2018). Networks are well
143 suited to represent trajectories as well (Von Landesberger et al., 2016; Huang et al., 2016; Chen
144 et al., 2015), allowing the use of graph visualisation methods (Vehlow et al., 2015; Beck et al., 2014).
145 Our proposed method builds upon this literature. We leverage a network-based representation
146 that removes the requirement for consistency in the measurement regions. Each region in time
147 corresponds to a different node. Instead of a collection of time-series, the data is represented as a
148 dynamic network.

149 Networks have been used to represent census data for clustering purposes (Dias and Nonato,
150 2015; Setiadi et al., 2017), but these works did not explore temporal evolution, where they are
151 particularly powerful. Networks allow a natural representation of these inconsistent regions, with
152 both spatial and temporal connections. There are other possible representations that have similar
153 properties, but we adopted networks to allow the use of the vast existing literature and methods.

154 3.3 Data clustering

155 Data clustering is one of the elementary processes for data analysis, simplifying the data into a
156 smaller number of homogeneous sets that can be interpreted in the same way. There is no shortage of
157 contributions for this problem (Fahad et al., 2014), but most neighbourhood related applications still
158 rely on k-means (Jain, 2010; Delmelle, 2016) and, to a lesser extent, Self Organising Maps (Delmelle,
159 2017; Ling and Delmelle, 2016).

However, a method for geographic data analysis should not ignore the geographic component of the data, and we extend research exploring ways to incorporate it. One straightforward option for agglomerative methods (Han et al., 2001) is to consider only nearby clusters for merging (Chavent et al., 2017), which can also be done for k-means (Soor et al., 2018). Alternatively, the spatial distance could be directly added to the inter-cluster metric (Chavent et al., 2017) via a mixing parameter, which adds flexibility to the method, but introduces the problem of finding the correct application-dependent values.

Indeed, one crucial step in most clustering algorithms is the definition of the number of clusters. We avoid this problem by considering hierarchical methods (Soille and Najman, 2012), where the result is not a partition of the data, but a tree of partitions, similar to a dendrogram. This approach is interesting for interactive methods, because it allows the user to change the number of clusters on the fly. Since our data is represented as a network, we opted for an heuristic variation of the maximum weighted matching algorithm called *sorted maximal matching* (Dias et al., 2017), which merges clusters based on the weights of the edges between pairs of clusters.

3.4 Cluster characterisation

While visualisation has gained prominence as a crucial component of scientific discovery, justification, and communication Tufte et al. (1998), visually representing evolving spatial data is a challenging old problem (Monmonier, 1990; Andrienko et al., 2003; Ferreira, 2015; Zheng et al., 2016).

Most geographic data is naturally bidimensional and maps work well in this case (Zheng et al., 2016; Ward et al., 2015), but the additional temporal dimension cannot be so naturally represented. One straightforward option is to leverage tridimensional plots (Andrienko et al., 2014; Tominski and Schulz, 2012), but this can lead to visual obstructions or scaling problems unless a tridimensional display device is used. A simpler, well adopted, option is to display a map that corresponds to a subset of the temporal information, allowing the user to change the time with an associated control (Chen et al., 2017; Valdivia et al., 2015; Dal Col et al., 2018; Doraiswamy et al., 2014). Small multiples can be used (Von Landesberger et al., 2016), but only when there are few temporal snapshots. However, none of these options is suitable to represent many variables at the same time.

Using data clustering, we can represent the region's cluster instead of all the its variables (Dal Col et al., 2018; Valdivia et al., 2015; Von Landesberger et al., 2016). While this simplifies the geographic portion of the visualization, it introduces the problem of how to summarise the contents of each

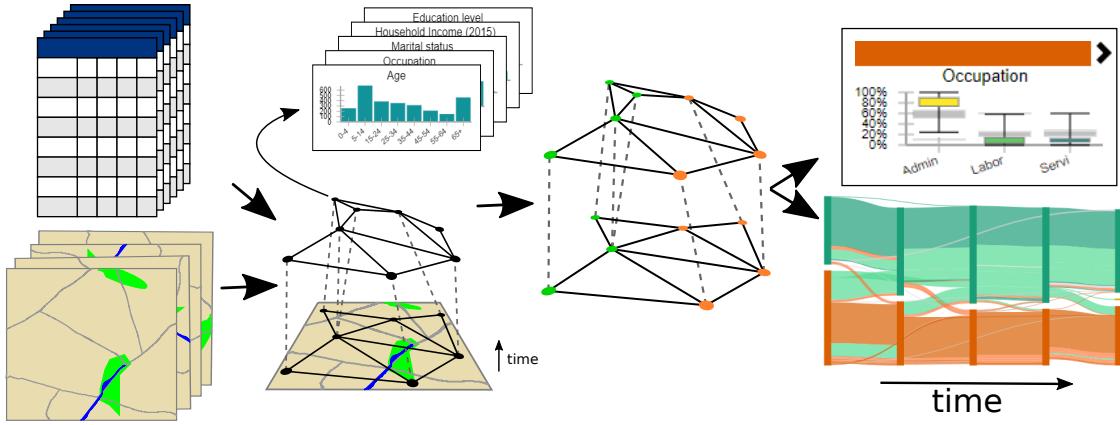


Figure 2: Overview of the proposed method. A network is generated combining the original census data, encoding the changing geographical information. The network is partitioned into an hierarchy (Dias et al., 2017). The characteristics and evolution of the clusters are then visually represented.

cluster. One traditional approach is to use parallel coordinates plot (Ferreira et al., 2015), but these they can get cluttered representing similar clusters over several variables. Further, for demographic applications, the clusters are usually strongly characterised by a small subset of values (Delmelle, 2016, 2017). Therefore, in the proposed method, we identify the variables that are most relevant to the characterisation of each cluster. The distribution of values on that variable is then represented using a boxplot, a well known statistical plot displaying basic properties of the distributions.

4 Visualising the demographic spatio-temporal evolution

Figure 2 presents an overview of the processing steps of the proposed method, illustrating how the nodes of the network are used to represent the regions. The following sections elaborate this figure and explain the main features of the interface we built to visualise and explore the evolution of neighbourhoods on the basis of our proposed method.

4.1 Census methodology and data representation

Census data is disseminated in a tabulated form for aggregation areas: whole country, state/province, metropolitan region, and so on. To allow for a more meaningful comparison of the data, we aggregated related variables (e.g. White, Black, Asian, Other) into what we called an *aspect* (e.g. Race). The aspects are represented using normalised histograms. This normalisation is crucial for direct

comparison. In essence, it is a generalisation of the standard method of comparing percentages, 206
since each aggregation area has a different total population. 207

Each area of each census year is represented as a node, and edges are placed between nodes if the 208
corresponding regions share geographic borders in the same year. Further, edges are placed between 209
nodes if the corresponding regions belong to sequential years and there is geographical overlap 210
between them. This approach leads to a single network representing the whole spatio-temporal 211
space of the data. Our objective then becomes to identify partitions of this network such that the 212
nodes of each partition are more similar between themselves than to the other nodes. 213

4.2 Geographic content clustering

 214

Having tied the regions together into a network, we can now partition it to identify similar sets 215
of regions. We start by adopting a distance function between the nodes, measuring the difference 216
between the data of the regions. This value is then associated with the edges, leading to a weighted 217
dynamic network. Every node has a collection of histograms, each representing the distribution of 218
certain aspect in the population. 219

Let $G = (V, E)$ be a network, where $V = \{v_1, v_2, \dots, v_n\}$ is the set of nodes and $E = \{(v_i, v_j), i \neq 220$
 $j \text{ and } i, j \in [1, n]\}$ is the set of edges. A function H associates each node to a set of K histograms. 221
We define the distance D between two nodes v_i and v_j as: 222

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

where d is a distance metric between histograms and w is a sequence of non-negative weights associated 223
with each aspect, $\sum_{k \in [1, K]} w_k = 1$. While any histogram metric can be used, we adopted a 224
euclidean distance between the vectors, because it led to reasonable results with reduced computational 225
cost. Therefore the distance between two nodes is defined as the weighted average distance 226
between its associated histograms, where the weights can be adjusted by the user. 227

Once the distances are associated to the edges, we use watershed cuts (Cousty et al., 2009) to 228
create an initial clustering, which is then refined into a hierarchy using the Sorted Maximal Matching 229
(SMM) (Dias et al., 2017) with median linkage. The initial watershed step is performed to create 230
an initial clustering and reduce the running time of the SMM. We introduced one new parameter to 231
this method: a maximum distance threshold for the merges, to avoid the early merging of outliers. 232

233 We refer the reader to the original paper (Dias et al., 2017) for more details, including a complete
234 performance evaluation using several metrics.

235 Each resulting cluster is contiguous in the network. This means that two similar, but non-
236 contiguous, sets of areas will be classified into two different clusters, which can be counter-intuitive.
237 To overcome this issue, we *augment* the network with two new edges per node from a nearest
238 neighbours graph (Pedregosa et al., 2011) using only the distances between the histograms. These
239 edges connect nodes with similar content, if they are not already connected, providing a path for
240 the algorithm to group similar nodes.

241 4.3 Cluster characterisation and variable relevance

242 A crucial step in understanding neighbourhood change is to characterise the evolving clusters. The
243 composition of each cluster is represented here by simple statistical measures, considering each aspect
244 separately. We compute the minimum, maximum, median, 25%, and 75% quantiles for each variable
245 of each aspect for all clusters in the hierarchy. While interpreting these values is more complex than
246 interpreting just the average, they provide far more information about the underlying distribution.

247 We also use these statistical measurements to discover what characterises each cluster, that is,
248 what makes it different from the others. We define the *relevance* of a variable of an aspect based
249 on the distance between the interquartile ranges (IQR) of the clusters in the same hierarchical
250 level. If the IQRs overlap for all clusters, that variable is not relevant to the characterisation of the
251 cluster, but if the IQRs are distant, it means that this specific range of values is something that only
252 occurs in this cluster. Examining IQRs therefore provides users a straightforward visual method for
253 determining what variables most clearly define a given cluster.

254 4.4 Clusters and trajectories

255 While the partition of the data into different clusters helps the user to understand what groups exist
256 and where they are, we are also interested in the evolution of these groups. To examine this process
257 of evolution directly, we introduce the concept of *trajectories*. Trajectories are composed by regions
258 classified into the same sequence of clusters over the considered period. This enables direct access to
259 regions that evolved in the same manner. While individual census tracts remain interesting, these
260 trajectories are the main unit of exploration in this work.

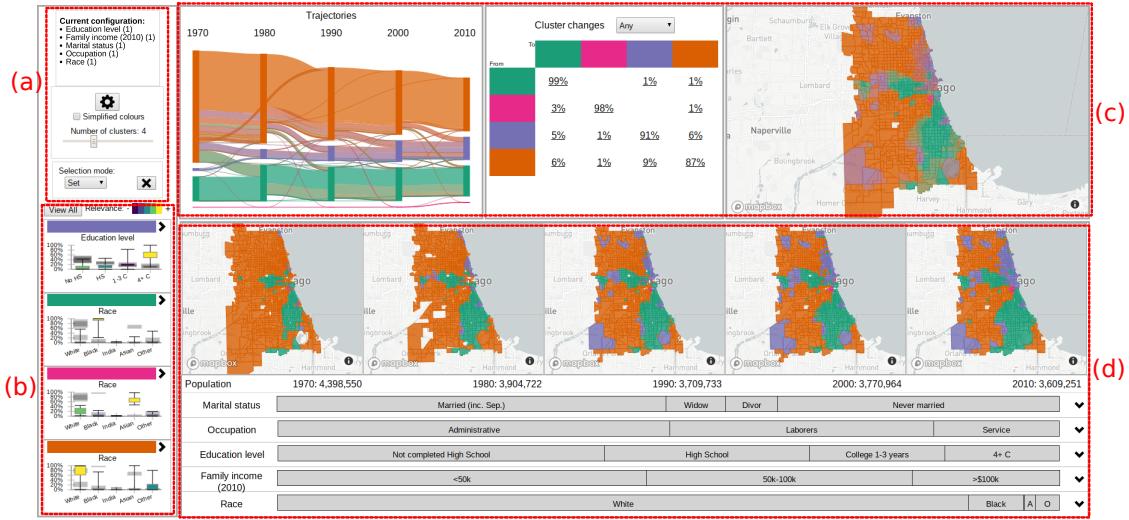


Figure 3: Initial interface of our method showing the demographic evolution of Chicago. **(a)**: Configuration panel with the current clustering parameters and controls. **(b)**: Cluster overview illustrating the most relevant aspect for each cluster. **(c)**: Trajectories overview and the general evolution 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.

4.5 User interface

To validate and explore the results of our methodology, we built a user interface, illustrated in Figure 3, considering census tract (CT) level data from the Chicago region between 1970 and 2010. This region is known for its entrenched racial divide and the emergence of a '*young urban*' population with a higher education level (Delmelle, 2016, 2017). More details are presented in Section 5.

As illustrated by Figure 3, our proposed interface heavily relies on colour to express cluster-related information. We adopted this convention because colours can be used in all our visual tools in a coherent manner. However, there is a limit on the number of distinct colours that can be used. We limited the number of clusters to eight because this was the largest number of colours that we could reliably and accessibly use, derived from the 8-class Dark2 set from ColorBrewer (Harrower and Brewer, 2003).

The configuration panel, on top left in Figure 3, displays which aspects were used 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 colour option. The gear button allows access to the other configuration options that do require further processing, such as changing location, aspects, and weights.

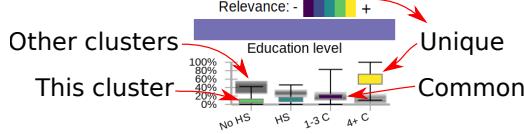


Figure 4: Enhanced boxplot of the clusters' characteristics allows a quick comparison to the other clusters.

277 The cluster overview panel, on the bottom left in Figure 3, displays a brief summary of each
 278 cluster, based on the distance between the IQRs, as detailed in Section 4.3. The *View all* button
 279 opens a new panel where all aspects are included, while the chevron at the side lets the user expand
 280 each cluster separately.

281 We adopted an *enhanced* version of the traditional boxplot, which includes the IQRs for the other
 282 clusters, in slightly larger and faded black rectangles. We also colour the current IQR according to
 283 its relevance. For instance, the boxplot that summarises the purple cluster illustrated in Figure 3,
 284 detailed on Figure 4, illustrates that this cluster is best defined by the proportion of the population
 285 with four or more years of college. The user can quickly see that this is relevant because the
 286 corresponding IQR is coloured with the highest relevance present in the legend. It is also clear
 287 that, while this cluster includes CTs that have between 10% to 90% of people in this variable,
 288 approximately, half of them have about 60% of the population with four or more years of college.
 289 Since all the other IQRs are well separated, this is a defining characteristic of this cluster. Conversely,
 290 the proportion of the population with one to three years of college is not relevant, as indicated by
 291 black fill in the rectangle representing the IQR of this cluster, in overlapping position with the
 292 rectangles of the other clusters. By clicking on the coloured bar above the boxplot, the user can
 293 select all trajectories that contain this cluster at any point in time.

294 The trajectories overview aims to convey basic information about the trajectories, where they are,
 295 and what changes are involved. This is done using three sub panels. The first, on the left, contains
 296 a Sankey diagram illustrating the evolution of the clusters over time. The widths are proportional
 297 to the population involved. In our example in Figure 3, the orange and green clusters contain most
 298 of the population and are fairly stable over time. The pink cluster is small and mostly stable. The
 299 purple cluster is increasing, mostly by incorporating areas that were previously orange. Since the
 300 purple group corresponds to the emergent 'young urban' group, this corroborates the findings of
 301 Delmelle (Delmelle, 2016, 2017), showing that our network-based method can recover results from

the traditional data processing approach. 302

In the next panel, illustrated in the top middle of Figure 3, is a transition matrix between the 303
clusters. It indicates a rounded percentage of the population whose area changed between each pair 304
of clusters. This kind of table can be found in the related literature (Delmelle, 2016), so it is familiar 305
to the advanced users. It not only informs the proportional changes, but allows the selection of the 306
corresponding trajectories for further analysis. 307

The panel in the top right of Figure 3 is a map of the region under analysis, summarising the 308
geographical evolution of the clusters over time. The colours are derived from the clusters involved 309
in each trajectory, which are consistent across the linked views. 310

The bottom part of the interface contains the details for the selected trajectories, or for the whole 311
city if nothing is selected, as illustrated in Figure 3. This panel contains two main regions: the small 312
multiple maps, depicting the clusters at each year, and the stacked bar plots that summarise the 313
overall composition of these regions. In this example, the maps show the transition from orange to 314
green and purple in several regions over time. Clicking on a region in these maps will bring up a 315
new panel with the original census data of this specific region. The actual population numbers are 316
below the maps. 317

Each aspect is represented by a stacked bar plot, where the width of each rectangle corresponds 318
to the average percentage of that variable over the considered period. In this case, about half of 319
the people in Chicago in the considered period are married, and the percentage that are Widowers 320
or Divorced is roughly similar. About half of the population work in Administrative jobs, a third 321
never completed high-school, approximately half have gross family income below 50,000USD per 322
year. The vast majority identify as white. Placing the mouse over one of the bars will open a small 323
panel with the temporal evolution of that specific variable, and clicking on the chevron on the right 324
side expands the corresponding aspect, showing details of the temporal evolution of each variable 325
and also the corresponding IQRs for the whole city. 326

5 Illustrative scenarios 327

In this section we present three illustrative scenarios, using decennial census data from the United 328
States (Manson et al., 2017) and Canada², tabulated by CTs, from 1970 to 2010. The prototype 329

²<http://datacentre.chass.utoronto.ca/census/>

330 interface allows access to 41 regions, 29 in the US and 12 in Canada. New York City was split into
331 its boroughs to avoid memory crashes on the client browser due to the high number of CTs. We used
332 five aspects for the USA: Education level, Family income, Marital status, Occupation, and Race;
333 and seven for Canada: Age, Education level, Home language, Household Income, Marital status,
334 Occupation, Place of birth, and Religion.

335 While our method does not require geographic harmonisation, it requires matching the variables
336 over time. The supplementary material contains the details of which census columns were used for
337 each aspect. Income is slightly inaccurate, even though we did correct for the official inflation. We
338 grouped the original ranges into three larger ranges, but they do not match precisely.

339 These results are meant to demonstrate the utility of the interface for understanding the evolu-
340 tionary dynamics of urban neighbourhoods. They also show the face validity of the results generated
341 by our novel network-based approach.

342 5.1 Chicago

343 Our first scenario examines Chicago, focusing on a region loosely following the City's administrative
344 borders. Its demographic composition is well explored in the literature, with reports of racial divide
345 and gentrification (Delmelle, 2016, 2017; Hwang and Sampson, 2014). Our goal is therefore to
346 confirm that our novel method can reproduce existing findings.

347 The initial state of the prototype is illustrated in Figure 3. The specific workflow used to identify
348 the existence and details of the gentrification process is illustrated in Figure 5. The first step is to
349 identify the compositions of each cluster from the boxplots, so orange is associated with majority
350 of White population, green with majority Black, and purple with higher proportion of four years of
351 college or more (high education level). The expanded version of the boxplots for the purple cluster
352 shows a higher income level and majority of occupations in administrative jobs, therefore the purple
353 cluster identifies gentrified regions.

354 The trajectories plot illustrates the process of gentrification, progressively absorbing regions
355 from the orange cluster (White). This corroborates results from the literature reporting that Black
356 neighbourhoods are less likely to gentrify (Hwang and Sampson, 2014). Moreover, this process is
357 unlikely to be reversed, as indicated by the limited number of trajectories leaving the purple stream.
358 Next, we select the region that is gentrified in 2010, by clicking on the corresponding rectangle in
359 the trajectories plot, updating the information on the maps and the details portion of the interface.

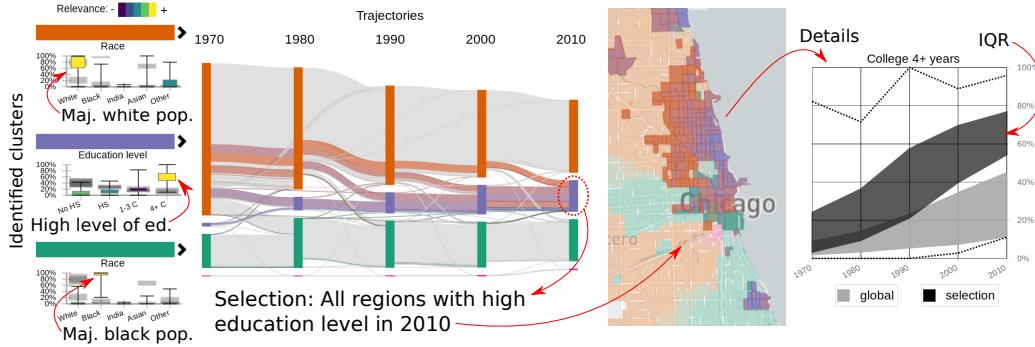


Figure 5: Workflow to discover gentrification in Chicago: the purple cluster corresponds to high education / income. Its population is increasing over time, absorbing from the majority White cluster (orange). By selecting the purple cluster in 2010, the region is highlighted in the maps. The proportion of people with 4+ years of college is increasing in the whole city (grey IQRs), but significantly more in this region (black).

The corresponding regions are highlighted in the maps, where the spatial pattern is clear, corresponding exactly to previous findings in the literature based upon harmonisation (Hwang and Sampson, 2014). Further, we can also identify the regions that gentrified earlier on the small maps that depict the involved regions over time. Since the most relevant aspect is Education, specifically "Four or more years of college", we can expand the details of this aspect, as illustrated in the right-most portion of Figure 5, which is increasing for the whole city (grey band), but faster and to a higher level in this region (black band).

5.2 Toronto

We consider a region that is approximately the administrative border of the current city of Toronto, using all seven available aspects with equal weights. The results are summarised in Figure 6, considering eight clusters.

The population with low percentage of University degrees is represented in orange, mostly anglophone population in green, Asian immigrants in yellow, high percentage of income in the highest bracket in purple, high percentage of Jewish people in light green and brown, high percentage of Eastern Non-Christian religion in pink, and high concentration of single people in dark grey. From the trajectories plot, we can see that Toronto is more dynamic than Chicago, with one cluster constantly shrinking. In the 1970s, the city was divided into four clusters: low number of university degrees, Jewish population, majority anglophones, and high income. Interestingly, the more recent

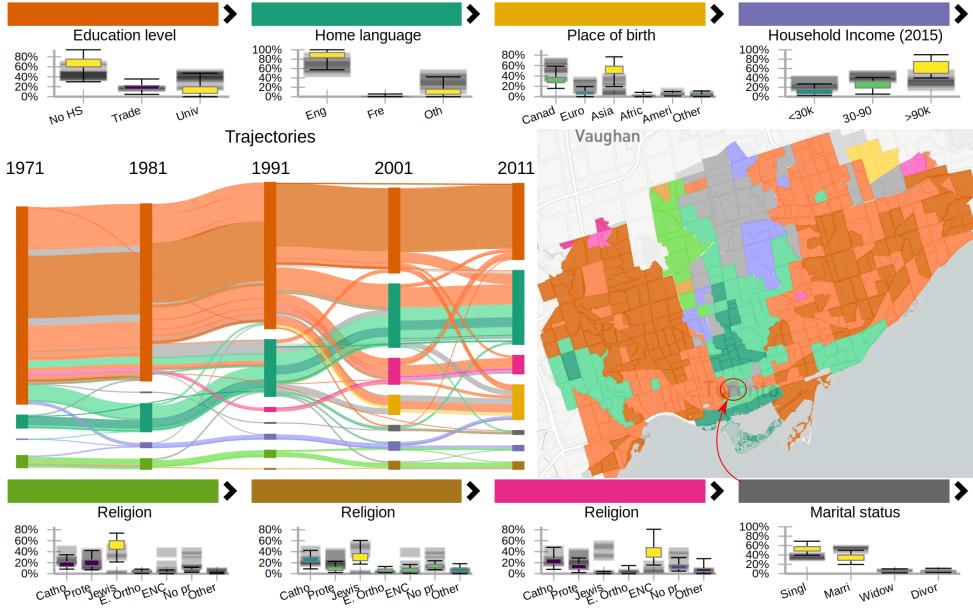


Figure 6: Clustering results for Toronto, with eight clusters, including clusters representing Jewish population, high and low income, low education, and Asian immigration.

378 clusters that absorbed regions from the orange cluster have similar education profiles and are differ-
 379 entiated by other aspects. In this sense, the city is growing diverse, changing from a common low
 380 education profile to a higher level of education with more diversity in religion (pink) and immigration
 381 (yellow).

382 Indeed, the growing Asian population is visible starting in the 1980s and building thereafter,
 383 leading to the yellow and pink clusters. While both include a high percentage of people born in
 384 Asia, the pink is more defined by religion, with low percentage of university degrees, and contains the
 385 lowest percentage of people in the highest income bracket for these clusters; the yellow is less defined
 386 by religion, and has higher education and income, geographically corresponding to the Markham
 387 region, known for its Chinese population. A similar division also happens for the two Jewish clusters,
 388 where the light green cluster has lower education and income levels than the brown cluster. The
 389 purple cluster of high income is somewhat stable. Until 2011 the cluster included the Bridle Path
 390 neighbourhood, known for its wealthy population. In 2011 it was classified into the yellow cluster
 391 of Asian immigration, since about 35% of the population for this CT were then born in Asia. The
 392 income distribution did not change, with 85% of the population with an income of 90k CAD or
 393 more.

394 The most significant indicator of Toronto's dynamism is the presence of grey regions on the map.

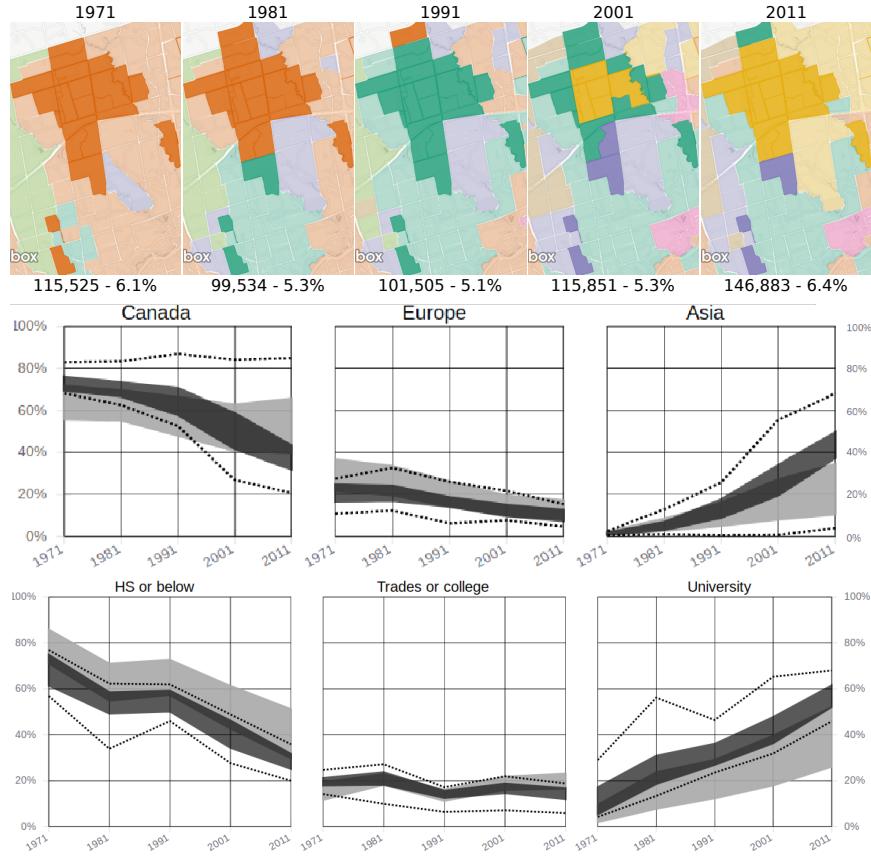


Figure 7: Details for some regions of Toronto that were classified into 3 or more clusters over time.

These represent regions associated to three or more clusters over this five census period. Using the 'Add' mode for the trajectory selection, we select their trajectories, and a subset of the details is illustrated in Figure 7. 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, 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 neighbouring parts of the volatile region, and the Asian born population increases sharply, as illustrated by the appearance of the yellow cluster. This cluster 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 neighbouring Bridle Path area.

The geographical borders of the clusters obtained using our method are similar to the regions

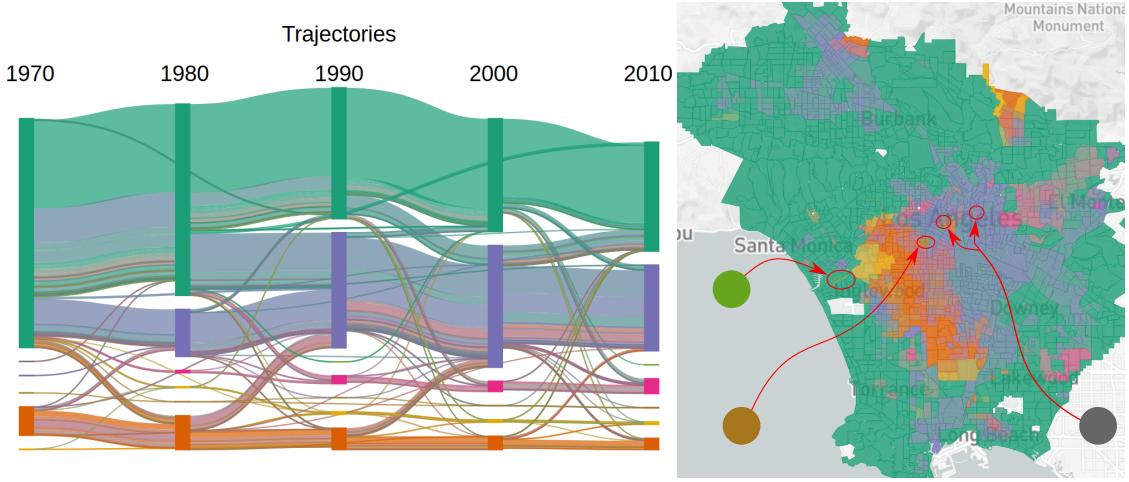


Figure 8: Result for Los Angeles with 8 clusters, including three small and ephemeral clusters. Cluster characterisation is displayed in Figure 9.

408 presented by previous studies considering Toronto (Hulchanski, 2007), but our interface provides a
 409 deeper insight into their demographic composition, since we consider more data than solely Average
 410 Income, which appears to be a good proxy variable nonetheless. This scenario showcases the ability
 411 of our method and interface to capture and understand the sources of urban volatility.

412 5.3 Los Angeles

413 We selected for this scenario a region around the metropolitan area of Los Angeles (LA), following
 414 urban density. The summary of the results using all aspects with equal weights and eight clusters
 415 is illustrated in Figure 8. The full statistical description of the clusters is illustrated in Figure 9,
 416 where the most relevant aspect of each cluster is highlighted. From the trajectories plot, we can see
 417 that there is a large but shrinking cluster, depicted in green, one increasing cluster in purple, an
 418 almost constant orange cluster, a smaller but increasing pink cluster, and three other small clusters.
 419 The corresponding map illustrates where these clusters are located, and that they are somewhat
 420 geographically stable, with some movement between the green, orange, and purple clusters.

421 From Figure 9, we can see that the green cluster is characterised by a high percentage of White
 422 population, low percentage of population in the lowest income bracket, mostly administrative occu-
 423 pations, and about 30% of the population with four or more years of college. The orange cluster
 424 is characterised by a high percentage of Black population, with few people in the highest range of
 425 income and education. The purple cluster corresponds to a high concentration of "Other" in race,

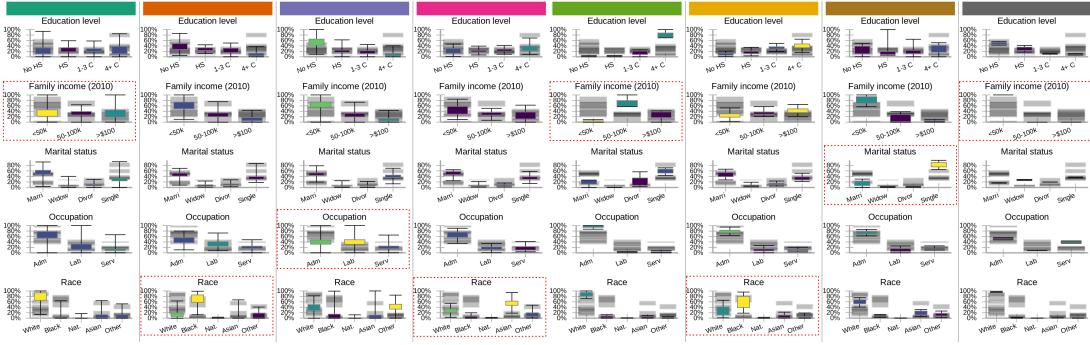


Figure 9: Full characterisation of the eight clusters found for LA. The red rectangles indicate the most relevant aspects for each cluster.

which includes Hispanic for this dataset, high concentration of Labourers, and low education and income. The pink cluster contains a high percentage of Asian population and about 30% of the population with four or more years of college. The light green cluster contains very few people in the lower income bracket, mostly White population, with the highest percentage of population with four or more years of college, working administrative jobs, and a high concentration of singles. The yellow cluster represents Black population, with higher level of education and income, mostly working administrative jobs. The brown cluster represent a majority of single population, working administrative jobs with mostly low income. The dark grey cluster is characterised by all its population in the lowest income bracket, low education level, with a majority of White population. Since the extremes in the boxplots of the grey cluster are not significantly different, we can also surmise that this cluster is either small or homogeneous.

The green, orange, and purple clusters present a significant intra-cluster variance in most variables, as indicated by extreme whiskers of the boxplots. While fifty percent of the CTs in the green cluster have between 20% and 40% of people in the lowest income bracket, that cluster also includes CTs where none and all the population belongs to that bracket. This might indicate that this cluster represents different groups of people that are not different enough to be separated at this level of the hierarchy. Conversely, the light green, brown, and dark grey clusters are different enough to be separated into their own clusters at this level, despite being small and ephemeral, including only a few CTs.

The orange area in the map in Figure 8 presents movement, indicated by the presence of green and purple tones mixed with the orange, which may warrant further exploration. By clicking on the orange bar above the boxplot, we select all trajectories that contain the orange cluster. The

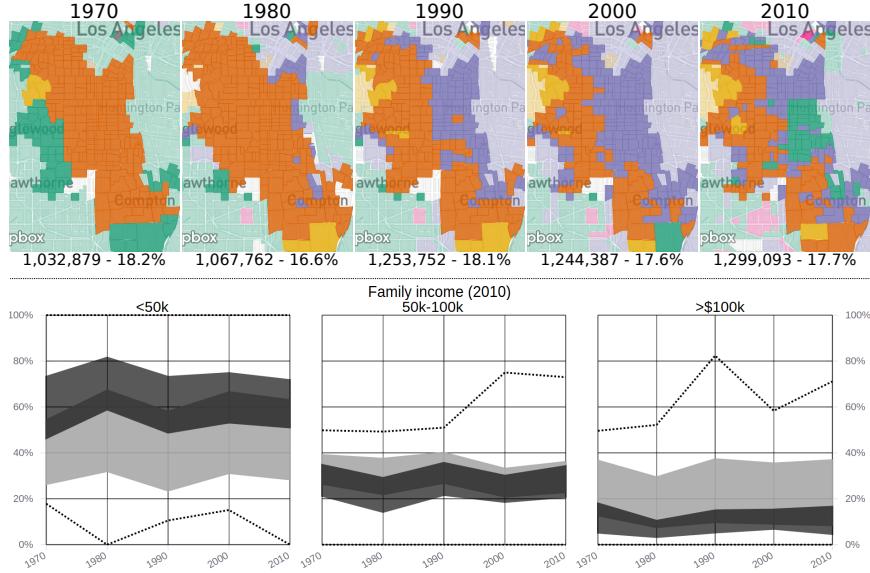


Figure 10: Top: Geographic changes in the majority Black population cluster (orange) and Labourers cluster (green). Bottom: Income evolution for this region (black) and the whole city (grey).

corresponding details are illustrated in Figure 10. This shows a location change, where the orange cluster is progressively replaced by the purple cluster on its east side, and in turn expanding to the west. Interestingly, the population increased, the racial profile changed, but the distribution of income was reasonably stable, with a higher amount of the population in the lowest income range and very few people in the highest income range. Indeed, the income difference is significant when compared to the city-wide distribution.

A portion of this region is classified into the green cluster in 2010, indicating a majority white population. To further understand that change, we clear the current selection, and select all regions that changed from orange in 1970 to green in 2010, using the transition matrix. A portion of the resulting region, near the Florence-Graham region, is depicted in Figure 11, along with the temporal evolution of Race. Despite this difference, the other aspects are similar to the ones from the region in Figure 10, with slightly lower income and education profiles. While the racial aspect changed considerably, the economic and educational aspects stayed the same.

While Toronto is more dynamic than Los Angeles, possibly due to size differences, the volatile regions shown in Figure 7 did not change as quickly or dramatically as the ones shown in Figure 11, 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

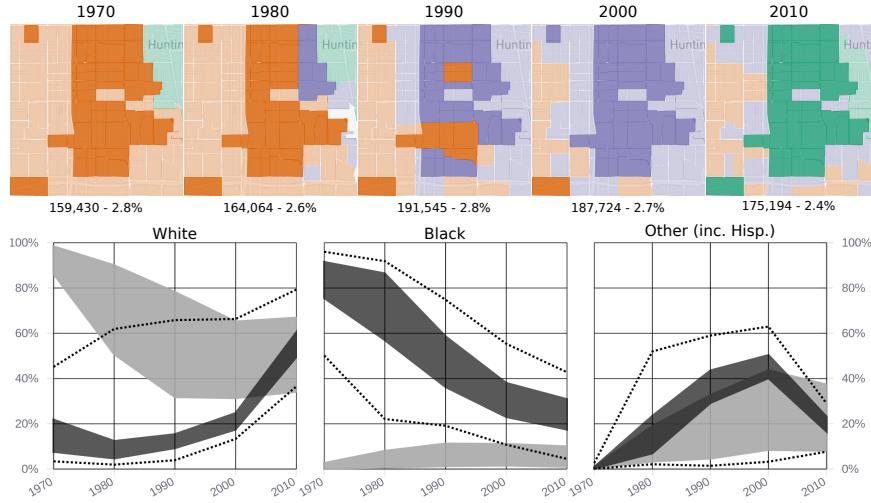


Figure 11: Details for a volatile region contained in the area of Figure 10. This region went from Black to Hispanic to White.

more general stability, but quicker changes in smaller scales. The supplementary material contains
465 brief summaries of all the regions accessible in this prototype.
466

If Toronto illustrates a relatively smooth process of whole-sale diversification, Los Angeles shows
467 the ability of our system and method to uncover more intensive smaller-scale processes of neighbour-
468 hood change. This capacity to examine and classify different types of change operating at different
469 scales demonstrates the flexibility of the method and tool for advancing the more open-ended notion
470 of "neighbourhood" envisioned by recent neighbourhood research.
471

6 Expert feedback

As our method and tool are novel to the field, and somewhat exotic, we subjected them to the
473 critical scrutiny of experts. We contacted academic and industry experts in sociology and urban
474 sciences to solicit their evaluation of our methodology. They had access to the prototype tool, a
475 descriptive documentation of the features (included in the supplementary material), and a sequence
476 of documentation videos illustrating how to perform specific tasks. The documentation explains
477 which datasets are used and how the data is represented and processed, noting explicitly that there
478 is no geographic harmonisation. We focused our inquiries on the results obtained, asking if they
479 found anything interesting in the data. The message sent and their full response is included in the
480 supplementary material. Each of the five experts is identified by a letter, from A to E.
481

482 The overall overall response of the experts was positive, mentioning that the prototype allows
483 them to analyse census data without the additional work of obtaining and cleaning the data (A, B,
484 E), and it allows the inclusion of geographic visual analysis tools in their research process (D). It
485 enables the users to tell different stories about neighbourhoods/cities and their changes (A), visualise
486 the relationship between key urban variables over time (D), offering a quick way to identify particular
487 neighbourhoods that one may be interested in studying more in depth around a particular issue or
488 efficiently understanding the context of an area (E). Indeed, the experts identified gentrification
489 processes in Manhattan (B) and Dallas (E), reinforced a hypothesis for occupational clustering (D),
490 and highlighted how the method can be used to compare neighbourhoods and cities (A). In summary,
491 their view was that the proposed methodology can be a viable alternative for the visual analytics of
492 evolving demographic data.

493 The interface was "easy to navigate" (B), but it was also considered "overwhelming" (A), "in-
494 timidating" (E), and "tricky to interpret" (C), possible side-effects of our effort to increase repre-
495 sentational accuracy, where we avoided using simplified representation or labels. Identifying clusters
496 by their most relevant variables was welcome, but the overlap of information from different clusters
497 in the boxplot was "a bit confusing" (C) when colour was not present. Further, most clusters can
498 be sufficiently characterised using only the most relevant aspect, but this is not generally true.

499 While the map of trajectories was mentioned as a "good summary map", how it related to the
500 clustering method was unclear (C). The methods include different options on how the colours are
501 used, but both are works in progress since reliably representing several distinct entities using colours
502 is humanly unfeasible. Indeed, the number of distinguishable colours was a significant constraint,
503 we found indications that more clusters should be used in some cases, even if eight clusters is more
504 than what is traditionally considered in these analyses. Conversely, increasing the number of clusters
505 would also complicate the interpretation of the results.

506 The experts also mentioned the poor responsiveness of the method when changes in the clustering
507 parameters required server-side processing (B,D). Indeed, the current implementation can take a few
508 minutes to cluster regions with high number of CTs, like Los Angeles or Brooklyn. Server-processing
509 reduced the amount of data transferred to client, but it might increase the response time under load.
510 We implemented a cache policy that greatly improved the performance, but fully pre-processing the
511 results is not practical due to size of the parameter space.

512 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. 513
They also mentioned the need to export subsets of data, plots, and maps to be used in reports and 514
publications (C, D, E). More importantly, while these experts were aware that our method does 515
not perform geographic harmonisation, none of them mention it. We did not specifically ask if this 516
difference led to unexpected results, but rather if they found interesting insights. Most experts found 517
phenomena corroborated by the specialised literature, indicating that our methodology produces 518
equivalent results, with a fraction of the effort. We interpret the fact that most of them were 519
interested in the next steps as confirmation of the accuracy of the method. 520

7 Discussion and limitations

Our objective was to leverage a network based data representation and visualisation methods for 522
the exploration of geographically inconsistent region-based data. While we successfully replicated 523
and corroborated results from the literature, this method still has significant limitations. 524

Removing the need for geographical harmonisation greatly reduces the amount of work necessary 525
to explore demographic data, but the method still requires consistent variables across the years. 526
Matching the variables can be trivial for some aspects (Age), but challenging for others (Income). 527
The divulged income ranges vary over time and the actual values change due to inflation. Some 528
variables were not considered in earlier censuses, such as Race in Canada, or Hispanic population 529
in the USA, hampering its use when they are available. Since this is only a prototype, we matched 530
few aspects, but a proper demographic analysis would benefit from all available information. 531

The limitation on the number of displayed clusters because of the limited number of distinguishable 532
colours was significant. While increasing the number of clusters would further complicate an 533
already complex analysis, it might be warranted for some regions. Colour is a fundamental and 534
intuitive tool for information representation that can be coherently used across different plots, so 535
we opted to use it, even if in a limited way. With eight colours, there was overlap between some 536
clusters, the relevance gradient, and the colour combination. 537

The cognitive load on the user is significant, as we compromised simplicity for accuracy. While 538
other works labelled the clusters, as 'young urban', 'struggling', and so on (Delmelle, 2016, 2017), 539
we show the statistical characteristics of the clusters, which are harder to interpret, as the data may 540
have subtle nuances that labels would otherwise hide. This also led to a crowded interface, mitigated 541

542 somewhat the use of pop-up panels and collapsible sections. For some cities, especially if they are
543 small and stable, the panels can appear redundant, but each provide a different way to interact with
544 the information that can ease the exploration process for larger and dynamic cities.

545 8 Conclusion

546 The objective of this work was to demonstrate that longitudinal studies of evolving regions can
547 be performed without geographical harmonisation. We proposed an alternative methodology that
548 robustly considers the data in its original geography, without the creation of arbitrary artificial data
549 points.

550 This methodology was then used to create a publicly accessible system, with an interactive and
551 intuitive interface, allowing a transparent evaluation and replication of our results. We used this
552 interface to corroborate results from the literature and we hope that it will be used to corroborate
553 future results as well.

554 The feedback from experts was positive and most of them were able to extract insight from the
555 prototype while indicating interest in using it for their research efforts. Indeed, the experts also
556 demonstrated further interest in similar tools, indicating that visual analytics methods that leverage
557 user interaction can be valuable in this field. Since our interface can be used by non-experts as well,
558 we also contributed to scientific dissemination and stakeholder transparency in urban sciences.

559 More importantly, we introduced a new idea that, apparently, was never considered in the liter-
560 ature. While significant resources were used trying to improve geographical harmonisation, nobody
561 questioned if it was really necessary. We proved that it is not necessary, at least for the vast majority
562 of demographic studies, especially for neighbourhood effects and neighbourhood dynamics.

563 References

- 564 Abbott, A. (1997, 06). Of Time and Space: The Contemporary Relevance of the Chicago School.
565 *Social Forces* 75(4), 1149–1182.
- 566 Allen, J. and Z. Taylor (2018). A new tool for neighbourhood change research: The canadian longi-
567 tudinal census tract database, 1971–2016. *The Canadian Geographer / Le Géographe canadien*.
- 568 Andrienko, G., N. Andrienko, H. Schumann, and C. Tominski (2014). Visualization of trajectory

- attributes in space–time cube and trajectory wall. In *Cartography from Pole to Pole*, pp. 157–163. Springer. 569
- Andrienko, N., G. Andrienko, and P. Gatalsky (2003). Exploratory spatio-temporal visualization: 570
an analytical review. *Journal of Visual Languages & Computing* 14(6), 503–541. 571
- Beck, F., M. Burch, S. Diehl, and D. Weiskopf (2014). The State of the Art in Visualizing Dynamic 572
Graphs. In R. Borgo, R. Maciejewski, and I. Viola (Eds.), *EuroVis - STARs*. The Eurographics 573
Association. 574
- Chavent, M., V. Kuentz-Simonet, A. Labenne, and J. Saracco (2017). Clustgeo: an r package for 575
hierarchical clustering with spatial constraints. *Computational Statistics*, 1–24. 576
- Chen, W., F. Guo, and F.-Y. Wang (2015). A survey of traffic data visualization. *IEEE Transactions 577
on Intelligent Transportation Systems* 16(6), 2970–2984. 578
- Chen, W., Z. Huang, F. Wu, M. Zhu, H. Guan, and R. Maciejewski (2017). Vaud: A visual 579
analysis approach for exploring spatio-temporal urban data. *IEEE Transactions on Visualization 580
& Computer Graphics*. 581
- Cousty, J., G. Bertrand, L. Najman, and M. Couprise (2009, Aug). Watershed cuts: Minimum 582
spanning forests and the drop of water principle. *IEEE Transactions on Pattern Analysis and 583
Machine Intelligence* 31(8), 1362–1374. 584
- Dal Col, A., P. Valdivia, F. Petronetto, F. Dias, C. T. Silva, and L. G. Nonato (2018). Wavelet- 585
based visual analysis of dynamic networks. *IEEE Transactions on Visualization and Computer 586
Graphics PP(99)*, 1–1. 587
- Delmelle, E. C. (2016). Mapping the dna of urban neighborhoods: Clustering longitudinal sequences 588
of neighborhood socioeconomic change. *Annals of the American Association of Geographers* 589
106(1), 36–56. 590
- Delmelle, E. C. (2017). Differentiating pathways of neighborhood change in 50 u.s. metropolitan 591
areas. *Environment and Planning A: Economy and Space* 49(10), 2402–2424. 592
- Dias, F. and L. G. Nonato (2015). Some operators from mathematical morphology for the visual 593
analysis of georeferenced data. In *Workshop on Visual Analytics, Information Visualization and 594
Scientific Visualization - SIBGRAPI*. 595

- 597 Dias, M. D., M. R. Mansour, F. Dias, F. Petronetto, C. T. Silva, and L. G. Nonato (2017, Oct). A
598 hierarchical network simplification via non-negative matrix factorization. In *2017 30th SIBGRAPI*
599 *Conference on Graphics, Patterns and Images (SIBGRAPI)*, pp. 119–126.
- 600 Diez-Roux, A. V., F. J. Nieto, C. Muntaner, H. A. Tyroler, G. W. Comstock, E. Shahar, L. S.
601 Cooper, R. L. Watson, and M. Szkoł (1997). Neighborhood environments and coronary heart
602 disease: a multilevel analysis. *American journal of epidemiology* 146(1), 48–63.
- 603 Dmowska, A. and T. F. Stepinski (2018). Spatial approach to analyzing dynamics of racial diversity
604 in large u.s. cities: 1990–2000–2010. *Computers, Environment and Urban Systems* 68, 89 – 96.
- 605 Dmowska, A., T. F. Stepinski, and P. Netzel (2017, 03). Comprehensive framework for visualizing
606 and analyzing spatio-temporal dynamics of racial diversity in the entire united states. *PLOS*
607 *ONE* 12(3), 1–20.
- 608 Doraiswamy, H., N. Ferreira, T. Damoulas, J. Freire, and C. T. Silva (2014, Dec). Using topological
609 analysis to support event-guided exploration in urban data. *IEEE Transactions on Visualization*
610 *and Computer Graphics* 20(12), 2634–2643.
- 611 Eicher, C. L. and C. A. Brewer (2001). Dasymetric mapping and areal interpolation: Implementation
612 and evaluation. *Cartography and Geographic Information Science* 28(2), 125–138.
- 613 Fahad, A., N. Alshatri, Z. Tari, A. Alamri, I. Khalil, A. Y. Zomaya, S. Foufou, and A. Bouras (2014,
614 sep). A survey of clustering algorithms for big data: Taxonomy and empirical analysis. *IEEE*
615 *Transactions on Emerging Topics in Computing* 2(3), 267–279.
- 616 Ferreira, N. (2015). *Visual analytics techniques for exploration of spatiotemporal data*. Ph. D. thesis,
617 Polytechnic Institute of New York University.
- 618 Ferreira, N., M. Lage, H. Doraiswamy, H. Vo, L. Wilson, H. Werner, M. Park, and C. Silva (2015).
619 Urbane: A 3d framework to support data driven decision making in urban development. In *Visual*
620 *Analytics Science and Technology (VAST), 2015 IEEE Conference on*, pp. 97–104. IEEE.
- 621 Firebaugh, G. and C. R. Farrell (2016, Feb). Still large, but narrowing: The sizable decline in racial
622 neighborhood inequality in metropolitan america, 1980–2010. *Demography* 53(1), 139–164.
- 623 Galster, G. C. (2019). *Making our neighborhoods, making our selves*. University of Chicago Press.

Gotway, C. A. and L. J. Young (2002). Combining incompatible spatial data. <i>Journal of the American Statistical Association</i> 97(458), 632–648.	624 625
Hallisey, E., E. Tai, A. Berens, G. Wilt, L. Peipins, B. Lewis, S. Graham, B. Flanagan, and N. B. Lunsford (2017, Aug). Transforming geographic scale: a comparison of combined population and areal weighting to other interpolation methods. <i>International Journal of Health Geographics</i> 16(1), 29.	626 627 628 629
Han, J., M. Kamber, and A. K. Tung (2001). Spatial clustering methods in data mining. <i>Geographic data mining and knowledge discovery</i> , 188–217.	630 631
Harrower, M. and C. A. Brewer (2003). Colorbrewer.org: An online tool for selecting colour schemes for maps. <i>The Cartographic Journal</i> 40(1), 27–37.	632 633
Huang, X., Y. Zhao, C. Ma, J. Yang, X. Ye, and C. Zhang (2016, Jan). Trajgraph: A graph-based visual analytics approach to studying urban network centralities using taxi trajectory data. <i>IEEE Transactions on Visualization and Computer Graphics</i> 22(1), 160–169.	634 635 636
Hulchanski, D. J. (2007). The three cities within toronto: Income polarization among toronto's neighbourhoods, 1970-2005.	637 638
Hwang, J. and R. J. Sampson (2014). Divergent pathways of gentrification: Racial inequality and the social order of renewal in chicago neighborhoods. <i>American Sociological Review</i> 79(4), 726–751.	639 640
Jain, A. K. (2010). Data clustering: 50 years beyond k-means. <i>Pattern recognition letters</i> 31(8), 651–666.	641 642
Kwan, M.-P. (2018). The limits of the neighborhood effect: Contextual uncertainties in geographic, environmental health, and social science research. <i>Annals of the American Association of Geographers</i> 108(6), 1482–1490.	643 644 645
Lee, A. C.-D. and C. Rinner (2015). Visualizing urban social change with self-organizing maps: Toronto neighbourhoods, 1996–2006. <i>Habitat International</i> 45, 92–98.	646 647
Li, Y. and Y. Xie (2018). A new urban typology model adapting data mining analytics to examine dominant trajectories of neighborhood change: A case of metro detroit. <i>Annals of the American Association of Geographers</i> 108(5), 1313–1337.	648 649 650

- 651 Ling, C. and E. C. Delmelle (2016). Classifying multidimensional trajectories of neighbourhood
652 change: a self-organizing map and k-means approach. *Annals of GIS* 22(3), 173–186.
- 653 Liu, X., Y. Song, K. Wu, J. Wang, D. Li, and Y. Long (2015). Understanding urban china with
654 open data. *Cities* 47, 53 – 61. Current Research on Cities (CRoC).
- 655 Logan, J. R., B. J. Stults, and Z. Xu (2016). Validating population estimates for harmonized census
656 tract data, 2000–2010. *Annals of the American Association of Geographers* 106(5), 1013–1029.
- 657 Logan, J. R., Z. Xu, and B. J. Stults (2014). Interpolating us decennial census tract data from
658 as early as 1970 to 2010: A longitudinal tract database. *The Professional Geographer* 66(3),
659 412–420.
- 660 Looker, B. (2015). *A nation of neighborhoods: imagining cities, communities, and democracy in*
661 *postwar America*. University of Chicago Press.
- 662 Manson, S., J. Schroeder, D. V. Riper, and S. Ruggles (2017). Ipums national historical geographic
663 information system: Version 12.0 [database].
- 664 Monmonier, M. (1990). Strategies for the visualization of geographic time-series data. *Cartographica:*
665 *The International Journal for Geographic Information and Geovisualization* 27(1), 30–45.
- 666 Pedregosa, F., G. Varoquaux, A. Gramfort, V. Michel, B. Thirion, O. Grisel, M. Blondel, P. Pretten-
667 hofer, R. Weiss, V. Dubourg, J. Vanderplas, A. Passos, D. Cournapeau, M. Brucher, M. Perrot,
668 and E. Duchesnay (2011). Scikit-learn: Machine learning in Python. *Journal of Machine Learning*
669 *Research* 12, 2825–2830.
- 670 Poorthuis, A. (2018). How to draw a neighborhood? the potential of big data, regionalization,
671 and community detection for understanding the heterogeneous nature of urban neighborhoods.
672 *Geographical Analysis* 50(2), 182–203.
- 673 Reades, J., J. D. Souza, and P. Hubbard (2019). Understanding urban gentrification through machine
674 learning. *Urban Studies* 56(5), 922–942.
- 675 Rey, S., E. Knaap, S. Han, L. Wolf, and W. Kang (2018). Spatio-temporal analysis of socioeco-
676 nomic neighborhoods: The open source longitudinal neighborhood analysis package (oslnap). In
677 *Proceedings of the 17th Python in Science Conference (SciPy 2018)*, pp. 121–128.

- Sampson, R. J. (2012). *Great American city: Chicago and the enduring neighborhood effect*. University of Chicago Press. 678
679
- Sandryhaila, A. and J. M. Moura (2013). Discrete signal processing on graphs. *IEEE transactions on signal processing* 61(7), 1644–1656. 680
681
- Setiadi, T., A. Pranolo, M. Aziz, S. Mardiyanto, B. Hendrajaya, and Munir (2017, Oct). A model of geographic information system using graph clustering methods. In *2017 3rd International Conference on Science in Information Technology (ICSITech)*, pp. 727–731. 682
683
684
- Shelton, T. and A. Poorthuis (2019). The nature of neighborhoods: Using big data to rethink the geographies of atlanta’s neighborhood planning unit system. *Annals of the American Association of Geographers* 0(0), 1–21. 685
686
687
- Shuman, D. I., S. K. Narang, P. Frossard, A. Ortega, and P. Vandergheynst (2013). The emerging field of signal processing on graphs: Extending high-dimensional data analysis to networks and other irregular domains. *IEEE Signal Processing Magazine* 30(3), 83–98. 688
689
690
- Soille, P. and L. Najman (2012). On morphological hierarchical representations for image processing and spatial data clustering. In *Applications of Discrete Geometry and Mathematical Morphology*, pp. 43–67. Springer. 691
692
693
- Soor, S., A. Challa, S. Danda, B. S. Daya Sagar, and L. Najman (2018, July). Extending k-means to preserve spatial connectivity. In *IGARSS 2018 - 2018 IEEE International Geoscience and Remote Sensing Symposium*, pp. 6959–6962. 694
695
696
- Stepinski, T. F. and A. Dmowska (2019). Imperfect melting pot—analysis of changes in diversity and segregation of us urban census tracts in the period of 1990–2010. *Computers, Environment and Urban Systems* 76, 101–109. 697
698
699
- Stone, C. N., R. P. Stoker, J. Betancur, S. E. Clarke, M. Dantico, M. Horak, K. Mossberger, J. Musso, J. M. Sellers, E. Shiau, et al. (2015). *Urban neighborhoods in a new era: Revitalization politics in the postindustrial city*. University of Chicago Press. 700
701
702
- Tominski, C. and H.-J. Schulz (2012). The Great Wall of Space-Time. In M. Goesele, T. Grossch, H. Theisel, K. Toennies, and B. Preim (Eds.), *Vision, Modeling and Visualization*. The Eurographics Association. 703
704
705

- 706 Tufte, E. R., S. R. McKay, W. Christian, and J. R. Matey (1998). Visual explanations: images and
707 quantities, evidence and narrative.
- 708 Valdivia, P., F. Dias, F. Petronetto, C. T. Silva, and L. G. Nonato (2015, Oct). Wavelet-based
709 visualization of time-varying data on graphs. In *2015 IEEE Conference on Visual Analytics*
710 *Science and Technology (VAST)*, pp. 1–8.
- 711 Vehlow, C., F. Beck, and D. Weiskopf (2015). The State of the Art in Visualizing Group Structures in
712 Graphs. In R. Borgo, F. Ganovelli, and I. Viola (Eds.), *Eurographics Conference on Visualization*
713 (*EuroVis*) - STARS. The Eurographics Association.
- 714 Von Landesberger, T., F. Brodkorb, P. Roskosch, N. Andrienko, G. Andrienko, and A. Kerren
715 (2016). Mobilitygraphs: Visual analysis of mass mobility dynamics via spatio-temporal graphs
716 and clustering. *IEEE transactions on visualization and computer graphics* 22(1), 11–20.
- 717 Ward, M. O., G. Grinstein, and D. Keim (2015). *Interactive data visualization: foundations, tech-*
718 *niques, and applications*. AK Peters/CRC Press.
- 719 Zheng, Y., W. Wu, Y. Chen, H. Qu, and L. M. Ni (2016, Sept). Visual analytics in urban computing:
720 An overview. *IEEE Transactions on Big Data* 2(3), 276–296.