

Article

Historical collaborative geocoding

Rémi Cura^{1,4,†} , Bertrand Dumenieu^{3,4,†} , Nathalie Abadie², Benoit Costes², Julien Perret^{2,3,4} , Maurizio Gribaudo^{3,4}

¹ IGN; remi.cura@gmail.com

² IGN; first.last@ign.fr

³ EHESS; last@ehess.fr

⁴ GeoHistoricalData; first@geohistoricaldata.org

* Correspondence: julien.perret@ign.fr

† These authors contributed equally to this work.

Abstract: The latest developments in digital humanities have increasingly enabled the construction of large data sets which can easily be accessed and used. These data sets often contain indirect localisation information, such as historical addresses. Historical geocoding is the process of transforming the indirect localisation information to direct localisation that can be placed on a map, which enables spatial analysis and cross-referencing. Many efficient geocoders exist for current addresses, but they do not deal with temporal information and are usually based on a strict hierarchy (country, city, street, house number, etc.) that is hard, if not impossible, to use with historical data. Indeed, historical data are full of uncertainties (temporal, textual, positional accuracy, confidence in historical sources) that can not be ignored or entirely resolved. We propose an open source, open data, extensible solution for geocoding that is based on gazetteers composed of geohistorical objects extracted from historical topographical maps. Once the gazetteers are available, geocoding an historical address is a matter of finding the geohistorical object in the gazetteers that is the best match to the historical address searched by the user. The matching criteria are customisable and include several dimensions (fuzzy string, fuzzy temporal, level of detail, positional accuracy). As the goal is to facilitate historical work, we also propose web-based user interfaces that help geocode (one address or batch mode) and display over current or historical topographical maps, so that geocoding results can be checked and collaboratively edited. The system has been tested on the city of Paris, France, for the 19th and the 20th centuries. It shows high response rates and is fast enough to be used interactively.

Keywords: Historical dataset; geocoding; localisation; geohistorical objects; database; GIS; collaborative; citizen science; crowd-sourced; digital humanities

1. Introduction

1.1. Context

In historical sciences, cartography and spatial analysis are extensively used to reveal the spatial organisations that hide within data with textual indirect spatial references like placenames or postal addresses. Mapping such data requires each item to be geocoded, *i.e.* assigned with coordinates on the Earth surface by matching the indirect spatial reference with entities identified in a reference geographical datasource (*e.g.* a topographic map georeferenced in a well-known coordinate reference system) [1]. Problems emerge when such spatial references are obsolete due to the temporal gap between the data to be geocoded and the reference datasource: locating the London Crystal Palace (destroyed by fire in 1936) on a today map would be rather tricky. Worse still, it might create ambiguities and possibly lead to erroneous geocoding, as the Crystal Palace refers nowadays to a South London residential area. Whereas manual geocoding can deal with such cases, the constantly growing volume of historical data, which results from the multiplication of initiatives in the digital humanities, calls for

35 automatic approaches. Despite the existence of highly efficient (atemporal) geocoding tools and API, a
36 truly historical geocoder is still at stake [2,3].

37 1.2. Related work

38 Geocoding is an inevitable step in any spatially-based study with considerable bodies of data,
39 which makes it a critical process in various contexts: public health, catastrophe risk management,
40 marketing, social sciences, etc. Many geocoding web services have been developed to fulfil this need,
41 originating from private initiatives (Google Geocoding API, Mapzen¹), public agencies (the French
42 National Address Gazetteer²) or from the open-source community (OSM Nominatim³, Gisgraphy⁴).
43 These services can be characterized in terms of their three main components [1,4]: input/output data,
44 reference dataset and processing algorithm. The *input* is the textual description the user wants to
45 refine into coordinates (e.g. "13 rue du temple, Paris, France"). The *reference dataset* contains geographic
46 features associated with textual descriptions of addresses. The *processing algorithm* consists in finding
47 the best match between the latter and the former input description. Finally, the output usually contains
48 a geographical feature and its similarity score (perfect match or approximate for instance).

49 These tools answer millions of geocoding queries each day with great success. Indeed, the
50 quality of geocoding services can be estimated via two very important criteria [5]. The first is the
51 database quality: how complete and up to date is the reference database? The second is the results
52 characterization: how spatially accurate is each result and what is its associated confidence?

53 Despite their quality, such geocoding approaches can not be used for (geo)historical data for
54 three main reasons. The first is that existing geocoding services do not take into account the temporal
55 aspect of the query or the dataset they rely on. Indeed, they usually rely on current data, such as
56 *OpenStreetMap*⁵ data, continuously updated. As such, they implicitly work on a valid time that is the
57 present (or possibly the interval between the beginning of the database construction and the present
58 time). The second reason is that they rely on an exhaustive, strongly hierarchical database whose
59 accuracy can be checked against ground truth (i.e. there is always a way to check the actual localization
60 of an address). Unfortunately, historical data are not easily verifiable: one has to check them with
61 difference available (geo)historical sources (possibly incomplete and conflicting) and, often, making
62 assumptions or hypotheses. Such hypotheses are in their turn continuously challenged and updated
63 by new discoveries, and there is no way to give a truly definitive answer. Indeed, primary sources
64 may also be wrong or misleading. The third reason is that historical sources available to construct a
65 gazetteer are sparse (both spatially and temporally), heterogeneous, and complex. We believe that
66 all these specificities call for a dedicated approach. Similar observations have already been made
67 in the context of archival data by the UK National Archives for instance [6]. Large historical event
68 gazeteers already exist [7] and provide an important basis for the construction of the reference dataset.
69 Nevertheless, we found few related articles [2,8], and for all of them geocoding was not a focus.

70 We could not find an historical geocoding approach that considers the characterization of
71 geocoding results. Yet, this aspect is essential for historical geocoding because of the very unprecise
72 and sparse nature of geohistorical data. Indeed, geocoding results have to be validated and/or edited
73 manually. Considering the large amount of addresses (> 100000 *for Paris*) and the potential complexity
74 of the task, this is clearly a lot of work. Fortunately, several projects such as *OpenStreetMap* have lead the
75 way for what is usually called *Volunteered Geographical Information* (VGI) [9] of *crowdsourcing geospatial*
76 *data* [10]. This approach consists in using a collaborative approach to solve the problem collectively,
77 usually by implying citizens in the process. As suggested in a recent typology of participation in

1 mapzen.com

2 adresse.data.gouv.fr/api/

3 nominatim.openstreetmap.org

4 gisgraphy.com

5 openstreetmap.org

78 citizen science and VGI [11], different levels of participation can be defined. These levels go from
 79 “crowdsourcing”, where the cognitive demand is minimal, to “extreme citizen science” or “collaborative
 80 science”, where citizens are involved in all stages of the research (problem definition, data collection
 81 and analysis). In the rest of this article, we propose a collaborative historical geocoding approach
 82 that opens a way for a simpler participation of citizens in geohistorical research (thanks to dedicated
 83 interactive tools), but also for a more collaborative geohistorical science (thanks to a reproducible
 84 research approach [12–15], open source tools and open data).

85 1.3. Approach and contributions

86 In this article, we focus on the historical geocoding problem. Following the classical approach, we
 87 will present in particular the construction of a geohistorical database and the development of matching
 88 (data linkage) methods that fully use the temporal aspects of the geohistorical data and the input query.
 89 The main contributions of this article consist in (1) a formalisation of the historical geocoding problem,
 90 (2) a minimal model of geohistorical objects that can be easily re-used and extended, (3) an open source
 91 geocoding tool that is powerful, easy to use and can be extended with any geohistorical data, (4) a
 92 graphical tool to control and edit the geocoding results, which can then optionally be used to enrich
 93 the geohistorical database, (5) qualifications of geocoding results in term of textual, spatial, temporal
 94 aspects.

95 2. Approach overview

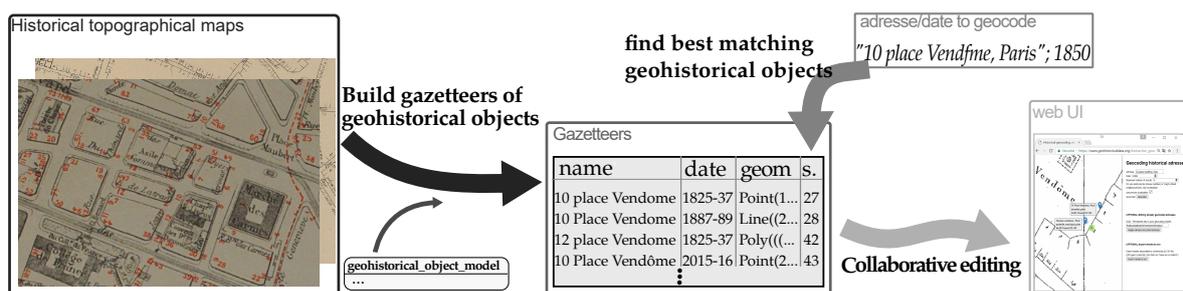


Figure 1. Gazetteers of geohistorical objects are created based on information extracted from georeferenced historical maps. Geocoding an historical address is then finding the best matching object based on a customised function (semantic, temporal aspect, spatial precision, etc.). Results can be displayed through a dedicated web interface for collaborative editing.

96 Based on historical sources and historical topographical maps, we extract geographic features
 97 that are gathered into gazetteers. These (geo) historical feature are modeled in a generic way into a
 98 Relational Database Management System (RDBMS). Geocoding an historical address is then finding
 99 the geohistorical object in the gazetteer that best match this historical address, thus may be done by
 100 means of distances that can be customised by the user. Lastly the results can be displayed via a web
 101 mapping interface over current or historical topographical maps, and further checked and edited
 102 collaboratively. Figure 1 illustrates this approach.

103 2.1. Building a geohistorical gazetteer: Extracting geohistorical objects from historical topographical maps

104 The starting point to build gazetteers is information extracted from historical topographical maps.
 105 The first part of extraction is to digitize the maps and georeferencing the map in a defined coordinate
 106 reference system. This maps are historical sources, and as such an historical analysis is performed to
 107 estimated the probable valid time (temporalization), positional accuracy, completeness, confidence,
 108 relation to other historical maps, etc. The whole process is carefully designed and explained in detail

109 in [16]. Then geohistorical objects are extracted from the referenced historical maps, manually (in a
110 collaborative way), or with the help of computer vision techniques.

111 2.1.1. General consideration about building a spatio-temporal database

112 Extracting information from topographic historical maps amounts to building a spatio-temporal
113 database. There are several approaches to do so, and we stress that we do not attempt to create a
114 continuous spatio-temporal database. Instead, we store representations of the same space at multiple
115 moments in history, the well-known snapshot model [17]. The main advantage is that for a given
116 moment in time we can have several conflicting snapshots coexisting. This is essential, as solving the
117 conflict may not be possible, and reporting these several conflicting geocoding results to historians
118 may help appreciate these results. The drawbacks of this model, *i.e.* information redundancy and its
119 inability to store the changes themselves, can be overcome during the geocoding process.

120 2.1.2. Historical topographic maps as geohistorical sources

121 In our approach, we focus on historical topographic maps as the main sources for two main
122 reasons:

- 123 • the way they portray spatial information is close to today topographic mapping, making the
124 integration of the information they convey in a Geographical Information System (GIS) easy;
- 125 • the main goal of topographic maps is to provide a reliable depiction of the shape and location of
126 geographical features.

127 Although this choice seriously reduces the number of possible sources and therefore lessens the
128 quantity of accessible spatial information, it aims at efficiency. Indeed, topographic maps are a
129 good compromise between their reliability, the quantity of spatial information they contain and the
130 complexity of extracting information.

131 We create a snapshot for each historical topographic map by relying on three steps:

- 132 1. georeferencing the map in a well-defined coordinate reference system supported by common
133 GIS tools,
- 134 2. assigning the map a *valid time*,
- 135 3. extracting geographical features from the map.

136 2.1.3. Georeferencing topographic historical maps

137 We have to establish a correspondence between each pixel of the historical maps and geographical
138 coordinates. To do so, we first choose a common spatial reference system (SRS). Then we identify
139 common geographic features between historical maps and current maps: so-called ground control
140 points (GCP). Last, we compute a warping transform that will respect at best the matching points.
141 Finding GCPs between current maps and historical maps can be increasingly difficult as we go back in
142 time, because there are less and less perennial GCPs. Consider for instance the city of Paris, where the
143 French Revolution and its consequences combined with the 19th century transformations (including
144 the so-called Haussmannian transformations) resulted in massive changes in the shape of the city. To
145 this end, we can start by georeferencing *e.g.* 20th maps to current maps, then georeference *e.g.* 19th
146 maps to 20th maps, and so on for even older maps.

147 *Common spatial reference system (SRS)*

148 The choice of a SRS is not easy, as each SRS induces projection errors that depend on the
149 covered area. Whatever the choice of SRS, it is essential that the implied accuracy is well-known and
150 documented in order to qualify the absolute accuracy of each geo-referenced map. We restrain ourselves
151 to SRS using meters as base units (opposed to SRS using degrees), as they are much closer in nature to
152 those used in historical topographical maps.

153 *Selection of ground control points*

154 The identification of pairs of GCPs is a critical step because the number, distribution and quality
155 (*i.e.* positional accuracy, reliability, confidence) of the points strongly influence the quality of the
156 georeferencing. While the quality of the selected points depends on each map, a simple rule of thumbs
157 is to select, as much as possible, homogeneously distributed points to go further [18]. Three parameters
158 have to be considered: the geometric type of the features carrying the ground control points, their
159 nature and the method used to identify them. Usually, features chosen as ground control points are
160 represented by 2D points; lines or surfaces may also be used, and possibly even curves [19].

161 On historical maps, the positionnal accuracy of mapping themes can vary greatly either due to
162 the purpose of the map, or the mapmaking process itself. Optionally, one can rely on geodetic features
163 drawn on the map such as meridian or parallels provided one can fully characterised the geodetic
164 characteristics of these lines.

165 The actual identification of GCPs can be achieved by automatic or manual processes. Automatic
166 approaches are notably used for historical aerial photographs, where feature detection and matching
167 algorithms are well fitted [20]. Common GIS tools offer georeferencing softwares allowing to manually
168 select pairs of *ground control points* identified in both the input and the reference maps. Such tools are
169 often used for historical maps georeferencing because: (1) they are easy to utilize and (2) they allow
170 historians to control the quality and reliability of the identified points using co-visualization between
171 both maps.

172 *Choosing a geometric transformation model*

173 Once an acceptable set of paired features has been identified, the last step is to compute the
174 transformation from the input map to the reference. Several transformation models have been
175 proposed: global transforms (affine, projective), global with local adaptations (polynomial-based) and
176 local transforms (rubbersheeting, Thin Plate Spline, kernel-based approaches, etc.). Studies have been
177 conducted to assess the relevance of these transformation for historical maps [18,21,22]. They show
178 that choosing a model is mostly a matter of compromise between the final spatial matching between
179 the feature pairs (*i.e.* the expected residual error) and the acceptable distortion of the map regarding
180 its legibility. Exact or near-perfect matching between features can be achieved with local transforms
181 and high order polynomials, whereas the internal structure of the map is most preserved by global
182 transformations. Low order polynomials offer a compromise between both constraints.

183 2.1.4. Temporalization: locating geohistorical sources in time

184 Georeferencing is a way of locating multiple maps in the same reference space. Similarly,
185 *temporalization* is the process of locating each geohistorical source in time. When building
186 spatio-temporal snapshots from historical maps, the key problem is to determine the moment where
187 the map is representative of the actual state of the area it portrays, *i.e.* the *valid time* of the map. We
188 considered the valid time of each map as the period starting with the beginning of the topographic
189 survey and ending with the publication of the map, which are often uncertain. Representing uncertain
190 or imprecise periods of time is a common issue when dealing with historical information and many
191 authors relied on the fuzzy set theory to represent and reason on imperfect temporal knowledge [23,24].
192 We model imprecise valid times as trapezoidal fuzzy sets. They are trapezoidal functions of time with
193 values ranging from 0 (the source provides no information at this time) to 1 (geographical entities
194 portrayed in the map are regarded as existing and tangible at this time). We rely on the pgSFTI⁶
195 postgres extension to store and manipulate such temporal fuzzy information. For instance, Figure 2

⁶ <https://github.com/OnroerendErfgoed/pgSFTI>

196 illustrates the *valid time* of a map whose topographic survey started in year 1775, ended between 1779
 197 and 1780 and which was engraved late 1780.

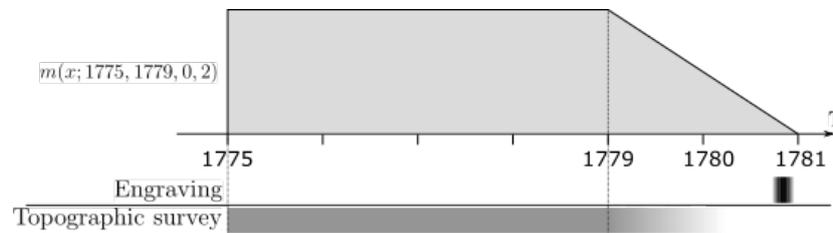


Figure 2. An uncertain valid time modelled as a trapezoidal fuzzy set function

198 2.1.5. Extracting information from maps

199 Once the historical maps have been georeferenced and temporalized, their cartographic objects
 200 can be extracted to produce geohistorical objects. The most classical way to extract information from
 201 maps is by human action with a classical GIS software (e.g. QGIS). However each one historical map
 202 of Paris contains a large amount of information to be extracted (e.g. tens of thousands of street names,
 203 hundreds of thousands of building number, etc.). A first solution is then to use computer vision and
 204 machine learning methods to create automatic extraction tools. These tools can process the whole map
 205 in a few hours. Regrettably, such tools are difficult to design, are very specific to each historical map,
 206 and may produce low quality results (see Figure 3). Recently, collaborative approaches have shown to
 207 be very efficient for building large geographical databases in a relatively short period (OSM⁷, NYPL⁸).



Figure 3. In this example, hand written text is automatically detected and extracted (red) from an historical map.

208 2.2. Modelling geohistorical objects

209 Information extracted from historical maps is used to create gazetteers. Those are made of
 210 geohistorical objects. To this end, we design a geohistorical objects model with all necessary attributes
 211 and also flexibility to adapt to the great variety of geohistorical object types and sources. Our goal is to
 212 provide a generic minimal (geo)historical object model that can be used by other and easily extended
 213 when necessary.

⁷ <http://www.openstreetmap.org>

⁸ <http://buildinginspector.nypl.org/>

214 2.2.1. modelling choices

215 Geohistorical data are extremely diverse, both in terms of historical sources and in terms of how
216 the sources were dealt with by historians. As such, historians use complex tailored models. We do not
217 aim at modelling all geohistorical data in all their complexity. Instead, we propose to model the bare
218 minimal common properties of all geohistorical objects, and offer mechanisms so this model can be
219 easily extended and tailored to the specificities of the data. To define the bare minimal model, we start
220 from the very nature of a geohistorical object, that is both an historical object and a geospatial object.
221 The extension mechanism is provided via a database-object oriented design using table inheritance,
222 and is packaged into a PostgreSQL extension⁹.

223 2.2.2. geohistorical objects model

224 Geohistorical objects have both an historical and a geospatial part. We stress that modelling
225 historical source and numerical origin process of a geohistorical object is an essential part. The detail
226 of the model are illustrated in Figure 4.

227 Historical aspect

228 In our model, an historical object is defined by its name, source and temporalization.

- 229 • *Name*. By name, we mean the historical name that was used to identify the object in the historical
230 source, and the current name that is used by historians to identify the object in the current context.
231 For instance, the historical name of the Eiffel tower in Paris may be "tour de 300 mètres", but,
232 today, it is referenced as "tour Eiffel".
- 233 • *Source*. A historical object is defined by a primary historical source (document), where the object
234 is referenced. Beside the historical source, the way the object was digitized in this source is also
235 essential. For instance, a street name may have the Jacoubet topological map as historical source,
236 and would have been digitized via collaborative editing on the georeferenced map.
- 237 • *Temporalization*. Any historical source is associated with temporal information (fuzzy dates),
238 which is the period during which the source is likely relevant. Beside the historical source
239 temporal information, a historical object can also have its own temporal information. For
240 instance, a street may have been extracted from a historical map having been drawn between
241 1820 and 1842. Besides this information, using other historical documents allow to narrow the
242 probable existence of this street to 1824-1836.

243 Geospatial aspect

244 A geohistorical object is also defined by geospatial information: a direct spatial reference
245 (geometry) and its positional accuracy metadata.

- 246 • *Geometry*. A feature has a geometry which follows the OGC standard¹⁰. It may be a point,
247 polyline, polygon, or a composition of any number of those, in a specified SRS. The geometry is
248 extracted from the historical source (in a manual or automatic way).
- 249 • *Positional accuracy*. Historical features have positional accuracy information. This precision
250 expresses the spatial uncertainty of the historical source (the person drawing the map may have
251 made mistakes) and the spatial imprecision of the digitizing process (the person editing the
252 digitised map may have made a mistake). One historical source may contain several accuracy
253 metadata, one for each geohistorical object type it contains. For instance, a historical map may
254 contain buildings and roads. Buildings may have a different positional accuracy (5 metres) than
255 road axis (20 metres). Besides, the digitising process precision may have been of 5 metres.

⁹ https://github.com/GeoHistoricalData/geohistorical_objects

¹⁰ <http://www.opengeospatial.org/standards>

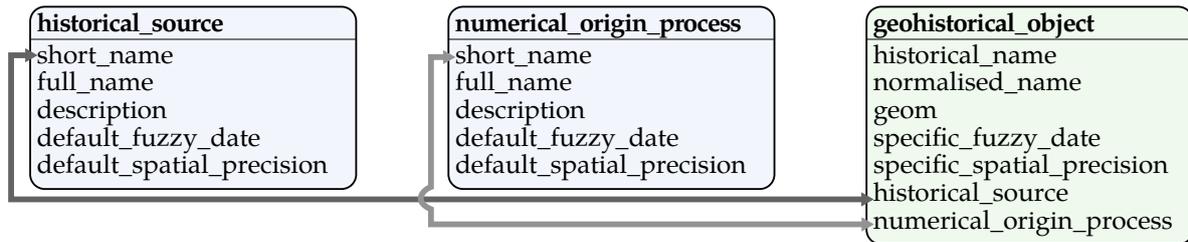


Figure 4. The geohistorical object model, where each object is characterized by its historical source (for instance the historical map the object was described in) and a numerical origin process, which is the process through which the object was digitized. Besides source and origin process, an object is also described by a fuzzy date, a text and a geometry.

256 2.2.3. A database of geohistorical objects

257 We define a conceptual schema for geohistorical objects, which is based on two names, a source, a
 258 capture process, fuzzy dates and a geometry. This defines the core of a generic geohistorical object.
 259 Yet this geohistorical object model is easily extendible using the table inheritance mechanism, an
 260 object-oriented design mechanism that is available in PostgreSQL (see Figure 5).

261 Table inheritance

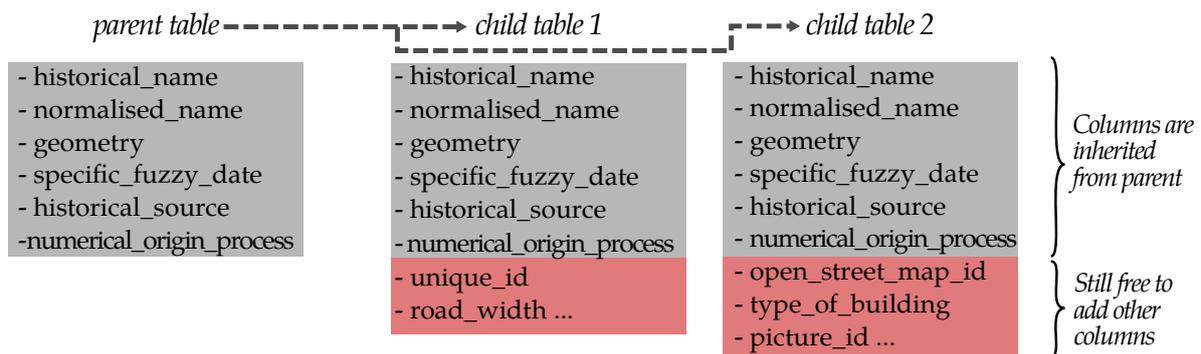


Figure 5. The table inheritance mechanism: a child table inheriting from a parent table inherits all the parent columns, and can also have its own.

262 The concept of table inheritance is simple. When a table *child* is created as inheriting from a table
 263 *parent*, *child* will have at least the columns of *parent*, but can also have other columns (provided there
 264 is no name/type collision). This means in our case that a table of geohistorical objects will inherit from
 265 the main geohistorical object table, *i.e.* will have all the core columns of geohistorical objects (names,
 266 sources, temporal aspect, spatial aspect), but can also have its own tailored column, providing the
 267 necessary flexibility.

268 Another key aspect of table inheritance is that the *parent* table is queried, the query will be
 269 executed on not only the rows of *parent* table, but also on the rows of all *child* table. This means that
 270 all tables using the geohistorical object model will be virtually grouped and accessible from one table.

271 Simulated inheritance of index and constraints

272 The PostgreSQL table inheritance mechanism is however limited in some aspects, because
 273 constraints and index can not be inherited. Constraints are essential, because they are used to
 274 guarantee that any geohistorical object will correctly use existing sources from the sources tables
 275 ("historical_source" and "numerical_origin_process"). Indexes are also essential, because when using

276 hundred of thousand of geohistorical objects, they are needed to help speed the queries. We index
277 not only names, but all geohistorical object core columns (names, sources, temporal aspect, spatial
278 aspect). We propose a registering function that the user can execute only once when creating a new
279 geohistorical object table.

280 *Modelling a geohistorical object from the user perspective*

281 The practical steps to create geohistorical objects are simple:

- 282 1. Add the historical source and numerical origin process in the source and process tables.
- 283 2. Create a new table inheriting geohistorical objects and containing your additional custom
284 columns
- 285 3. Use the registering function with this table name
- 286 4. Insert your data in the table.

287 *2.3. Geocoding historical addresses with geohistorical object gazetteers*

288 In the previous section, we explained how we create gazetteers of geohistorical objects from maps.

- 289 1. an historical map is scanned,
- 290 2. scans are georeferenced using hand picked control points,
- 291 3. historical work allow to estimate temporal information and spatial precision of the map,
- 292 4. roads name and axis geometry is extracted from the scan (manually or automatically),
- 293 5. building numbers are extracted from the scan (manually or automatically),
- 294 6. in some cases, building numbers can be generated from the available data,
- 295 7. normalised names are created from historical names,
- 296 8. geohistorical objects are created.

297 The next step is to use these gazetteers to geocode historical addresses.

298 *2.3.1. Historical geocoding concept*

299 In our method, geocoding something is finding the most similar geohistorical objects within the
300 available gazetteers, which then provides the geospatial information. This approach relies on two key
301 components: gazetteers of geohistorical objects, and a metric to find the best matches. This approach
302 allows to perform geocoding in a broad sense, as it does not rely on a structured address (number,
303 street, city, etc.), but rather on a non constrained name.

304 *2.3.2. Creating geohistorical object gazetteers for geocoding*

305 geohistorical object gazetteers are key for the geocoding. These objects are extracted from
306 topographical historical maps and inserted into geohistorical objects tables. Each table form a gazetteer.

307 *Database architecture for geocoding*

308 We again use the PostgreSQL table inheritance mechanism. To this end, we create two tables
309 dedicated to geocoding. Now gazetteers tables that will be used in geocoding must inherits from these
310 two tables. "precise_localisation" table is for building number geohistorical objects, e.g. "12 rue du
311 temple, Paris". "rough_localisation" table is for road axis, neighbourhood, cities geohistorical objects.
312 We chose to have two separate tables for ease of use and performance. Geocoding queries are then
313 performed on the two parents tables, but thanks to inheritance, these parents tables virtually contains
314 all the gazetteers table containing the actual geohistorical objects, as illustrated in Figure 6.

315 *2.3.3. Finding the best matches*

316 Once geohistorical objects gazetteers describing precise and rough localisation are available,
317 geocoding is finding the best match between the input query and the objects.

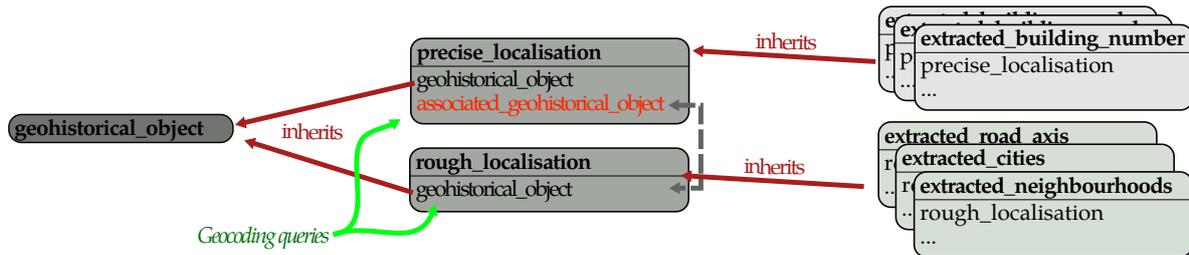


Figure 6. Geocoding table architecture. Two tables of `geohistorical_object` are the support for geocoding queries. Because all extracted geohistorical objects tables inherits from these two tables, they both virtually contains all the objects.

318 Concept

319 We call the potential matches "candidates", and the problem is then to rank the candidates from
 320 best to worst. The user can chose how many candidates he wants, depending on the application. For
 321 an automated batch geocoding, the best match (top candidate) is optimal. For a human analysis of
 322 data, several matches may be more interesting (top 10 candidates for instance). What can be qualified
 323 as "best" depends on the user expectations. We provide a number of metrics than can be combined by
 324 a user into a tailored ranking function. The function is expressed in SQL, with access to all postgres
 325 math functions. We describe the available metrics and give example of such function.

326 Example

327 For instance when a user geocodes the address "12 rue de la vannerie, paris" in 1854, user may be
 328 more interested into geohistorical objects that are textually close (e.g. a geohistorical object "12 r. de la
 329 vannerie Paris", 1810), or maybe geohistorical objects that are temporally close (e.g. "12 r. de la Tannerie
 330 Paris",1860).

331 Metric: string distance w_d

332 We use the string distance provided by the PostgreSQL Trigram extension (pg_trgm¹¹), which
 333 compares two strings of characters by comparing how many successive sets of 3 characters are shared.
 334 For instance "12 rue du temple" will be farther away from "12 rue de la paix" than from "10 r. du
 335 temple".

336 Metric: temporal distance t_d

337 Both the address query and the geohistorical object are described by fuzzy dates. In order to
 338 compare such temporal information, we propose a simple fuzzy date distance that casts fuzzy dates
 339 into polygons. The x axis is the time, and the y axis is the probability of existence of the object. Then
 340 the distance between twon dates A and B is computed as $\text{shortest_line_length}(A,B) + \text{Area}(A) - \text{Area}(A \cap B)$.
 341 Note that this distance is asymmetric.

342 Metric: building number distance b_d

343 To get building number distance, a function first extracts the building number both from the input
 344 address query (b_i) and from the geohistorical object (b_d). If b_i and b_d have same parity, the distance
 345 is $|b_d - b_i|$. If parity is different, the distance is $||b_d - b_i| + 10|$. In France, building numbers have
 346 in general the same parity on each side of the street (e.g. Left : 1,3,5,.. ; Right: 2,4,6..). We analysed

¹¹ <https://www.postgresql.org/docs/current/static/pgtrgm.html>

347 current building number in Paris and determined that on average, given a building number b_i , the
 348 closest building number with a different parity has a 10 number difference.

349 *Metric: positional accuracy s_p*

350 Another way to rank the geohistorical objects is to use their positional accuracy. The positional
 351 accuracy of a geohistorical object is either the positional accuracy computed for this objects when it is
 352 available, or the default positional accuracy of its geohistorical source.

353 *Metric: level of detail distance s_d*

354 Providing localisation information at different level of detail, depending on the user requirement
 355 is an important quality issue for our geocoder. For instance if the level of detail of the user's query data
 356 is the city, there is no need to perform a more precise geocoding. Therefore the user can specify a target
 357 scale range (S_l, S_h). Then given a geohistorical object whose geometry is buffered ($geom_b$) with its
 358 spatial precision, the scale distance is defined by $least(|\sqrt{area(geom_b)} - S_l|, |\sqrt{area(geom_b)} - S_h|)$.
 359 The formula $\sqrt{area(geom_b)}$ gives an idea of the spatial scale of the geohistorical object.

360 *Metric: geospatial distance g_d*

361 The user may provide an approximate position for the area he is interested in. For instance in
 362 France both city "Vitry-le-François" (East) and "Vitry-sur-Seine" (near Paris) exist, but are very spatially
 363 far away. A user expecting results in the Paris area may provide a geometry (a point for instance)
 364 near Paris. Then the classical geodesic distance is computed between the provided geometry and the
 365 candidates geohistorical object.

366 *Example of matching function*

The different metrics can be weighted and combined depending on the user needs. Equation 1
 gives an example that favour good string similarity, but not at the price of a large temporal distance.

$$100 * w_d + 0.1 * t_d + 10 * n_d + 0.1 * s_p + 0.01 * s_d + 0.001 * g_d \quad (1)$$

367 2.4. Collaborative editing of geohistorical objects

368 The geocoding approach we have presented in the previous section works inside a PostgreSQL
 369 database. Given an input address and fuzzy date, plus a set of parameters, it returns the geohistorical
 370 objects that matches the input the best. Yet the geocoding results are only as good as the gazetteers
 371 used. The geohistorical objects within the gazetteers may be spatially imprecise, mistakenly named
 372 or simply missing. Given that the volume of geohistorical objects is large (for Paris, approximately
 373 50 k building number per historical map), we create a collaborative platform to facilitate geocoding,
 374 visualising the results and editing the geospatial objects when necessary. To this end, we create a
 375 dedicated web application so collaborative editing is possible without having to install specific tools.

376 2.4.1. About collaborative editing

377 Given the complexity of calibrating automatic extraction tools on specific maps and their relative
 378 reliability, the collaborative digitisation of vector objects from maps is a safe alternative. For instance,
 379 we used such an approach in order to extract the main feature of the Cassini maps (18th century
 380 France) [25]. Furthermore, the results of the collaborative extraction of features can then be used to
 381 test, calibrate or train automatic extraction algorithms.

382 2.4.2. Collaborative editing architecture

383 Figure 7 outlines the architecture used for collaborative editing.

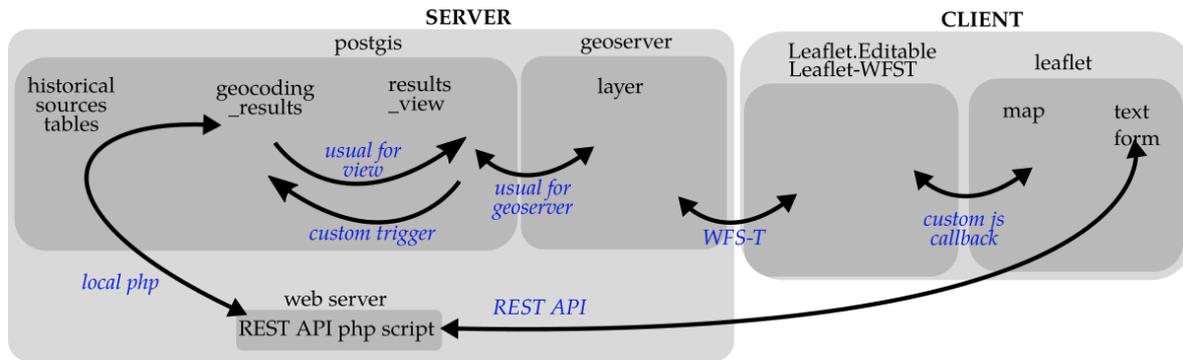


Figure 7. Conceptual architecture for interactive display and edit of geocoding results. The stack contains only standard components.

384 Architecture

385 The heart of the architecture is a PostgreSQL database server, which contains the geohistorical
 386 objects gazetteers that will be used for geocoding as well as the geocoding function. A webservice can
 387 geocode addresses and return results via a REST API. However, the webservice has another option
 388 where the results are not returned, but instead written in a result table along with a random unique
 389 identifier (RUID). The RUID is then the key that permit to display and edit the results. To this end, a
 390 geoserver can access (read and edit) the result table via the WFS-T protocol. A web application based
 391 on Leaflet then acts as a user interface to display and edit the results via the geoserver.

392 Persistence of geocoding results and edits

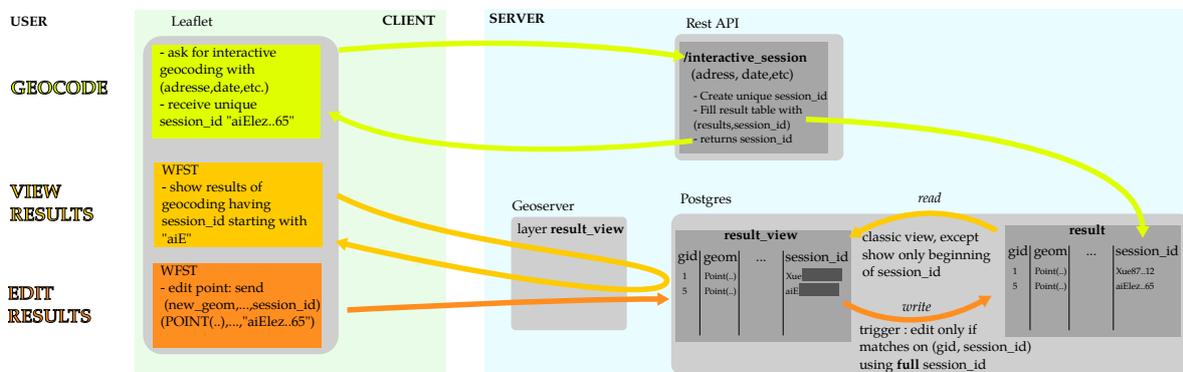


Figure 8. Collaborative display and edit is achieved through a mix of standards (REST, WFS-T) and custom solution (triggers) that enable the sharing of a basic public/private key.

393 The architecture that allows persistence of results is illustrated in Figure 8. When using the
 394 RUID mechanism, each geocoding result (that is the found geohistorical object from the gazetteers) is
 395 associated to this RUID. That way the user can always access its results, regardless of the computer
 396 session or browser cache issues.

397 To edit, a specific mechanism is used. The user does not directly edit the result table, as he could
 398 potentially edit other people results. Instead, the user edit a dedicated result_view that acts like a
 399 bouncing. It allow edit only if the edit is occurring on a row that has the user RUID. User edit of the
 400 geospatial objects do of course not affect the source data, for a tracking purpose.

401 Instead, a new user edit automatically creates an edited copy of the geohistorical object in a
 402 dedicated table "user_edit_added_to_geocoding" that is a gazetteer and is used by the geocoding

403 process. In this table are inserted the edited geohistorical objects. The objects retain their
 404 "historical_source", but their "numerical_origin_process" is changed to properly document the fact that
 405 they are the result of a collaborative editing.

406 2.4.3. Collaborative editing user interface

407 We consider that building an efficient user interface is very important for historical geocoding. In
 408 particular, many end users are specialised on history rather than on computer science, and thus an
 409 easy access to geocoding is essential. All our interfaces are web-based for a maximum of compatibility.
 410 We propose three interfaces whose results are shared.

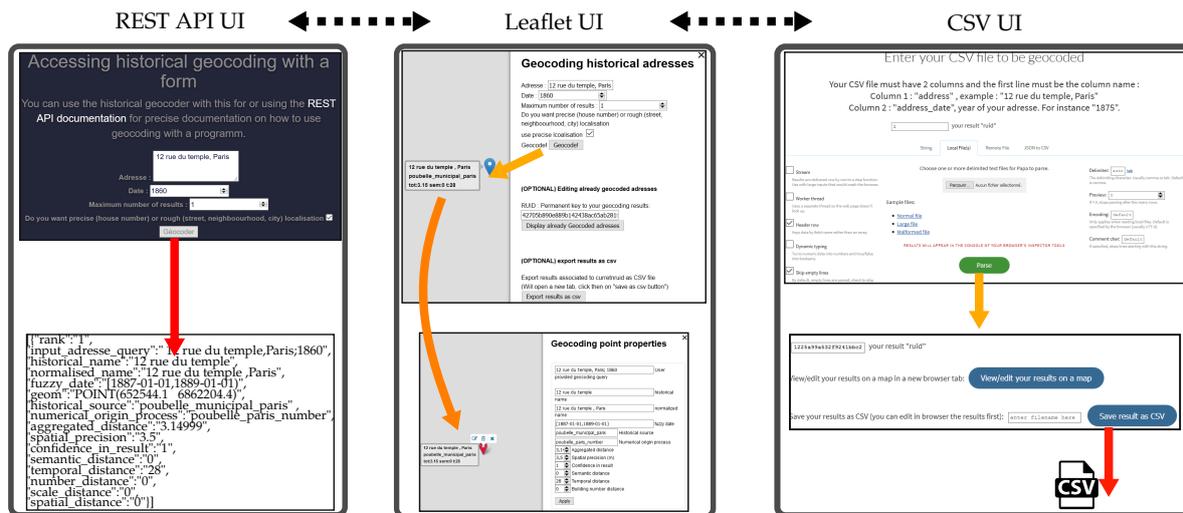


Figure 9. Various Web User Interface for use of historical geocoding.

411 Interface for a REST API.

412 The simplest interface we propose is a form that helps build the necessary REST API parameters.
 413 Indeed, a REST API works via URL containing precise parameters, and it can be tedious to manipulate.
 414 For instance:

415 [https://www.geohistoricaldata.org/geocoding/geocoding.php?adresse=2012ruedutemple,Paris&](https://www.geohistoricaldata.org/geocoding/geocoding.php?adresse=2012ruedutemple,Paris&date=1860&number_of_results=1&use_precise_localisation=1)
 416 [date=1860&number_of_results=1&use_precise_localisation=1](https://www.geohistoricaldata.org/geocoding/geocoding.php?adresse=2012ruedutemple,Paris&date=1860&number_of_results=1&use_precise_localisation=1)

417 This interface is designed to be used in an automated way, for batch geocoding.

418 Interface for batch geocoding via CSV files.

419 In our experience historian often work with spreadsheet files, where each line will be a potential
 420 historical object, along with an address and a date. To facilitate the geocoding of these addresses,
 421 we propose a User Interface that can read Coma Separated Value (CSV) files (which is a standard
 422 spreadsheet format), and geocode the address and date within. This interface is built around
 423 PapaParse¹² Javascript framework. Then, the geocoding results can be either downloaded as a
 424 CSV file, or displayed and edited in a web application.

¹² <http://papaparse.com>

425 Interface for display and edit of results.

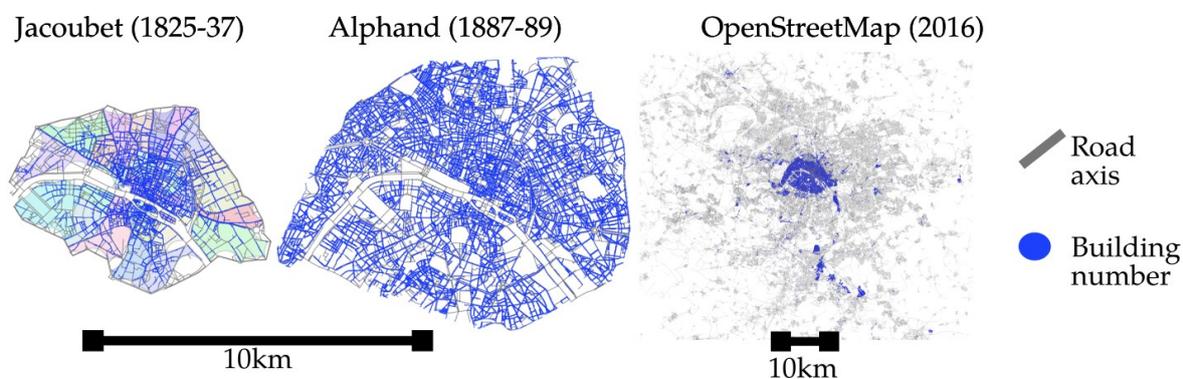
426 The most complex interface we propose is based on the Leaflet¹³ Javascript framework. There, the
 427 user can geocode an address, or use already geocoded address via the RUID mechanism (see Section
 428 2.4.2.1), be it from previous sessions or from geocoded CSV files. The geocoding results are displayed
 429 on top of a relevant historical map, and can be edited. User can edit results geometry as well as results
 430 names (historical and normalised). We stress that although such edit are stored in the database, and
 431 used by further geocoding queries, they do not affect source data, by design.

432 3. Results

433 We perform several experiments to validate our approach. First we use the geohistorical model to
 434 integrate objects extracted from historical topographical maps from the 19th century for the city of
 435 Paris, and the current OpenStreetMap road axis and building numbers for Paris city surroundings.
 436 We successfully integrate the road axis, building numbers, and neighbourhoods to the geocoder
 437 sources. We then perform multiscale geocoding of dozens of thousand of historical addresses extracted
 438 manually by historians and extracted automatically by automatic process. For one of our datasets, we
 439 ask the historian to manually correct the automated geocoding results, so as to evaluate the quality of
 440 our method. Last, we test the collaborative editing of geohistorical object in two scenarios: analysis
 441 (several results for one address), and edit (efficiency of check/edit top results for several addresses).

442 3.1. geohistorical objects sources

443 We mainly use three historical sources of geohistorical objects to perform geocoding. The first two
 444 are Historical topographic maps of Paris from the 19th century. These maps are georeferenced then
 445 street axis (and possibly building numbers) are manually extracted. The third historical sources are
 road axis and building number for Paris surrounding extracted from current Open Street Map data.



446 **Figure 10.** geohistorical objects used from geocoding extracted from the source maps.

447

447 3.1.1. Historical topographic maps used

448 We integrated two major French atlases of Paris from the 19th century as geohistorical sources.
 449 The first one is the "Atlas municipal de la Ville, des faubourgs et des monuments de Paris"¹⁴ created
 450 at the scale of 1 : 2000 between 1827 and 1836 by Theodore Simon Jacoubet, an architect who was
 451 working for the municipal administration of Paris. The second atlas is the 1888 edition of the "Atlas
 452 municipal des vingt arrondissements de la ville de Paris"¹⁵. For legibility reasons, we refer to the first

¹³ <http://leafletjs.com>

¹⁴ Municipal atlas of the city, suburbs and monuments of Paris.

¹⁵ Municipal atlas of the 20 districts of Paris

453 atlas as the "Jacoubet atlas" and the second as the "Alphand atlas"¹⁶. The Jacoubet atlas depicts a city
 454 standing between the housing development following the sale of the properties confiscated during the
 455 French Revolution and the majors changes in the urban structure arising from the emergence of the first
 456 train stations in 1837-1840 and the so-called Haussmannian transformations.

457 The Alphand atlas is a portrayal of Paris at the scale of 1 : 5000, after most of the Haussmannian
 458 transformations (major rework of Paris urbanism in the 19th century) and after the city was merged
 459 with 11 of its neighboring municipalities in 1860. Both atlases contain large scale topographic views of
 460 Paris, separated in several sheets (54 and 16 respectively) and portray the urban street network with
 461 each street named, building of public purposes and religious buildings (see Figure 11). In addition, the
 462 house numbers are specified for most of the streets in the city, although the Alphand atlas pictures only
 463 the numbers at the start and end of each street section. Both atlases are also built upon triangulation
 464 canvas covering the entire city, allowing us to expect a high positional accuracy of the geographical
 465 features they contain.

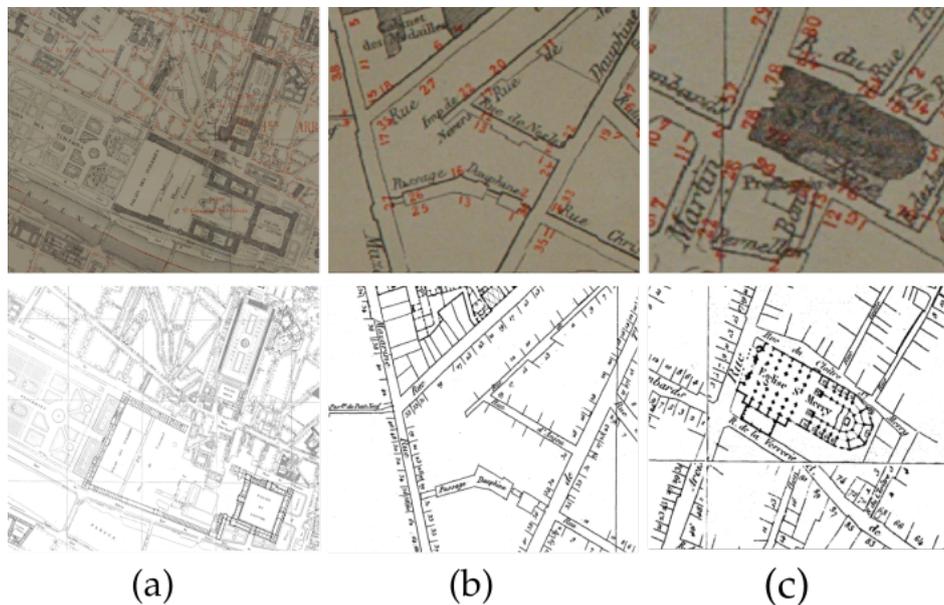


Figure 11. Samples of the georeferenced Alphand Atlas (2nd row) and Jacoubet Atlas (1st row) at different scales: district (a) and urban islet (b). Column (c) shows how buildings are portrayed in the maps.

466 We georeferenced the two atlases using the grids drawn on the maps, which are aligned on the
 467 Paris meridian, as a pseudo-geodetic objects to identify feature pairs. The dimensions of the grid cells
 468 also appear on the maps, allowing us to reconstruct the grids in a geographic reference system. We
 469 have chosen to georeference the maps in the Lambert I conformal conic projection, which uses the
 470 Paris meridian as prime meridian and rely on the NTF (Nouvelle Triangulation Française) geodetic
 471 datum. The main advantage of this projection is that it is locally close to the planar triangulation of
 472 Paris used in the atlases. Thus, the projection of the maps can be reasonably approximated by the
 473 Lambert I projection, making the reconstruction of the grids in the target coordinate reference system
 474 straightforward. In addition, since both maps are at high scale and are reliable because they are official
 475 maps with high positional accuracy, we used rubbersheeting as the geometric transform model. The
 476 georeferencing process applied for each atlas was the following:

¹⁶ From the name of Jean-Charles Alphand who was at the time the director of the department of public works of Paris.

- 477 • reconstruct the meridian-aligned grid with Lambert I coordinates;
- 478 • in each sheet, mask the non-cartographic parts out (cartouche, borders,etc.);
- 479 • for each sheet, set pairs of ground control points at each intersection between the vertical and
- 480 horizontal lines of the grids in the map and in the reconstructed grid;
- 481 • transform each sheet with a rubbersheeting transform based on the ground controls points
- 482 previously indentified on the grids.

483 geohistorical objects extraction

484 Based on these atlases, vectorial road axis are manually drawn and the road name inputted For
 485 Alphand map, the building number at the beginning and end of ech street segment is also inputted.
 486 For Jacoubet, the building numbers from a previous map (Project Alpage, Vasserot map, [26]) are
 487 adapted to fit the Alphand map. Multiple series of successive checking and editing are performed
 488 using ad hoc visualisations and tools.

489 For Alphand, building numbers are then generated based on the available information (for each
 490 street segment, for each side, beginning and ending number) by linear interpolation, and an offset.
 491 The size of the offset is estimated by using current Paris road width when the road has not changed to
 492 much.

493 3.1.2. Other geohistorical sources

494 We also use current data from OpenStreetMap. We use the version of the data that has been
 495 transformed to be used by the Nominatim geocoder. Custom scripts extract road axis and building
 496 numbers. The dataset covers Paris city and its surroundings, and is dated to 2016.

497 3.2. Geocoding of Historical datasets

498 One of the end goal of our geocoding tool is to be useful for historians. Therefore, we contacted
 499 several historians working on Paris (19th century). They had been collecting historical addresses,
 500 which we geocoded by importing their data into the geocoding server. Figure 12 shows an extract of
 501 the thousands of geocoded addresses, while table 1 gives an overview of the number of successes and
 502 timing.

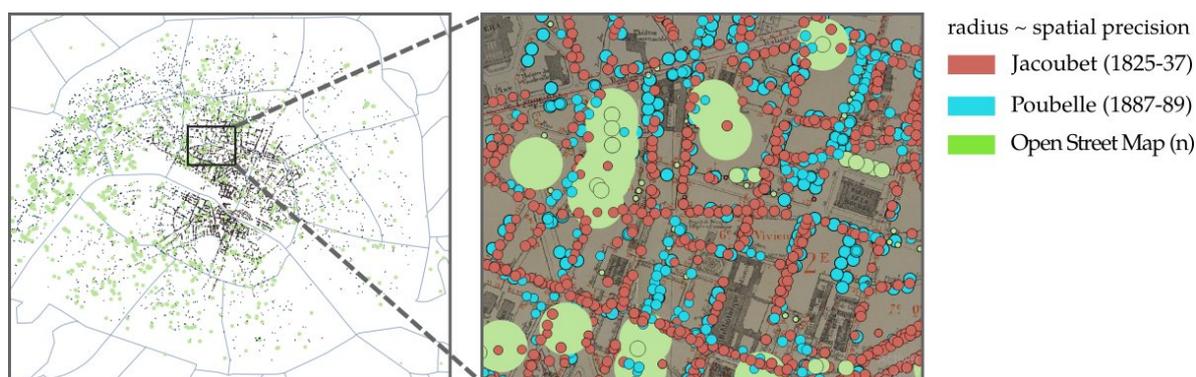


Figure 12. All geocoded historical datasets. Size is proportional to spatial precision.

503 3.2.1. Manually collected dataset

504 **South Americans:** South America immigrants living in Paris in 1926, manually input from census,
 505 collected by Elena Monges (EHESS).

506 **Textile:** Professionals of textile industry in Paris, manually input from the "Almanachs dy Commerce
 507 de Paris", from 1793 to 1845, collected by Carole Aubé (EHESS).

508 **Artists accommodations:** Addresses of artist studios and artists accommodations between 1791 and

Table 1. All geocoded historical datasets facts.

Dataset name	input addresses	response rate (rough)	secs/1000 addresses
South Americans	13991	13743 (250)	138
Textile	5777	5688 (16)	135
Textile 2	3070	3053 (2)	110
Artists accommodations	13907	10215 (2955)	244
Health administrators	1887	1698 (171)	316
Belle epoque (0.3)	6467	3880(337)	280
Belle epoque (0.5)	6467	6000	351

509 1831, collected by Isabelle Hostein (EHESS) to study their impact on Paris development.

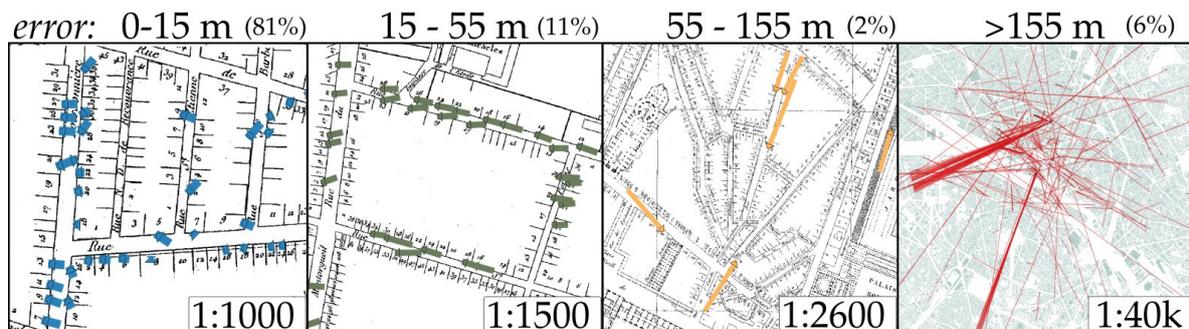
510 **Health administrators:** Addresses of public health and hygiene administrators in Paris between 1807
511 and 1919 ([27]), collected by Maurizio Gribaudo and Jacques Magaud (INED-EHESS).

512 3.2.2. Belle epoque

513 We geocode another set of addresses that are automatically extracted from directory of Paris
514 financial societies between 1871 and 1910. Directories are books referencing address of company (and
515 name and other information). The process of automatic extraction is complex in itself (Project Belle
516 Epoque, [28]), and is out of scope of this article. We only describe it briefly here.

517 First, each page of the directories of Paris for specific years are photographed. Pictures are then
518 straightened, and information is extracted via an OCR software which has been configured for the
519 directory specific layout. Further rule based processing parse the text into address fields. As a result
520 of this automatic process, the quality of addresses is often significantly lower than manually edited
521 addresses. Therefore, we test two settings by allowing a greater maximum string distance from 0.3 to
522 0.5 (over 1).

523 3.3. Manual editing of the geocoding results for evaluation

**Figure 13.** An historian manually move the geocoded addresses.

524 For one of the data set (Textile 1 and 2), the historian manually correct the geocoded results. We
525 then plot the segment between address point resulting of automated geocoding and address point
526 after manual editing. Results are presented in the table 2 and in Figure 13 We classify the results based
527 on the length of this segment (*i.e.* the error in meter the geocoding method made).

- 528 • When the edit move the adress point less than 15 meters, we can consider that the edit is mostly
529 about small moves , for instance centering the point on the building limit.
- 530 • Between 15 and 55 meters, the correct street has been found, but the building numbers are slightly
531 misplaced (a few numbers).
- 532 • Between 55 and 155 meters, in most cases the street is correct, but the building numbers are far
533 from their correct position.

Table 2. Evaluating the error of geocoded results, via the dist. (geographic distance) of edit, the percentage of the total 8823 addresses, the average aggregated distance score, the average string distance w_d , the average temporal distance t_d , and the subjective most common edit reason we encountered while browsing the data

dist. (m)	%	avg(agg)	avg(sem)	avg(tempo)	main edit cause (subjective)
0 - 15	81 %	9.4	0.07	19.5	moving point on building limit
15- 55	11 %	12.4	0.09	27.2	small numbering editing (same street)
55 - 155	2 %	23.7	0.14	41.2	large numbering editing (same street)
155 - 7.2k	6 %	26.9	0.18	49.1	editing street

- Above 155 meters, streets are wrong in most of the case.

We stress that given Paris building average size, and the lack of precise definition of an address (is it the position of the door, of the center of the building,...?) results up to 55 meters could be considered as very close to ground truth.

3.4. Collaborative editing

We propose several User Interface for easy geocoding, and collaborative editing of the geocoding results. We informally tested the interfaces and found that they facilitate geocoding, especially for the batch mode. We also test the collaborative editing in two use cases. In the first use case a specialised user geocodes a single address and display the top 3 results corresponding to this address. The user is expert and its goal is geocode an address and assess the reliability of the result at the same time. In the second use case, a user batch geocode several addresses (30), looking of the best result for each address. Then the user display the results on the map and check/edit the addresses.

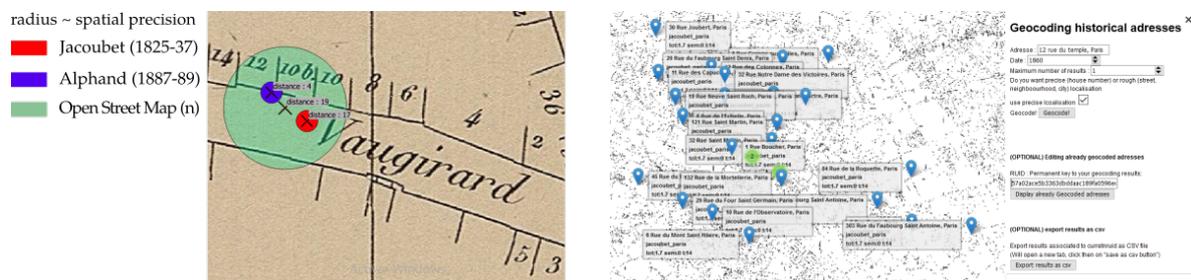


Figure 14. Two use cases: First use case, an expert geocodes an address and analyse the top 3 results to assess the reliability of the result. Second use case: a user batch geocodes 30 addresses (1 result per address) in Paris and check/edit the results.

545

3.4.1. Use case 1: top 3 results for one address

Using the web application, we geocode the address "10 rue de vaugirard, paris" for the date 1840, and ask for the top 3 results, as shown in first part of illustration 14. A matching building number geohistorical object exists in the three gazetteers extracted from the three maps. Based on the results, we can safely assume that this building number has not changed during the last 2 centuries.

3.4.2. Use case 2: batch geocoding of 30 addresses and check/edit

In this use case, a regular user is to check/correct 30 random addresses from the Jacoubet map using the web application. The task is performed quickly, the check and edit of each address is a matter of a few seconds. The main time consuming task is the loading of the background historical map, due to unfortunate hardware limitations. The edit speed seems to be on par with a desktop based edit solution (using QGIS).

546

557 4. Discussion

558 4.1. Genericity

559 Reaching a more generic geocoding service is important if we want to make it usable in other
560 contexts and to profit from the various sources of knowledge on past spaces.

561 4.1.1. Geohistorical sources and data

562 *Using external resources from the Web of data as new sources*

563 Besides features representing address points and streets, georeferenced features of other types
564 could be used with benefits by the geocoding service. As a matter of fact, people often refer to
565 places of interest, such as famous buildings, monuments like statues or fountains or even named
566 neighbourhoods to describe their position in space. We thus consider adding data about places of
567 interest to improve our geocoding service. Like the data that was used to build the geocoder, such data
568 could be gathered from ancient maps. But they may also come from existing gazetteers and knowledge
569 bases published on the Web of data, such as DBpedia¹⁷, Yago¹⁸, the Getty Thesaurus of Geographical
570 Names¹⁹ or the gazetteer of place names published by the French National Library²⁰.

571 *Widen the spectrum of cartographic sources*

572 We exploit Jacoubet and Alphanand maps, yet there are several more to be exploited toward the end
573 of the 19th century, and in the beginning of the 20th century. From the beginning of the 20th century,
574 Paris city administration produced a map per year. Of course, the main improvement direction would
575 be to add maps of other cities/countries! For France at least, major cities have often been mapped
576 starting from 1900.

577 Before the beginning of 19th century, the address system was very different in Paris. In mid 18th
578 century, the address system was in fact that each building would have a specific name (no number, no
579 notion of street name) in its neighbourhood. Our geocoding system has also been designed with this
580 type of addressing but it has not been tested yet. More generally, this type of indirect localisation is
581 very close to the field of web of knowledge.

582 *Diversity in geohistorical objects natures*

583 In this article several type of geohistorical objects were used for geocoding: building numbers,
584 streets axis, neighbourhood. Other datasets were investigated as well, such as the city limits extracted
585 by the project Geo Historical Data in a collaborative way from the Cassini maps [25]. In fact, a compiled
586 version of city limits (GeoPeople project [29]) from 1793 to 2010, created by EHES, has also been
587 tested. But building cadastres could also be integrated so as to have a building layout associated to
588 an address rather than a point, which would solve an old problem of address points. Indeed, there
589 is currently no consensus as to where a building number address point should be positioned: on the
590 entry door, on the letter box, etc. More excitingly, in some cases, more precise data is available, giving
591 the layout of apartments in buildings.

17 <http://wiki.dbpedia.org/>

18 <http://www.mpi-inf.mpg.de/departments/databases-and-information-systems/research/yago-naga/yago/#c10444>

19 <http://www.getty.edu/research/tools/vocabularies/tgn/>

20 <http://data.bnf.fr/>

592 4.1.2. Genericity in usages

593 *Named Entity Linking*

594 As we mentioned in section 4.1, people often refer to place names to describe their position in
595 space. The task of retrieving place names in a gazetteer or in a knowledge base, also known as (Spatial)
596 Named Entity Linking or toponym resolution, is a widely used way of desambiguating mentions
597 of spatial named entities extracted from texts by means of natural language processing approaches
598 for information retrieval, information extraction or document indexing purposes [30]. As we plan
599 to upgrade our geohistorical database with data about places of interest, we also have to adapt our
600 geocoding service in order to make it retrieve reference data stored in the database and corresponding
601 to place names mentions proposed by the users. Spatial Named Entity Linking implies solving issues
602 related to places names inherent ambiguity [31], such as the fact that a place may have several names
603 or the fact that several places may be designated by the same name. For each spatial named entity
604 mention to be disambiguated, unsupervised state of the art approaches first select candidates from
605 the gazetteer based on character string similarity. Then, they introduce additional criteria in order to
606 decide which candidate is the best reference for a given place name, usually taken from the textual
607 context of the mention [32,33]. In cases where textual context is very limited, like in tweets or location
608 descriptions extracted from directories, this step of candidate ranking reveals even more challenging
609 [34].

610 *Analysis tool of the cartographic sources content*

611 It is interesting to look at what historical sources were the most used for geocoding, although
612 the historical source are chosen based on a complex ranking function. If we take the example of the
613 over 10k geocoded addresses from the "Artists accommodations" dataset, we could expect all of the
614 results to be drawn from the Jacoubet map, as the dataset is between 1793 and 1836, and the Jacoubet
615 map is also in this range. Yet, analysing the results shows that if Jacoubet was used for 80% of the
616 addresses, Alphanth was used for 15%, although the map comes 30 years after. More surprisingly, the
617 OpenStreetMap current data is still used for 5% of addresses, although it is about 2 centuries after the
618 dataset.

619 Similar analysis on other datasets show similarly that all maps are always used, with of course a
620 focus on the temporally closest map. Interestingly, these results are in agreement with similar work as
621 presented in [35], chapter 4, where a prototype of multi-temporal geocoding is proposed. The approach
622 shows that for different datasets, all references maps (Jacoubet, Alphanth and BDAdresse (2010)) are
623 used, with proportions depending on the parameters considered and the weights of each criteria. We
624 think that this results are explained by the fact that historical maps miss some information, contain
625 errors, and do not have the same geographical coverage.

626 4.2. Quality of the geocoding

627 4.2.1. Increasing the quality of the gazetteers

628 *Collaborative enrichment*

629 We propose several ways to use the geocoding capabilities in an easy way through web based
630 User Interfaces. As we propose prototypes, the experiments are merely proofs of concepts for the
631 moment. For a real validation, a complete user study would be required, which is outside of the scope
632 of this article.

633 *Cross-referencing historical topographic maps*

634 One way to improve quality of available historical data is by using advanced crossreferencing.
635 Indeed, the process of linking and merging similar data from heterogeneous datasets, which is
636 called data conflation, enables to transfer information from one feature to the another, and thus may
637 brings additional knowledge about data imperfections without using ground truth data which are
638 non-existent for geohistorical data. For instance, [35,36] proposed an aggregated spatio-temporal
639 graph to merge and confront historical road networks. This process can reduce data heterogeneity and
640 allow to detect aberrations such as toponymic or numbering errors, or doubtful temporal trajectories
641 of objects like short disappearances, thereby leading to better data quality. Advanced cross-referencing
642 also makes it possible the construction of a genealogy of addresses by considering temporally linked
643 addresses, that can deal with toponymic evolution or changes in addressing systems or numbering of
644 buildings, thus paving the way for better spatio-temporal geocoding result.

645 4.2.2. Communicating the reliability of a geocoding

646 *Geocoding qualification and quality measures*

647 Modern geocoders are evaluated by how often they find a localisation, and how precise is the
648 localisation they return (see [37] for instance). The first criterion shows how able to retrieve an address
649 the geocoding algorithm is and also how exhaustive the gazetteer is. The second criterion refers to
650 the positional accuracy of the gazetteer. Using such quality evaluation measures that encompass
651 both the algorithm results and the gazetteer completeness makes the evaluation of their respective
652 quality impossible. Contrary to that, in the field of named entity linking, distinct quality evaluation
653 measures have been proposed for to the entity retrieval algorithm, like the measures proposed by [33]
654 and completed by [38], and for the reference knowledge base (see [39] for knowledge bases general
655 quality measures and [40] to evaluate the fitness of some knowledge bases for a given named entity
656 linking task).

657 *Geovisualisation*

658 The prototype of graphical user interface we propose could be improved in several ways. The
659 goal would be to efficiently provide information to user about the quality of geocoding, and the context
660 of results. First, the point displayed to represent the result could have a size proportional to estimated
661 spatial precision. This would help to visually assess the relevance of information. Second, the result
662 could be colour-coded to represent the temporal proximity with the input date. In a similar spirit, when
663 multiple results are proposed, a time slider would be most useful to graphically disambiguate between
664 result candidates. Third, the background historical map displayed in the prototype is currently set. Yet,
665 the most appropriate background map could be automatically displayed based on the input addresses
666 date provided by the user. Last, the current prototype becomes easily cluttered when displaying a
667 great amount of labels. Several strategies could be used, such as a better clustering of spatially close
668 results, shorter labels, or better labels placement.

669 4.2.3. Integrating user correction into historical sources

670 In collaborative editing, edit come from untrusted sources. Validating edits and solving conflicts
671 is then a classical problem. In our prototypes, every user edit is potentially used by the geocoder
672 (they are added to a dedicated gazetteer). We could use a voting scheme where edits are only taken
673 into account when a sufficient number of user have made them. However, we stress that due to
674 the number of data to edit (several hundred thousands building numbers), we prefer to rely on the
675 user benevolence, by considering that user spending time editing centuries old historical data are
676 committed to accurate editing.

677 4.2.4. Scalability

678 The main design choice of our geocoding architecture is to use a flat model for the address (an
679 address is any set of characters), as opposed to current geocoder which are highly hierarchical (an
680 address refers to a street, that refers to a neighbourhood, etc.). This modelling choice gives the freedom
681 that is necessary for data as incomplete as the historical ones, but also comes with a tradeoff regarding
682 scaling capabilities. Indeed, for strongly hierarchical data, it is possible to have separate databases for
683 each city for instance, thus preventing one database to grow too much, and ensuring a nice scaling
684 capability.

685 This is not however the case with our architecture. By using database indexes, we can theoretically
686 guarantee a fast geocoding time for up to few dozen of millions of geohistorical object used as sources.
687 The main bottleneck in this case is not the temporal aspect (it relies on PostGIS geometry, which enable
688 multiple theoretical solution for scaling), but the textual aspects (*i.e.* the address string itself). To scale
689 over dozens of million of addresses, specific architectures may be used to deal with the text search, for
690 instance distributed database (database sharding), in a similar spirit to the current software Elastic
691 Search. We stress however that given the current available amount of historical sources, such scaling
692 problem should not be an issue before a long time.

693 5. Conclusion

694 This article tackles the historical geocoding problem. As shown throughout the article, the
695 historical aspects bring major complications to the geocoding problem. The main difficulties come
696 from the nature of historical data (uncertainty, fuzzy date, precision, sparseness), which prevents
697 the use of current-address geocoding methods based on strong hierarchical modelling. Instead, we
698 propose a historical geocoding system based on a sound geohistorical object model. This model is
699 designed to cover the minimal features, and, by its genericity, modularity, and open source nature, can
700 easily be extended to feat other historical sources. Geohistorical objects from several historical sources
701 have been integrated into the database and coherently georeferenced and edited to form gazetteers.
702 Geocoding an address at a given time is then a matter of finding the best matching geohistorical
703 object in the gazetteers, if any. Our simple, coherent, historical geocoding system has been tested on
704 several real-life datasets collected by historians and can be easily used for other places/times/types of
705 localisations. Diverse historical sources covering two century for the city of Paris have been integrated
706 into the geocoder. The proposed geocoder is able to localise a large percentage of addresses with a fast
707 speed (about 200ms per address). Finally, the article describes a prototype of web-based User Interface
708 that demonstrates the interest of collaborative editing of localisation of addresses, and helps historians
709 and other digital humanities researchers use geocoding services.

710 **Supplementary Materials:** All the code and additional documentation are available on the project websites
711 <http://geohistoricaldata.org> and its associated code repository <https://github.com/Geohistoricaldata>. The code
712 for the geocoder itself it available here: https://github.com/GeoHistoricalData/historical_geocoding.

713 **Acknowledgments:** Thanks to the Belle Epoque project and Angelo Riva and Thierry Géraud, the Institut Louis
714 Bachelot for funding. Thanks to historians who contributed to creating the datasets, especially Benoit Costes for
715 the edit of Alphand map.

716

- 717 1. Goldberg, D.W.; Wilson, J.P.; Knoblock. From Text to Geographic Coordinates: The Current State of
718 Geocoding. *Journal of the Urban and Regional Information Systems Association* **2007**, *19*, 33–46.
- 719 2. St-Hilaire, M.; Moldofsky, B.; Richard, L.; Beaudry, M. Geocoding and Mapping Historical Census Data:
720 The Geographical Component of the Canadian Century Research Infrastructure. *Historical Methods: A
721 Journal of Quantitative and Interdisciplinary History* **2007**, *40*, 76–91. doi:10.3200/HMTS.40.2.76-91.
- 722 3. Daras, K.; Feng, Z.; Dibben, C. HAG-GIS: A spatial framework for geocoding historical addresses **2014**.
- 723 4. Hutchinson, M.J.; Veenendaal, B. An agent-based framework for intelligent geocoding. *Applied Geomatics*
724 **2013**, *5*, 33–44. doi:10.1007/s12518-011-0063-z.

- 725 5. Roongpiboonsopit, D.; Karimi, H.A. Comparative Evaluation and Analysis of Online Geocoding Services.
726 *International Journal of Geographical Information Science* **2010**, *24*, 1081–1100. doi:10.1080/13658810903289478.
- 727 6. Clough, P.; Tang, J.; Hall, M.M.; Warner, A. Linking archival data to location: A case study at the UK
728 national archives. *Aslib Proceedings* **2011**, *63*, 127–147. doi:10.1108/00012531111135628.
- 729 7. Mostern, R.; Johnson, I. From named place to naming event: creating gazetteers for
730 history. *International Journal of Geographical Information Science* **2008**, *22*, 1091–1108,
731 [<http://dx.doi.org/10.1080/13658810701851438>]. doi:10.1080/13658810701851438.
- 732 8. Smith, D.A.; Crane, G. Disambiguating Geographic Names in a Historical Digital Library. *International*
733 *Conference on Theory and Practice of Digital Libraries*. Springer, 2001, pp. 127–136.
- 734 9. Goodchild, M.F. Citizens as sensors: the world of volunteered geography. *GeoJournal* **2007**, *69*, 211–221.
735 doi:10.1007/s10708-007-9111-y.
- 736 10. Heipke, C. Crowdsourcing geospatial data. *{ISPRS} Journal of Photogrammetry and Remote Sensing* **2010**,
737 *65*, 550 – 557. {ISPRS} Centenary Celebration Issue, doi:<http://dx.doi.org/10.1016/j.isprsjprs.2010.06.005>.
- 738 11. Haklay, M., Citizen Science and Volunteered Geographic Information: Overview and Typology of
739 Participation. In *Crowdsourcing Geographic Knowledge: Volunteered Geographic Information (VGI) in Theory*
740 *and Practice*; Sui, D.; Elwood, S.; Goodchild, M., Eds.; Springer Netherlands: Dordrecht, 2013; pp. 105–122.
741 doi:10.1007/978-94-007-4587-2_7.
- 742 12. Fomel, S.; Claerbout, J.F. Guest Editors' Introduction: Reproducible Research. *Computing in Science*
743 *Engineering* **2009**, *11*, 5–7. doi:10.1109/MCSE.2009.14.
- 744 13. Aruliah, D.A.; Brown, C.T.; Hong, N.P.C.; Davis, M.; Guy, R.T.; Haddock, S.H.D.; Huff, K.; Mitchell, I.;
745 Plumbley, M.D.; Waugh, B.; White, E.P.; Wilson, G.; Wilson, P. Best Practices for Scientific Computing.
746 *CoRR* **2012**, *abs/1210.0530*.
- 747 14. Wilson, G.; Bryan, J.; Cranston, K.; Kitzes, J.; Nederbragt, L.; Teal, T.K. Good Enough Practices in Scientific
748 Computing. *ArXiv e-prints* **2016**, [[arXiv:cs.SE/1609.00037](https://arxiv.org/abs/1609.00037)].
- 749 15. Marwick, B. Computational Reproducibility in Archaeological Research: Basic Principles and a
750 Case Study of Their Implementation. *Journal of Archaeological Method and Theory* **2016**, pp. 1–27.
751 doi:10.1007/s10816-015-9272-9.
- 752 16. Duménieu, B. Un Système d'information Géographique Pour Le Suivi d'objets Historiques Urbains à
753 Travers l'espace et Le Temps. PhD thesis, Ecole des Hautes Etudes en Sciences Sociales, 2015.
- 754 17. Armstrong, M.P. Temporality in spatial databases. *GIS/LIS'88*, 1988, pp. 880–889.
- 755 18. Herrault, P.A.; Sheeren, D.; Fauvel, M.; Monteil, C.; Paegelow, M. A comparative study of geometric
756 transformation models for the historical "map of france" registration. *Geographia Technica* **2013**, pp. pp–34.
- 757 19. Fabbri, R.; Kimia, B. 3D curve sketch: Flexible curve-based stereo reconstruction and calibration. *Computer*
758 *Vision and Pattern Recognition (CVPR)*, 2010 IEEE Conference on. IEEE, 2010, pp. 1538–1545.
- 759 20. Cléri, I.; Pierrot-Deseilligny, M.; Vallet, B. Automatic Georeferencing of a Heritage of old analog aerial
760 Photographs. *ISPRS Annals of the Photogrammetry, Remote Sensing and Spatial Information Sciences* **2014**, *2*, 33.
- 761 21. Bitelli, G.; Cremonini, S.; Gatta, G. Ancient map comparisons and georeferencing techniques: a case study
762 from the Po River Delta (Italy). *E-perimetron* **2009**, *4*, 221–228.
- 763 22. Boutoura, C.; Livieratos, E. Some fundamentals for the study of the geometry of early maps by comparative
764 methods. *e-Perimetron* **2006**, *1*, 60–70.
- 765 23. De Runz, C.; Desjardin, E.; Piantoni, F.; Herbin, M. Anteriority index for managing fuzzy dates in
766 archaeological GIS. *Soft Computing* **2010**, *14*, 339.
- 767 24. Kauppinen, T.; Mantegari, G.; Paakkari, P.; Kuittinen, H.; Hyvönen, E.; Bandini, S. Determining
768 relevance of imprecise temporal intervals for cultural heritage information retrieval. *International journal of*
769 *human-computer studies* **2010**, *68*, 549–560.
- 770 25. Perret, J.; Gribaudo, M.; Barthelemy, M. Roads and cities of 18th century France. *Scientific Data* **2015**, *2*,
771 10.1038/sdata.2015.48.
- 772 26. Noizet, H.; Bove, B.; Costa, L. *Paris de Parcelles En Pixels*; Presses universitaires de Vincennes, 2013.
- 773 27. Gribaudo, M.; Magaud, J. L'action publique et ses administrateurs dans les domaines sanitaires et social en
774 France, 1800 à 1900, 1999.
- 775 28. Lazzara, G.; Levillain, R.; Géraud, T.; Jacquélet, Y.; Marquagnies, J.; Crépin-Leblond, A. The SCRIBO
776 Module of the Olena Platform: A Free Software Framework for Document Image Analysis. *Proceedings of*

- 777 the 11th International Conference on Document Analysis and Recognition (ICDAR); , 2011; pp. 252–258.
778 doi:10.1109/ICDAR.2011.59.
- 779 29. Plumejeaud-Perreau, C.; Grosso, E.; Parent, B. Dissemination and Geovisualization of Territorial Entities
780 History. *Journal of Spatial Information Science* **2014**, *8*, pp. 73–93. doi:10.5311/JOSIS.2014.8.119.
- 781 30. Shen, W.; Wang, J.; Han, J. Entity linking with a knowledge base: Issues, techniques, and solutions. *IEEE*
782 *Transactions on Knowledge and Data Engineering* **2015**, *27*, 443–460.
- 783 31. Overell, S. The problem of place name ambiguity. *SIGSPATIAL Special* **2011**, *3*, 12–15.
- 784 32. Mihalcea, R.; Csomai, A. Wikify!: Linking Documents to Encyclopedic Knowledge. Proceedings of the
785 Sixteenth ACM Conference on Conference on Information and Knowledge Management; ACM: New York,
786 NY, USA, 2007; CIKM '07, pp. 233–242. doi:10.1145/1321440.1321475.
- 787 33. Hachey, B.; Radford, W.; Nothman, J.; Honnibal, M.; Curran, J.R. Evaluating Entity Linking with Wikipedia.
788 *Artif. Intell.* **2013**, *194*, 130–150. doi:10.1016/j.artint.2012.04.005.
- 789 34. Zhang, W.; Gelernter, J. Geocoding location expressions in Twitter messages: A preference learning method.
790 *Journal of Spatial Information Science* **2014**, *2014*, 37–70.
- 791 35. Costes, B. Vers la construction d'un référentiel géographique ancien. Un modèle de graphe agrégé pour
792 intégrer, qualifier et analyser des réseaux géohistoriques. PhD thesis, Université Paris-Est, 2016.
- 793 36. Costes, B.; Perret, J.; Bucher, B.; Gribaudo, M. An aggregated graph to qualify historical spatial networks
794 using temporal patterns detection. 18th AGILE International Conference on Geographic Information
795 Science, 2015.
- 796 37. Zimmerman, D.L.; Fang, X.; Mazumdar, S.; Rushton, G. Modeling the Probability Distribution of Positional
797 Errors Incurred by Residential Address Geocoding. *International Journal of Health Geographics* **2007**, *6*, 1.
- 798 38. Brando, C.; Frontini, F.; Ganascia, J. REDEN: Named Entity Linking in Digital Literary Editions Using
799 Linked Data Sets. *CSIMQ* **2016**, *7*, 60–80. doi:10.7250/csimq.2016-7.04.
- 800 39. Zaveri, A.; Rula, A.; Maurino, A.; Pietrobon, R.; Lehmann, J.; Auer, S. Quality assessment for linked data:
801 A survey. *Semantic Web* **2016**, *7*, 63–93.
- 802 40. Brando, C.; Abadie, N.; Frontini, F. Linked Data Quality for Domain Specific Named Entity Linking.
803 Proceedings of the 1st Atelier Qualité des Données du Web, 16ème Conférence Internationale Francophone
804 sur l'Extraction et la Gestion de Connaissances, Reims, France, 2016.