

1 Article

2 Map Archive Mining: Visual-analytical Approaches 3 to Explore Large Historical Map Collections

4 Johannes H. Uhl ^{1,*}, Stefan Leyk ¹, Yao-Yi Chiang ², Weiwei Duan ² and Craig A. Knoblock ²

5 ¹ Department of Geography, University of Colorado Boulder, Boulder, Colorado, USA;

6 {johannes.uhl;stefan.leyk}@colorado.edu

7 ² Spatial Sciences Institute, University of Southern California, Los Angeles, California, USA;

8 {yaoyic;weiweidu;knoblock}@usc.edu

9 * Correspondence: johannes.uhl@colorado.edu; Tel.: +01-303-492-2631

10 **Abstract:** Historical maps constitute unique sources of retrospective geographic information.
11 Recently, several map archives containing map series covering large spatial and temporal extents
12 have been systematically scanned and made available to the public. The geographic information
13 contained in such data archives allows extending geospatial analysis retrospectively beyond the era
14 of digital cartography. However, given the large data volumes of such archives and the low
15 graphical quality of older map sheets, the processes to extract geographic information need to be
16 automated to the highest degree possible. In order to understand the salient characteristics, data
17 quality variation, and potential challenges in large-scale information extraction tasks, preparatory
18 analytical steps are required to efficiently assess spatio-temporal coverage, approximate map
19 content, and spatial accuracy of such georeferenced map archives across different cartographic
20 scales. Such preparatory steps are often neglected or ignored in the map processing literature but
21 represent highly critical phases that lay the foundation for any subsequent computational analysis
22 and recognition. In this contribution we demonstrate how such preparatory analyses can be
23 conducted using classical analytical and cartographic techniques as well as visual-analytical data
24 mining tools originating from machine learning and data science, exemplified for the United States
25 Geological Survey topographic map and Sanborn fire insurance map archives.

26 **Keywords:** map processing; retrospective landscape analysis; visual data mining, image retrieval,
27 low-level image descriptors, color moments, t-distributed stochastic neighborhood embedding,
28 USGS topographic maps, Sanborn fire insurance maps

29

30 1. Introduction

31 Historical maps contain valuable information about the Earth's surface in the past. This
32 information can provide a detailed understanding of the evolution of the landscape as well as the
33 interactions between and the drivers of geographic phenomena, such as human-made structures (e.g.,
34 transportation networks, settlements), vegetated land cover (e.g., forests, grasslands) or
35 hydrographic features (e.g., stream networks, water bodies). However, this spatial information is
36 typically locked in scanned map images and needs to be extracted to get access to the geographic
37 features of interest in machine readable data formats that can be imported into geospatial analysis
38 environments.

39 Map processing, or information extraction from digital map documents, is a branch of document
40 analysis that focuses on the development of methods for the extraction and recognition of information
41 in scanned cartographic documents. Map processing is an interdisciplinary field that combines
42 elements of computer vision, pattern recognition, geomatics, cartography, and machine learning. The
43 main goal of map processing is to "unlock" relevant information from scanned map documents to
44 provide this information in digital, machine-readable geospatial data formats as a means to preserve
45 the information digitally and facilitate the use of these data for analytical purposes [1].

46 Remotely sensed earth observation data from space and airborne sensors has been
47 systematically acquired since the early 1970s and provides abundant information for the monitoring
48 and assessment of geographic processes and how they interact over time. However, for the time
49 periods prior to operational remote sensing technology, there is little (digital) information that can
50 be used to document these processes. Thus, map processing often focuses on the development of
51 information extraction methods from map documents or engineering drawings created prior to the
52 era of remote sensing and digital cartography, thus expanding the temporal extent for carrying out
53 geographic analyses and landscape assessments to more than 100 years in many countries.

54 Information extraction from map documents includes the steps of *recognition* (i.e., identifying
55 objects in a scanned map such as groups of contiguous pixels with homogeneous semantic meaning),
56 and *extraction* i.e., transferring these objects into a machine-readable format (e.g., through
57 vectorization). Extraction processes typically involve image segmentation techniques based on
58 histogram analysis, color-space clustering, region growing or edge detection. Recognition in map
59 processing is typically conducted using template matching techniques involving shape descriptors,
60 cross-correlation measures or based on feature descriptors. Exemplary applications of map
61 processing techniques include the extraction of buildings [2-4], road networks [5], contour lines [6],
62 composite forest symbols [7] or the recognition of text from map documents [8,9]. Most approaches
63 rely on handcrafted or manually collected templates of the cartographic symbol of interest and
64 involve a significant level of user interaction, which impedes the application of such methods for
65 large-scale information extraction tasks where high degrees of automation are necessary to process
66 documents with high levels of variation in data quality.

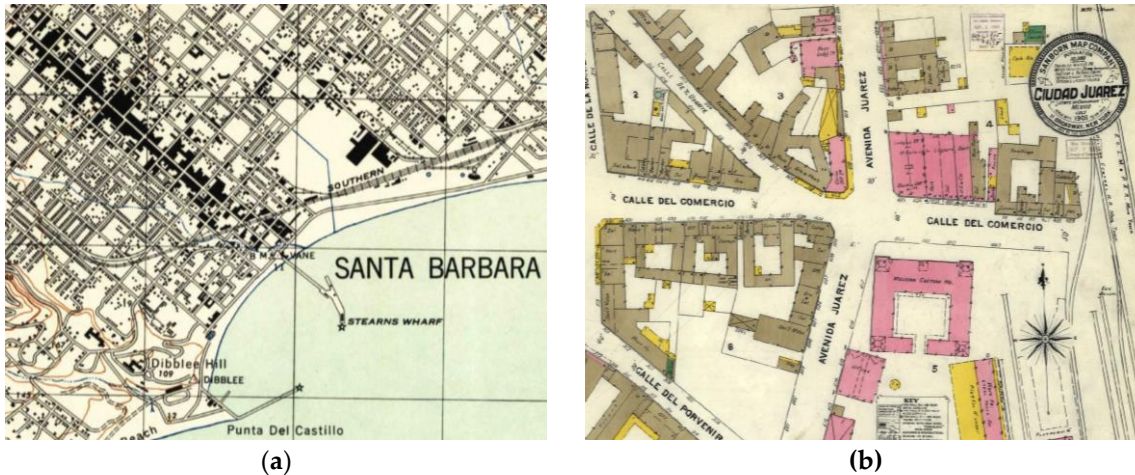
67 The availability of abundant contemporary geospatial data for many regions of the world offers
68 new opportunities to employ contemporary geospatial data as ancillary information to facilitate the
69 extraction and analysis of geographic content from historical map documents. Such approaches
70 include the use of contemporary spatial data for georeferencing historical maps [10], assessing the
71 presence of objects in historical maps across time at locations dictated by contemporary geospatial
72 vector data [11], or the automated collection of template graphics for cartographic symbols of interest
73 based on locations derived from modern geospatial data sources [12].

74 Most existing approaches for content extraction from historical maps still require a certain
75 degree of user interaction to ensure acceptable extraction performance for individual map sheets, e.g.
76 [13]. To overcome this persistent limitation, [14] and [15] propose the use of active learning and
77 similar interactive concepts for more efficient recognition of cartographic symbols in historical maps,
78 whereas [16] examine the usefulness of crowd-sourcing for the same purpose.

79 Moreover, the recent developments in deep machine learning in computer vision and image
80 recognition have catalyzed the use of such techniques for geospatial information extraction from
81 earth observation data [17-26]. This methodological development naturally projects into the idea of
82 applying state-of-the-art machine learning techniques for information extraction from scanned
83 cartographic documents, despite their fundamentally different nature compared to remotely sensed
84 data. Key in both cases is the need for abundant and representative training data which requires
85 automated sampling techniques. First attempts in this direction have used ancillary geospatial data
86 for the collection of large amounts of training data in historical maps [27-30].

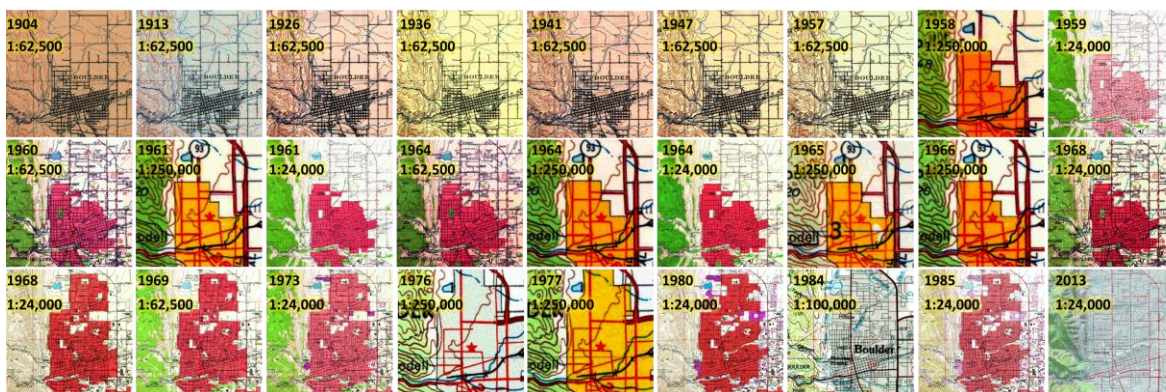
87 Besides this, several efforts have recently been conducted in different countries to systematically
88 scan, georeference, and publish whole series of topographic and other map documents. These
89 developments include efforts at the United States Geological Survey (USGS), that scanned and
90 georeferenced approx. 200,000 topographic maps published between 1884 and 2006 at different
91 cartographic scales between 1:24,000 and 1:250,000 [31] and the Sanborn fire insurance map collection
92 maintained by the U.S. Library of Congress, that contains more than 500,000 sheets of large-scale
93 maps of approximately 12,000 cities and towns in the U.S., Canada, Mexico, and Cuba, out of which
94 approximately 6,000 map sheets have been published as scanned map documents [32-34]. Figure 1
95 shows an example of a USGS topographic map sheet and a Sanborn map, respectively. Furthermore,
96 the United Kingdom Ordnance Survey scanned and georeferenced more than 200,000 topographic

97 map sheets and town plans for the United Kingdom dating back to the 1840s and provides many of
 98 them as seamless georeferenced raster layers [35,36].
 99



100 **Figure 1.** Examples of historical map documents: (a) USGS topographic map 1:31,680 from Santa
 101 Barbara (California, 1944) and (b) Sanborn fire insurance map from city center of Ciudad Juárez
 102 (Mexico, 1905).

103 These developments, alongside with advances in the processing, storage and distribution of
 104 large data volumes, offer great potential for automated, large-scale information extraction from
 105 historical cartographic document collections in order to preserve the contained geographic
 106 information and make it accessible for geospatial analysis. Because of the large amount of data
 107 contained in these map archives, high degrees of extraction automation are necessary. This
 108 constitutes a challenging task given the high variability in the content and quality of maps within an
 109 archive. Possible reasons for such variability are different conditions of the archived analogue map
 110 documents, differences in the scan quality, as well as changes in cartographic design best practices
 111 that may have resulted in different symbologies over multiple map editions (Figure 2).
 112



113 **Figure 2.** Example of the multi-temporal, multi-scale USGS topographic map archive, showing
 114 available map sheets covering Boulder, Colorado (USA) from 1904 to 2013 at various map scales.
 115

116 The challenges described above summarize some of the central tasks in map archive processing
 117 which include dealing with the sheer data volume, the differences in cartographic scales and designs,
 118 changes in content and cartographic representations and their degree of similarity in individual
 119 maps, the spatial and temporal coverage of the map sheets, and the spatial accuracy of the
 120 georeferenced map which dictates the degree of spatial agreement to contemporary geospatial
 121 ancillary data. While the previously described approaches represent promising directions towards
 122 higher levels of automation, they imply that the graphical characteristics of the map content to be
 123 extracted are known and that map scale and cartographic design remain approximately the same

124 across the processed map documents. Typically, many of these aspects are a priori unknown, since
125 such large amounts of data cannot be analyzed manually. However, these are relevant pieces of
126 information to better understand the data sources in order to design effective information extraction
127 methods.

128 The remote sensing community faces similar challenges. The steadily increasing amount of
129 remotely sensed earth observation data requires effective mining techniques to explore the content
130 of large remote sensing data archives. Therefore, visual data mining techniques have successfully
131 been used to comprehensively visualize the content of such archives. Such image information mining
132 (IIM) systems allow to discover and retrieve using available metadata, and based on the similarity of
133 the content of the individual datasets, or of patches of these [37-39] and guide exploratory analysis of
134 large amounts of data which facilitates the subsequent development of information extraction
135 methods. [40] implemented such a system for TerraSAR-X data, and [41] tested such approaches for
136 patches of Landsat ETM+ data and the UC Merced benchmark dataset. These systems are based on
137 spectral and textural descriptors precomputed at dataset or patch level that are then combined to
138 multidimensional descriptors characterizing spectral-textural content of the datasets or patches.
139 Other approaches include image segmentation methods to derive shape descriptors [42], include
140 spatial relationships between images into the IIM [43], or make use of structural descriptors to
141 characterize the change of geometric patterns over time across datasets within remote sensing data
142 archives [44]. Comparison of these descriptors facilitates the retrieval of similar content across large
143 archives. These approaches include methods for dimensionality reduction to visualize a whole
144 archive in a two or three-dimensional feature space based on content similarity.

145 Whereas in remote sensing data archives the spatio-temporal coverage of the data and their
146 quality is relatively well-known based on the sensor characteristics (e.g., time of operability,
147 satellite orbit, revisiting frequency, knowledge about physical parameters affecting data quality), this
148 may not always be the case for historical map archives, where metadata on spatial-temporal data
149 coverage might not be available or available in unstructured data formats only, preventing direct and
150 systematic analysis.

151 Thus, there is an urgent demand to develop a systematic approach to explore such digital map
152 archives, efficiently, prior to the extraction process, lending from similar efforts applied to remote
153 sensing data, but with a stronger focus on information obtained by metadata analysis. In this
154 contribution, we examine various techniques that could be used to build an image information
155 mining system for digital cartographic document archives in combination with metadata analysis.
156 These techniques allow to shed light on the following questions, which a potential user of such map
157 archives may ask prior to the design and implementation of information extraction methods:
158

- 159 • **What is the spatial and temporal coverage of the map archive content and does it vary across**
160 **different cartographic scales?** This is important because the coverage of the map data dictates
161 the spatial and temporal extent of the information that can be extracted from the map archive,
162 and thus, relates to the benefit of information extraction efforts and to the value of the extracted
163 data. Furthermore, such an analysis is useful to compare different map archives.
- 164 • **How accurate is the georeference of the maps contained in the archive? Does accuracy vary in**
165 **the spatio-temporal domain?** This constitutes a pressing question if ancillary geospatial data is
166 used for the information extraction, which requires certain degrees of spatial alignment between
167 map and ancillary data. For example, if it is possible to a-priori identify map sheets likely to
168 suffer from a high degree of positional inaccuracy, the user can exclude those map sheets from
169 template or training data collection, and thus, reduce the amount of noise in the collected
170 training data.
- 171 • **How much variability is there in the map content, regarding color, hue, contrast, and in the**
172 **cartographic styles used to represent the symbol of interest?** This is a central question affecting
173 the choice and design of a suitable recognition model. More powerful models or even separate
174 models for certain types of maps may be required if the representation of map content of interest
175 varies heavily across the map archive. Furthermore, knowledge about variations in map content

176 and about similarity between individual maps is useful regarding the training data sampling
 177 design, to ensure the collection of representative and balanced training samples.

178

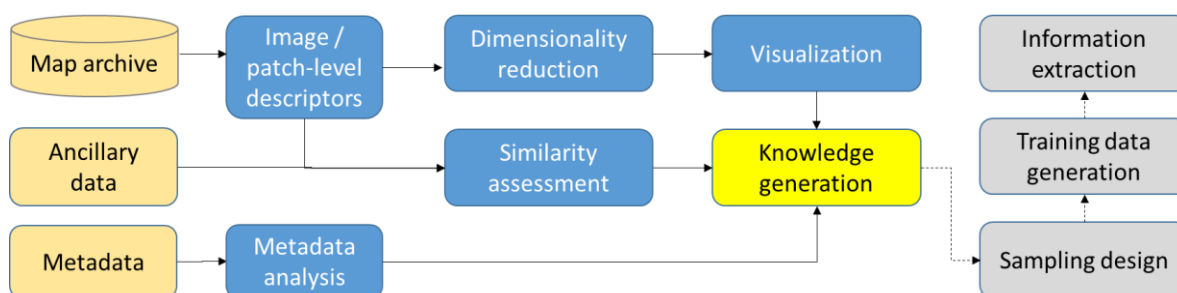
179 We present a set of methods that illuminate these questions, based on metadata analysis and
 180 descriptor-based visual data mining. Systematic mining approaches of relevant information about
 181 the map archive help to inform and educate the user community on critical aspects of data
 182 availability, quality and spatio-temporal coverage. Furthermore, these exploratory steps provide
 183 insights that are relevant for the implementation of large-scale information extraction methods from
 184 historical map archives and help to anticipate potential challenges involved. These methods have
 185 proven to provide valuable information highly relevant to design information extraction methods
 186 presented in [27,28], e.g., regarding the choice of training areas and classification methods. These
 187 methods can be generalized to other existing map archives in similar ways as well. Additionally, we
 188 aim to raise awareness about the importance of a-priori knowledge on large spatial data archives
 189 before using the data for information extraction purposes. Such a preprocessing stage is often
 190 neglected in published research that traditionally focuses on the extraction methods, specifically.
 191 However, this is important, non-trivial work highly relevant in the age of data intensive research on
 192 information extraction. We exemplify these methods using the USGS topographic map archive and
 193 the Sanborn fire insurance map collection.

194 2. Data, Methods and Results

195 In this contribution, we propose a set of methods that can be used to explore the spatial-temporal
 196 coverage of a historical map archive, its content, existing variations in cartographic design and to
 197 partially assess the spatial accuracy of the maps. The approaches range from pure metadata analysis
 198 to descriptor-based visual data mining techniques. *Metadata analysis* is conducted for the USGS
 199 topographic map archive exemplified for the states of California and Colorado (USA) based on
 200 structured metadata, as well as for the Sanborn fire insurance map archive in the United States based
 201 on unstructured metadata. *Content-based analysis* is demonstrated for the USGS topographic map
 202 archive covering the state of Colorado at different map scales, involving the use of image descriptors,
 203 dimensionality reduction, data visualization methods, and similarity assessment based on geospatial
 204 ancillary data. The USGS map archive includes 14,831 map documents in California, and 6,964 map
 205 sheets in Colorado,

206 Both metadata analysis and content-based analysis constitute preparatory steps yielding
 207 valuable information that facilitates the design and implementation of information extraction
 208 methods based on large map archives. Figure 3 shows how the proposed methods can be
 209 incorporated in information extraction workflows.

210



211

212 **Figure 3.** The methodology for metadata analysis and content-based knowledge generation on map
 213 archives to facilitate information extraction.

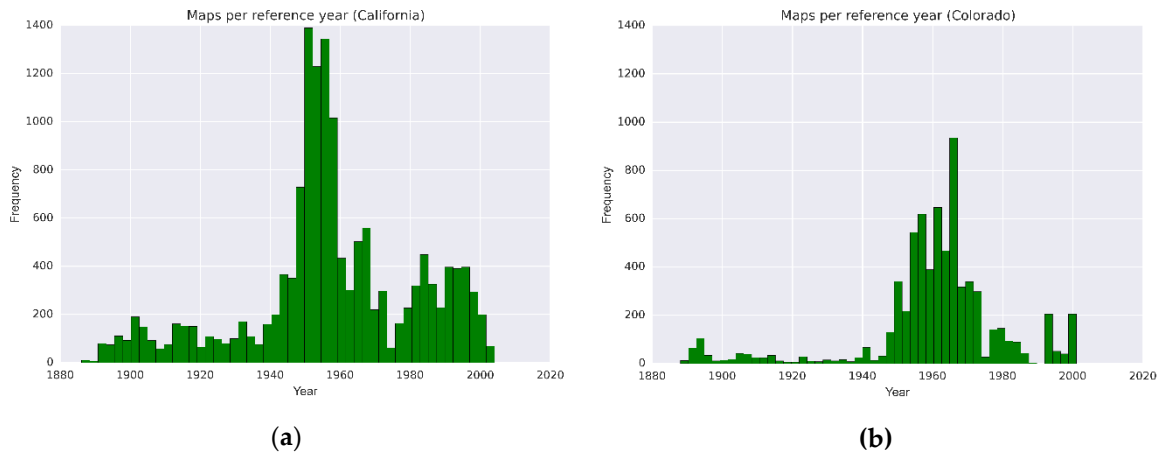
214

215

216 2.1. Metadata analysis

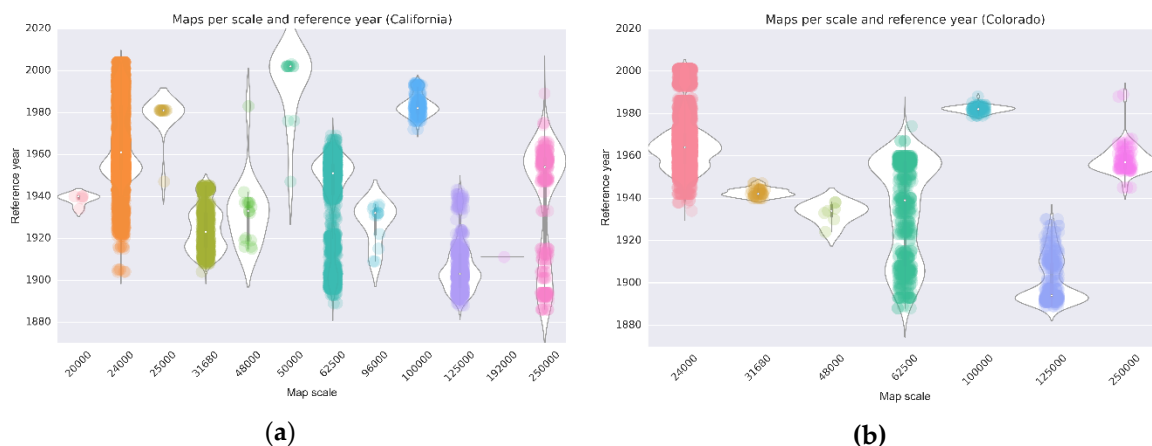
217 2.1.1. Metadata-based spatial-temporal coverage analysis

218 First, the temporal coverage of the map archives is analyzed. For the USGS map archive, which
 219 is accompanied by structured metadata (i.e., text files including unique identifiers for each map
 220 document), histograms based on the map reference year are created (Figure 4). It can be seen that the
 221 peak of map production was in California in the 1950s, and slightly later, in the 1960s in Colorado.



222 **Figure 4.** Histograms of USGS topographic maps (all available map scales) by reference year, (a) in
 223 California, and (b) in Colorado (USA).

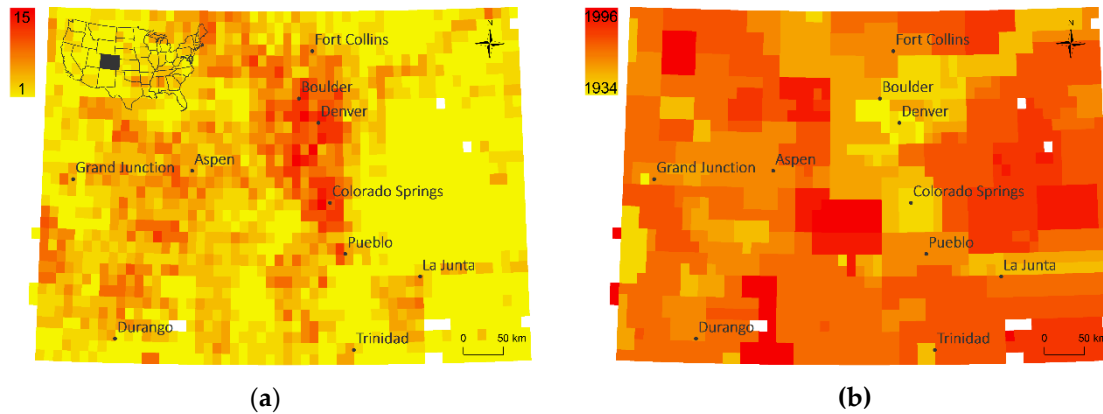
224 In addition to that, map production over time can be assessed in strata of map scales shown
 225 herein for the states of California and Colorado (Figure 5). These plots show the temporal distribution
 226 of published map editions (represented by the dots) and give an estimate of the underlying
 227 probability density (represented by the white areas) that indicates the map production intensity over
 228 time, separate and relative for each map scale. Such a representation helps to understand which time
 229 span can be covered with maps of various scales and thus can be used to determine which products
 230 to focus on for a particular purpose. This is important because maps of different scale contain
 231 different levels of detail resulting from cartographic generalization.
 232



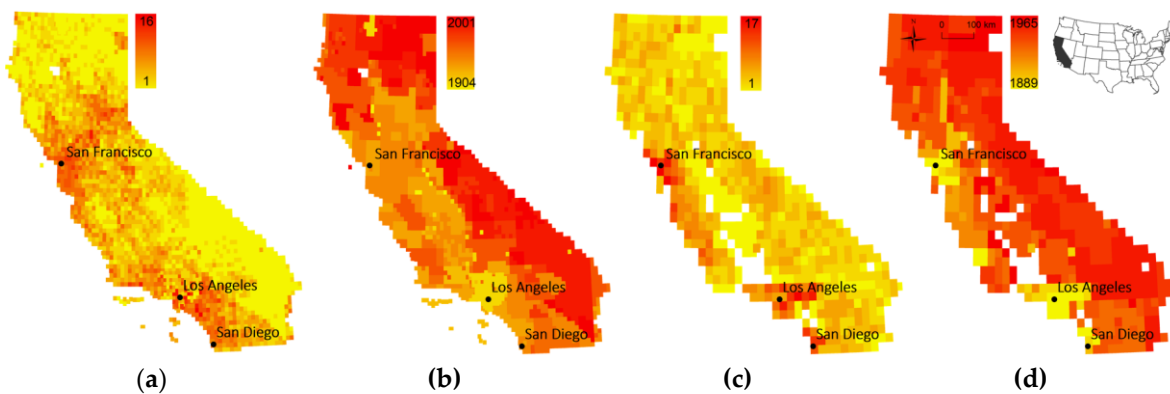
233 **Figure 5.** Produced USGS topographic maps per reference year and map scale (a) in California, and
 234 (b) in Colorado (USA).

235 In order to assess the spatial variability of map availability in a map archive over time, spatial map
 236 footprints (i.e., delineating map quadrangles) are generated based on USGS-delivered metadata. For
 237 each map footprint, statistics about available map sheets at those locations are computed. This allows
 238 the spatial visualization of the number of map editions and the earliest reference year available for
 239 each location, as shown in Figure 6 for the state of Colorado (scale 1:24,000), and for the map scales
 240 1:24,000 and 1:62,500 for the state of California in Figure 7, respectively. As can be seen, such

241 representations are useful to identify regions that have been mapped more intensively versus those
 242 for which temporal coverage is rather sparse. Furthermore, a user is immediately informed about the
 243 earliest map sheets for a location of interest to understand the maximum time period covered by
 244 these cartographic documents. Similar representations could be created for the average number of
 245 years between editions or the time span covered by map editions of a given map scale.



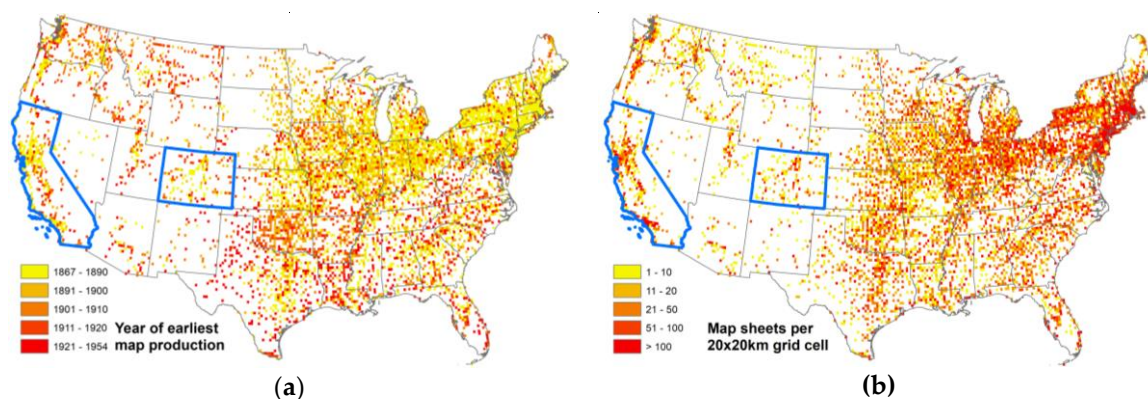
246 **Figure 6. (a)** Map edition counts and **(b)** earliest map production year per 1:24,000 map quadrangle
 247 in the state of Colorado (USA) based on metadata analysis.



248

249 **Figure 7. (a)** Map edition counts and **(b)** earliest map production year per 1:24,000 map quadrangle,
 250 **(c)** map edition counts and **(d)** earliest map production year per 1:62,500 map quadrangle in the state
 251 of California (USA) based on metadata analysis.

252 As a second example, the spatial-temporal coverage of the Sanborn fire insurance map archive
 253 is visualized. Since Sanborn map documents are commonly not offered as georeferenced datasets,
 254 this analysis is based on automatically extracted map locations (i.e., town or city name, county, and
 255 state) that were collected from unstructured metadata retrieved from HTML-based web content of
 256 the U.S. Library of Congress [45]. Additionally, the number of map sheets and their temporal
 257 coverage per location are extracted. The extracted data are geocoded using web-based geocoding
 258 services, which allows to visualize data availability and spatio-temporal coverage of Sanborn map
 259 documents. Figure 8 shows, similar to the above examples, the year of the first map production and
 260 the number of maps produced in total per location. The comparison of these visualizations for the
 261 highlighted states of California and Colorado to the previously shown Figures 6 and 7 shows the
 262 differences in spatio-temporal coverage between the two map archives, indicating a much sparser
 263 spatial coverage of the Sanborn map archive, but extending further back in time than the USGS map
 264 archive.
 265

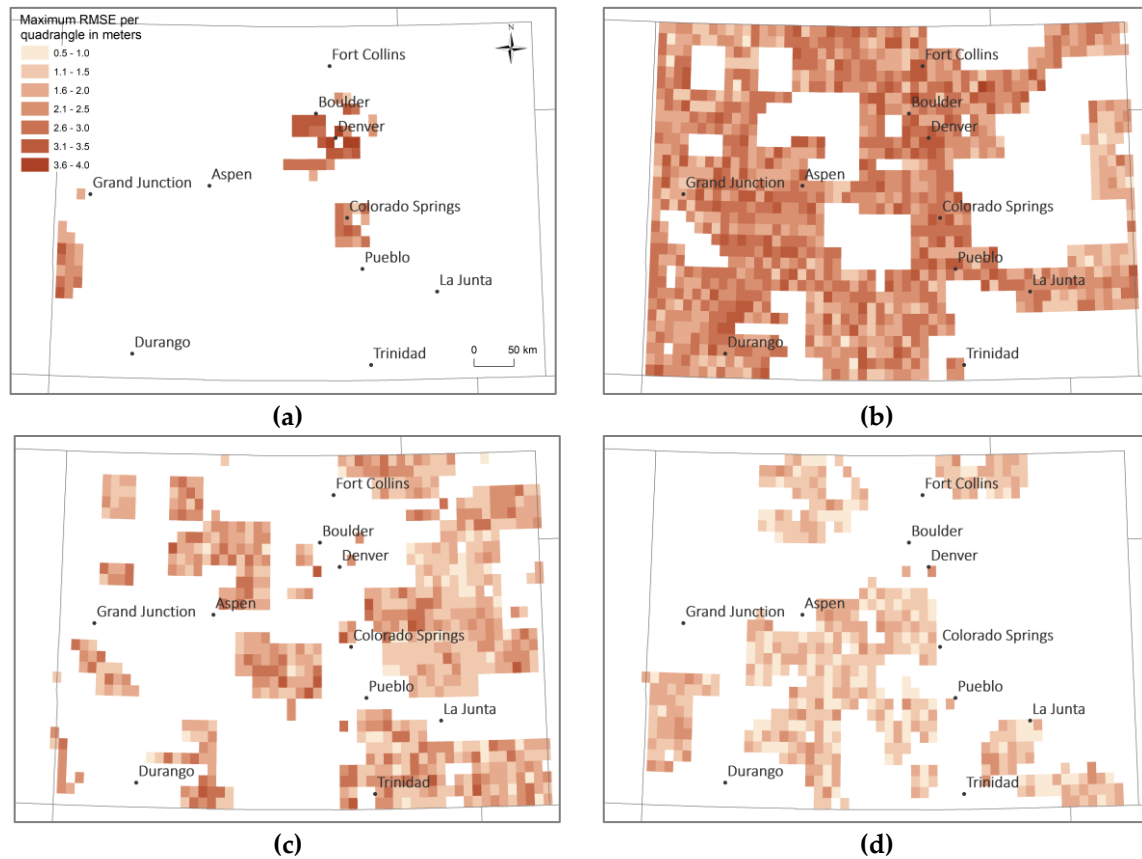


266 **Figure 8.** Sanborn fire insurance map archive coverage: (a) year of first map production per location
 267 and (b) number of available map sheets per location, both aggregated to grid cells of 10km for efficient
 268 visualization. Highlighted in blue the states of California and Colorado for comparison to the USGS
 269 map coverage shown in previous figures.

270 2.1.2. Metadata-based spatial-temporal analysis of positional accuracy

271 Positional accuracy of scanned maps can be caused by several factors, such as paper map
 272 distortions due to heat or humidity, the quality of surveying measurements on which the map
 273 production is based, deviations from the local geodetic datum at data acquisition time, cartographic
 274 generalization, and distortions introduced during the scanning and georeferencing process. While
 275 most of these effects cannot be reconstructed or quantified in detail, metadata delivered with the
 276 USGS topographic map archive contains information about the ground control points (GCPs) used
 277 for georeferencing the scanned map documents that allow for a partial assessment of these distortions
 278 and resulting positional inaccuracies.

279 The USGS topographic map quadrangle boundaries represent a graticule. For example, the
 280 corner coordinates for quadrangles of scale 1:24,000 are spaced in a regular grid of 7.5'x7.5'.
 281 Additionally, a finer graticule of 2.5'x2.5' is depicted in the maps. The intersections of this fine
 282 graticule are used by the USGS to georeference the maps. Therefore, pixel coordinates at those
 283 locations (i.e., the GCPs) are collected, and the corresponding known world coordinates of the
 284 graticule intersections are used to establish a second-order polynomial transformation based on least-
 285 squares adjustment. This transformation is used to warp the scanned document into a georeferenced
 286 raster dataset. The GCP coordinate pairs are reported in the metadata, as well as an error estimate
 287 per GCP that provides information on the georeference accuracy in pixels. Based on these error
 288 estimates given in pixel units and the spatial resolution of the georeferenced raster given in meters,
 289 the root mean standard error (RMSE) reflecting georeference accuracy in meters is calculated.
 290 Appending these RMSE values as attributes to the map quadrangle polygons allows to visualize
 291 georeferenced accuracy across the spatial-temporal domain. This is shown for the USGS maps of scale
 292 1:24,000 in the state of Colorado (Figure 9) for different time periods, visualizing the maximum RMSE
 293 per quadrangle and time period. Such temporally stratified representations are useful to examine if
 294 the georeference accuracy is constant over time. It can be seen that the earlier years in this example
 295 show higher degrees of inaccuracy than more recent map sheets. This has important implications for
 296 the user who is interested in using maps from different points in time that may exhibit different levels
 297 of inaccuracy.
 298

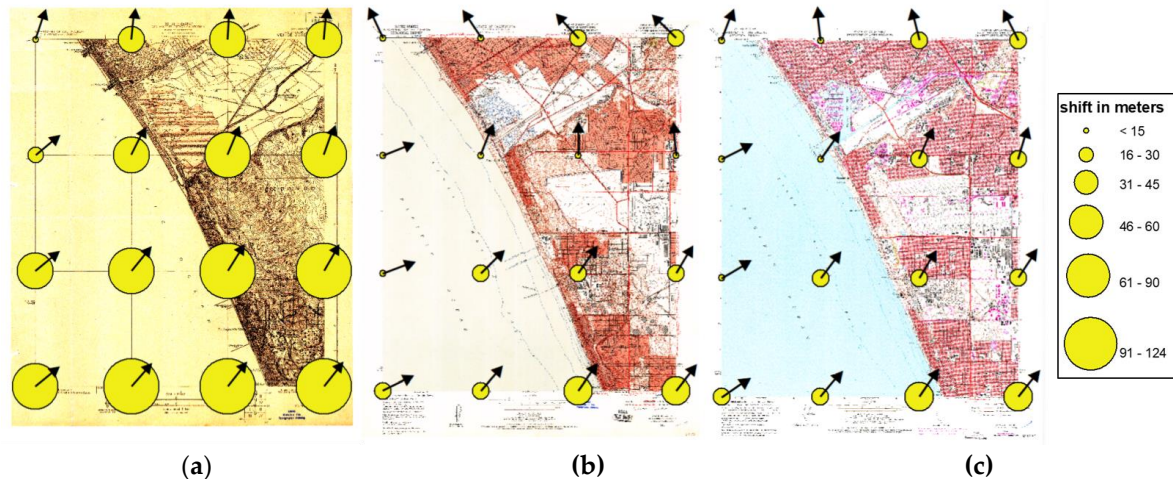


299 **Figure 9.** Spatio-temporal patterns of georeference accuracy of USGS topographic maps (1:24,000) in
 300 the state of Colorado (USA), for maps produced between (a) 1930-1950, (b) 1950-1970, (c) 1970-1990,
 301 and (d) 1990-2004.

302 Besides this information, the distortion introduced to the map by the warping process can be
 303 characterized by displacement vectors computed between the known world coordinates of each GCP
 304 (i.e., the graticule intersections) and the world coordinates corresponding to the respective pixel
 305 coordinates after applying the second-order polynomial transformation. These displacement vectors
 306 reflect geometric distortions and positional inaccuracy in the original map (i.e., *prior* to the
 307 georeferencing process) but are also affected by additional distortions introduced during
 308 georeferencing inaccuracies or through scanner miscalibration.

309 Assuming that objects in the map are affected by the same degree of inaccuracy like the graticule
 310 intersections, the magnitudes of these displacement vectors allow to estimate the maximum
 311 displacements to be expected between objects in the map and their real-world counterparts that may
 312 not be corrected by the second order polynomial transformation.

313 Figure 10 shows examples of these displacement vectors visualized for individual USGS map
 314 sheets at scale 1:24,000 from Venice (California) produced in 1923, 1957, and 1975. The magnitude of
 315 the local displacement is represented by the dot area, whereas the arrow indicates the displacement
 316 angle. This example shows similar patterns across the three maps, probably reflecting non-
 317 independent distortions between the maps since earlier maps are typically used as base maps for
 318 subsequent map editions, and some local variations due to inaccuracies introduced during
 319 georeferencing of the individual map sheets.



320

(a)

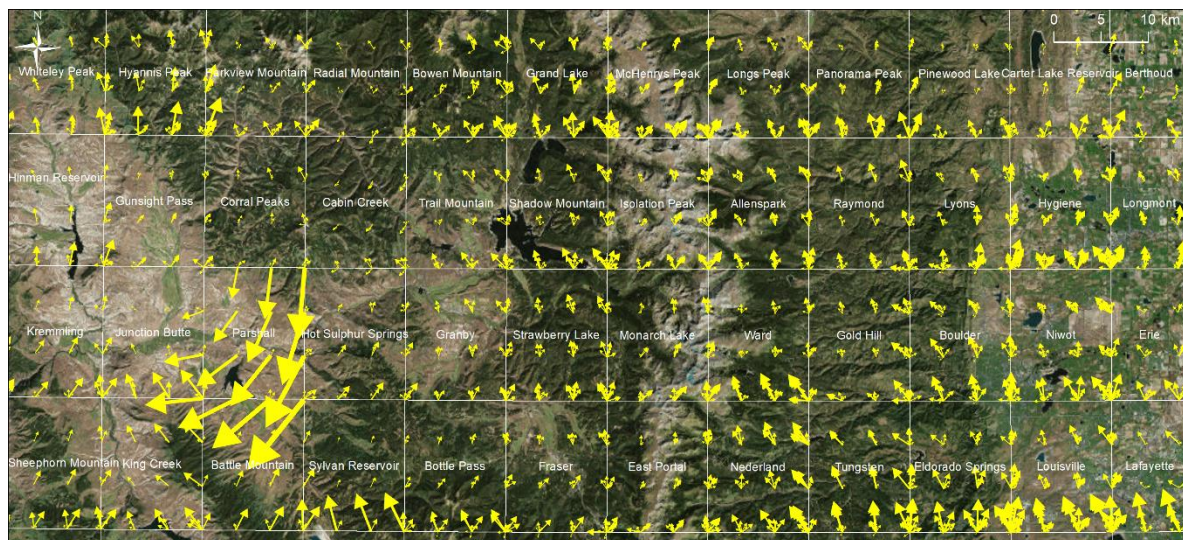
(b)

(c)

321 **Figure 10.** Displacement vectors at GCP locations characterizing the distortions introduced during
 322 the georeferencing of USGS topographic maps from Venice (California), produced in (a) 1923, (b) in
 323 1957, and (c) in 1975 (from left to right).

324 Additionally, these displacement vectors can be visualized as vector fields across large areas,
 325 allowing to identify regions, quadrangles, or individual maps of high or low positional reliability,
 326 respectively. Figure 11 shows the vector field of relative displacements for USGS maps of scale
 327 1:24,000 for a region Northwest of Denver, Colorado. Notable are the large displacement vectors in
 328 the Parshall quadrangle, indicating some anomalous map distortion, whereas the Cabin Creek
 329 quadrangle (Northeast of Parshall) seems to have suffered from very slight distortions only. Multiple
 330 arrows indicate the availability of multiple map editions in given quadrangles. Such visualizations
 331 may inform map users about the heterogeneity in distortions applied to the map sheets during the
 332 georeferencing process and may indicate different degrees of positional accuracy across a given study
 333 area.

334



335

336

337

338

Figure 11. Displacement vector field at GCP locations over multiple USGS map quadrangles of scale
 1:24,000, located North-west of Denver (Colorado), reflecting different types of distortions introduced
 to the map documents during the georeferencing process (Basemap source: [46]).

339 2.2. Content-based analysis

340

341

342

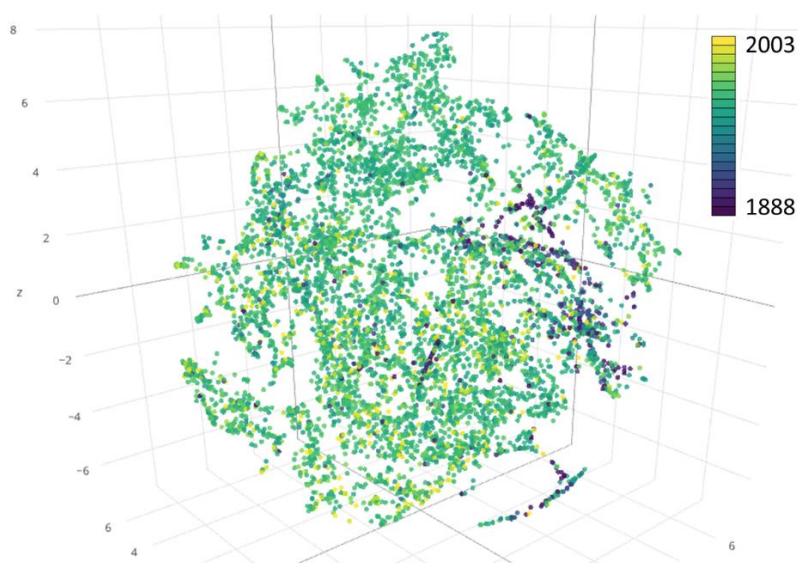
343

The presented metadata-based analysis provides valuable insights of spatial-temporal map
 availability, coverage, and spatial accuracy without analyzing the actual content of the map archives.
 However, it is important to inform the analyst about the degree of heterogeneity at content-level. In
 order to obtain detailed knowledge about the content of map archives, we propose a framework

344 based on low-level image descriptors computed for each map or map patches. Here, we employ color-
345 histogram based moments (i.e., mean, standard deviation, skewness and kurtosis, see [47]) computed
346 for each image channel in the RGB color space. Mean and standard deviation characterize hue,
347 brightness and contrast level of an image, skewness and kurtosis indicate the symmetry and flatness
348 of the probability density of the color distributions, and thus reflect color spread and variability of an
349 image. These four measures are computed for each channel of an image and stacked together to a 12-
350 dimensional feature descriptor, at image or patch level. In the case of scanned map documents, such
351 descriptors allow to retrieve maps or patches of maps of similar background color (depending on
352 paper type and scan contrast level), and maps of similar dominant map content, such as waterbodies,
353 urban areas, or forest cover. This similarity assessment is based on distances in the descriptor feature
354 space and can involve metadata (e.g., map reference year), or ancillary geospatial data, to assess map
355 content similarity across or within different geographic settings. Furthermore, approaches for
356 dimensionality reduction such as t-distributed stochastic neighborhood embedding (t-SNE, [48]) are
357 employed to visualize the data based on similarity in feature space. T-SNE allows to reduce the
358 dimensionality of high-dimensional data, where the relative distances between the data points in the
359 reduced feature space reflect the similarity of the data points in the original feature space. This
360 facilitates the visual or quantitative identification of clusters of similar map sheets and provides a
361 better understanding of the content of large map archives and their inherent variability. This kind of
362 similarity assessment and metadata analysis is useful in generating knowledge which can be used to
363 guide sampling designs to generate template or training data for supervised information extraction
364 techniques.

365 2.2.1. Content-based analysis at map level

366 Analyzing the content of the entire map archive with respect to similarities between the
367 individual map sheets is done by computing the image-moments based map descriptors. These 12-
368 dimensional features are transformed into a reduced 3-dimensional feature space that can be
369 visualized and interpreted intuitively. Figure 12 shows the 3D feature space for the 6,964 USGS maps
370 in the state of Colorado. The map reference year is used to color-code the points representing
371 individual map sheets. The clusters of dark blue points indicate fundamentally different color
372 characteristics of old maps in comparison to more recent maps represented by points colored in
373 green-yellow tones.

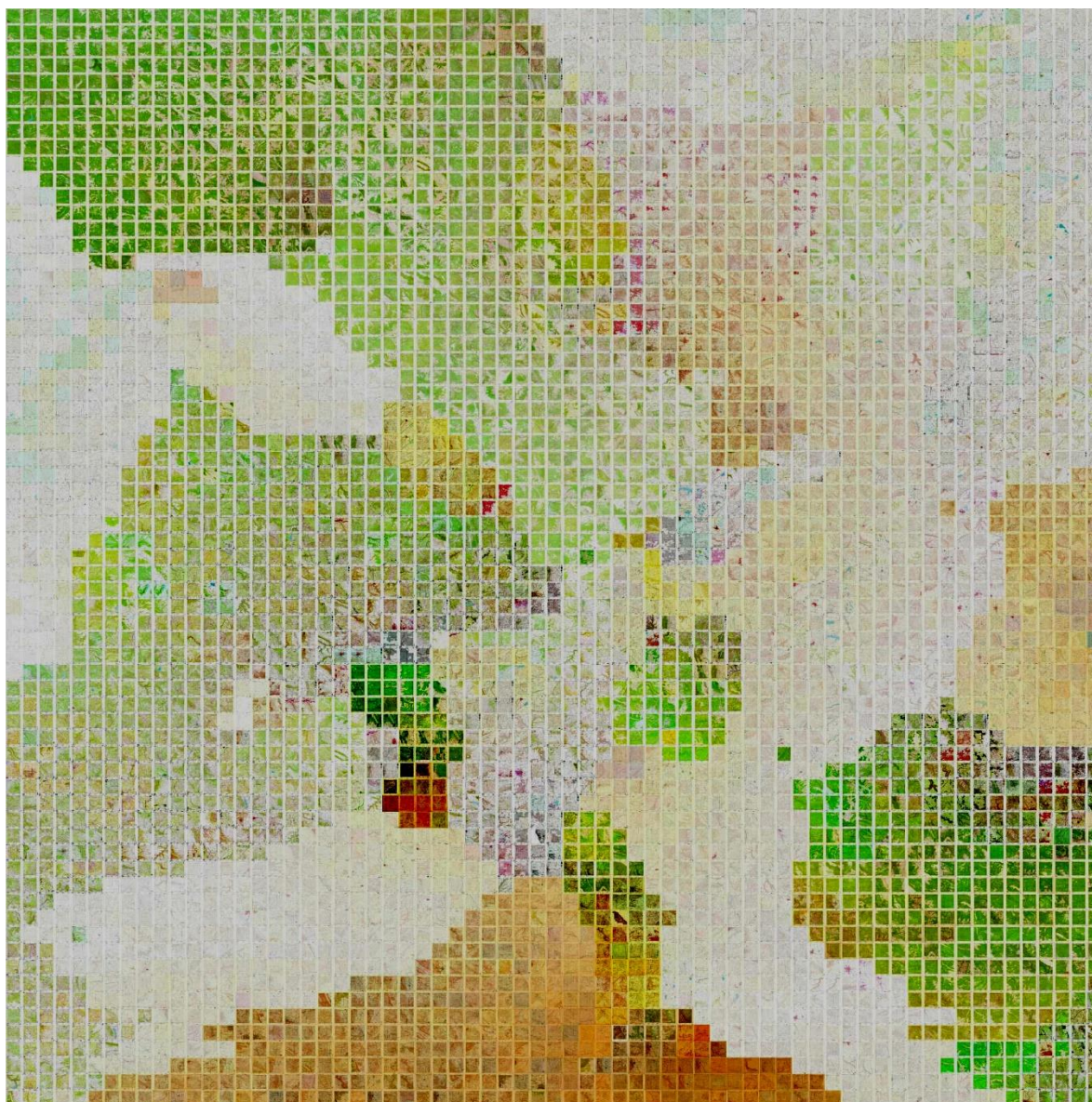


374

375 **Figure 12.** T-SNE visualization of the 6,964 USGS maps in the state of Colorado in a 3D feature space
376 based on 12-dimensional image descriptors obtained from channel-wise image moments.

377 In addition to color-coding the data points by the corresponding map reference year, the 12-
378 dimensional descriptors can be transformed into a 2D feature space and visualized using thumbnails

379 of individual maps corresponding to each data point in Figure 12. This transformation results in an
380 integrated visual assessment of map archives containing large numbers of map sheets. Figure 13
381 shows a t-SNE thumbnail visualization of a random sample (N=4,356) of the Colorado USGS maps in
382 a 2D feature space. Nearest neighbor snapping is used to create a rectified visualization. This is a very
383 effective way to visualize the variability in map contents, such as dominating forest area proportions.
384 It also illustrates the presence and abundance of different map designs and base color use, e.g., high
385 contrast and saturation levels in recent maps, compared to yellow-tinted map sheets from the
386 beginning of the 20th century centered at the bottom. The latter corresponds to the cluster of historical
387 maps located at the bottom of the point cloud in Figure 12.
388

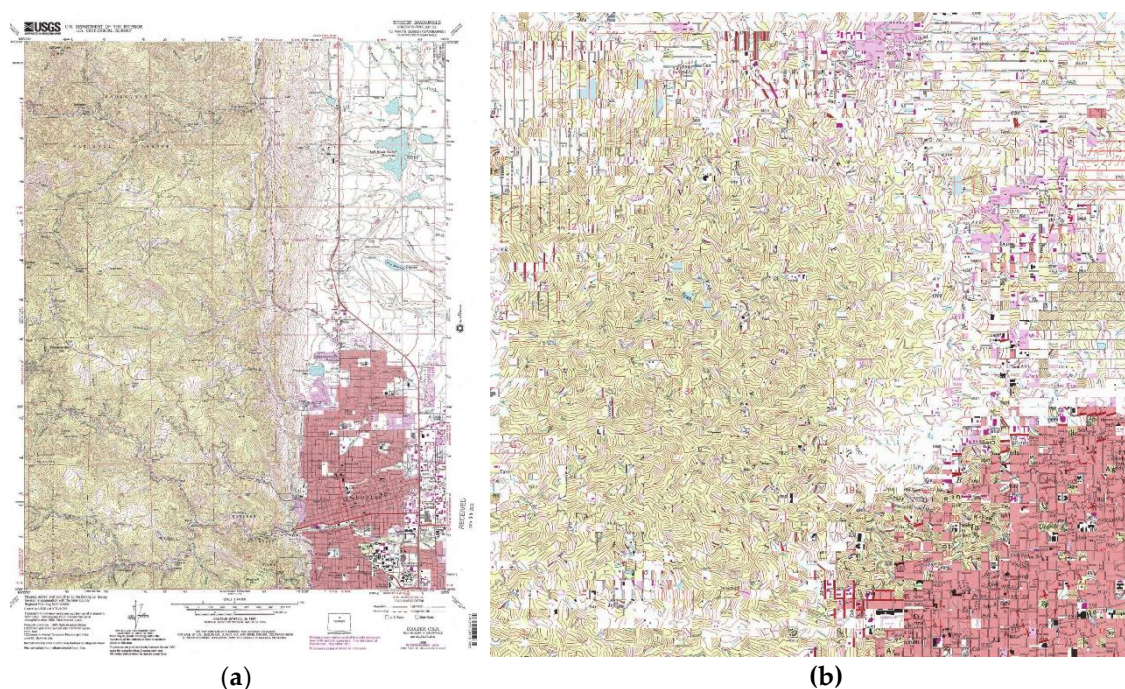


389
390 **Figure 13.** Thumbnail-based visualization of a subset of the USGS topographic maps in the state of
391 Colorado (USA) based on a 2D transformation of the 12-dimensional image descriptor feature space
392 using t-SNE.

393 2.2.2. Content-based analysis at within-map patch level

394 In order to assess the content within map sheets, map documents can be partitioned into tiles of
395 a fixed size (exemplified here for 100x100 pixels). Low-level descriptors based on image moments can
396 then be computed for each individual patch. However, if the patch size is chosen small enough, it is
397 computationally feasible to use the raw (or down-sampled) patch data (e.g., a line vector of all pixel

398 values in the patch) as a basis for t-SNE transformations. This can be useful if it is desired to introduce
399 a higher degree of spatiality and even directionality when assessing the similarity between the
400 patches. This method allows, for example, to rearrange a map document in patches based on patch
401 similarity, as shown for an example USGS map in Figure 14a. Quadrangle boundaries based on corner
402 coordinates delivered in the metadata can be used to clip the map contents and remove non-
403 geographic content in the map sheet edges. The clipped map content is then partitioned in tiles down-
404 sampled by factor 4, which results in a 1,875-dimensional feature vector per patch. These features are
405 then transformed into a 2D-space using t-SNE in order to create a similarity-based rearrangement of
406 the map patches (Figure 14b). This rearrangement highlights for example the groups of linear objects
407 of different dominant directions, such as road objects, or clusters of patches that contain contour lines
408 with diffuse directional characteristics. The incorporation of directionality may be useful to design
409 sampling schemes that generate training data allowing for rotation-invariant feature learning.
410



411 **Figure 14.** (a) USGS topographic map for Boulder, Colorado (1966), and (b) rearranged map patches
412 according to their similarity in a raw pixel value feature space using t-SNE.

413 2.2.3. Content-based analysis at cross-map patch level

414 If variations of specific cartographic symbols across large map archives are of interest and have to be
415 characterized, ancillary geospatial data can be employed to label the created map patches based on
416 their spatial relationships to the ancillary data. For example, it may be important to assess the
417 differences in cartographic representations of dense urban settlement areas across map sheets, in
418 order to design a recognition model for urban settlement. In such situations building footprint data
419 with built-year information and the respective spatio-temporal coverage can be employed to
420 reconstruct settlement distributions in a given map reference year (see [49]). Based on these reference
421 locations, building density surfaces can be computed for each map reference year. Using appropriate
422 thresholding allows to approximately delineate dense settlement areas for a given point in time.
423 Based on spatial overlap between map patches and these dense reference settlement areas, map
424 patches that are likely to contain urban area symbols can be identified across multiple maps. These
425 selected map patches can then be visualized in an integrated manner using t-SNE arrangement, as
426 exemplarily shown in Figure 15 for map patches collected across 50 USGS maps (1:24,000) in the states
427 of Colorado and California. This arrangement illustrates nicely the different cartographic styles that
428 are used to represent dense urban settlements across time and map sheets, and provides valuable
429 information useful for the design of a recognition model. Additional samples could be collected at

430 locations where no ancillary data is available, and their content can be estimated based on descriptor
431 similarity (i.e., patches of low Euclidean distance in the descriptor feature space) or using
432 unsupervised or supervised classification methods.



433
434 **Figure 15.** T-SNE arrangement of cross-map samples of patches likely to contain dense urban
435 settlement symbols.

436 3. Conclusions and Outlook

437 In this contribution we propose a set of methods for systematic information mining and content
438 retrieval in large collections of cartographic documents, such as topographic map archives. These
439 methods consist of pure metadata-based analyses, as well as content-based analyses using low-level
440 image descriptors such as histogram-based color moments, and dimensionality reduction methods
441 (i.e., t-SNE). We illustrate the proposed approach by exemplary analyses of the USGS topographic
442 map archive and the Sanborn fire insurance map collection. Our approach can be used to explore and
443 compare spatio-temporal coverage of these archives, the variability of positional accuracy, and
444 differences in content of the map documents based on visual-analytical tools. These content-based
445 map mining methods are inspired by image information mining systems implemented for remote
446 sensing data archives and have been applied to facilitate the design and implementation of
447 information extraction methods [27,28] Further work will include the identification of suitable
448 textural measures to be incorporated in the image descriptors. Additionally, the benefit of map
449 archive indexing based on low-level image descriptors will be tested in a prototypic map mining
450 framework. Moreover, these efforts will contribute to the design of adequate sampling methods to
451 generate representative training data for large-scale information extraction methods from historical
452 map archives.

453 **Acknowledgments:** This material is based on research sponsored in part by the National Science Foundation
454 under Grant Nos. IIS 1563933 (to the University of Colorado at Boulder) and IIS 1564164 (to the University of
455 Southern California).

456 **Author Contributions:** J.H.U. and S.L. conceived and designed the experiments; J.H.U. performed the
457 experiments; J.H.U. analyzed the data; J.H.U. wrote the paper.

458 **Conflicts of Interest:** The authors declare no conflict of interest.

459 **References**

- 460 1. Chiang, Y.-Y.; Leyk, S.; Knoblock, C.A. A survey of digital map processing techniques. *ACM Computing*
461 *Surveys* **2014**, *47*, 1-44. <http://dx.doi.org/10.1145/2557423>
- 462 2. Miyoshi, T.; Weiqing, L.; Kaneda, K.; Yamashita, H.; Nakamae, E. Automatic extraction of buildings
463 utilizing geometric features of a scanned topographic map. In *Proceedings of the 17th International Conference*
464 *on Pattern Recognition, 2004. ICPR 2004.*, IEEE: 2004. <http://dx.doi.org/10.1109/icpr.2004.1334607>
- 465 3. Laycock, S.D.; Brown, P.G.; Laycock, R.G.; Day, A.M. Aligning archive maps and extracting footprints for
466 analysis of historic urban environments. *Computers & Graphics* **2011**, *35*, 242-249.
467 <http://dx.doi.org/10.1016/j.cag.2011.01.002>
- 468 4. Arteaga, M.G. Historical map polygon and feature extractor. In *Proceedings of the 1st ACM SIGSPATIAL*
469 *International Workshop on MapInteraction - MapInteract '13*, ACM Press: 2013.
470 <http://dx.doi.org/10.1145/2534931.2534932>
- 471 5. Chiang, Y.-Y.; Leyk, S.; Knoblock, C.A. Efficient and robust graphics recognition from historical maps. In
472 *Graphics Recognition. New Trends and Challenges*, Springer Berlin Heidelberg: 2013; pp 25-35.
473 http://dx.doi.org/10.1007/978-3-642-36824-0_3
- 474 6. Miao, Q.; Liu, T.; Song, J.; Gong, M.; Yang, Y. Guided superpixel method for topographic map processing.
475 *IEEE Transactions on Geoscience and Remote Sensing* **2016**, *54*, 6265-6279.
476 <http://dx.doi.org/10.1109/tgrs.2016.2567481>
- 477 7. Leyk, S.; Boesch, R. Extracting composite cartographic area features in low-quality maps. *Cartography and*
478 *Geographic Information Science* **2009**, *36*, 71-79. <http://dx.doi.org/10.1559/152304009787340115>
- 479 8. Chiang, Y.-Y.; Moghaddam, S.; Gupta, S.; Fernandes, R.; Knoblock, C.A. From map images to geographic
480 names. In *Proceedings of the 22nd ACM SIGSPATIAL International Conference on Advances in Geographic*
481 *Information Systems - SIGSPATIAL '14*, ACM Press: 2014. <http://dx.doi.org/10.1145/2666310.2666374>
- 482 9. Chiang, Y.-Y.; Leyk, S.; Honarvar Nazari, N.; Moghaddam, S.; Tan, T.X. Assessing the impact of graphical
483 quality on automatic text recognition in digital maps. *Computers & Geosciences* **2016**, *93*, 21-35.
484 <http://dx.doi.org/10.1016/j.cageo.2016.04.013>
- 485 10. Tsorlini, A.; Iosifescu, I.; Iosifescu, C.; Hurni, L. A methodological framework for analyzing digitally
486 historical maps using data from different sources through an online interactive platform. *e-Perimetron*, **2014**,
487 *9(4)*, 153-165.
- 488 11. Hurni, L.; Lorenz, C.; Oleggini, L. Cartographic reconstruction of historic settlement development by
489 means of modern geo-data. *Proceedings of the 26th International cartographic conference*. Dresden,
490 Germany, 2013.
- 491 12. Leyk, S.; and Chiang, Y. Information extraction of hydrographic features from historical map archives using
492 the concept of geographic context. *Proceedings of AutoCarto 2016*, Albuquerque, New Mexico, USA, 2016.
- 493 13. Iosifescu, I.; Tsorlini, A.; Hurni, L. Towards a comprehensive methodology for automatic vectorization of
494 raster historical maps. *e-Perimetron* **2016**, *11(2)*, 57-76.
- 495 14. Budig, B.; van Dijk, T.C. Active learning for classifying template matches in historical maps. In *Discovery*
496 *Science*, Springer International Publishing: 2015; pp 33-47. http://dx.doi.org/10.1007/978-3-319-24282-8_5
- 497 15. Budig, B.; Dijk, T.C.V.; Wolff, A. Matching labels and markers in historical maps. *ACM Transactions on*
498 *Spatial Algorithms and Systems* **2016**, *2*, 1-24. <http://dx.doi.org/10.1145/2994598>
- 499 16. Budig, B.; van Dijk, T.C.; Feitsch, F.; Arteaga, M.G. Polygon consensus. In *Proceedings of the 24th ACM*
500 *SIGSPATIAL International Conference on Advances in Geographic Information Systems - GIS '16*, ACM
501 Press: 2016. <http://dx.doi.org/10.1145/2996913.2996951>
- 502 17. Maire, F.; Mejias, L.; Hodgson, A. A convolutional neural network for automatic analysis of aerial imagery.
503 In *2014 International Conference on Digital Image Computing: Techniques and Applications (DICTA)*, IEEE: 2014.
504 <http://dx.doi.org/10.1109/dicta.2014.7008084>
- 505 18. Castelluccio, M.; Poggi, G.; Sansone, C.; Verdoliva, L. Training convolutional neural networks for semantic
506 classification of remote sensing imagery. In *2017 Joint Urban Remote Sensing Event (JURSE)*, IEEE: 2017.
507 <http://dx.doi.org/10.1109/jurse.2017.7924535>
- 508 19. Audebert, N.; Le Saux, B.; Lefèvre, S. Semantic segmentation of earth observation data using multimodal
509 and multi-scale deep networks. In *Computer Vision – ACCV 2016*, Springer International Publishing: 2017;
510 pp 180-196. http://dx.doi.org/10.1007/978-3-319-54181-5_12

- 511 20. Marmanis, D.; Datcu, M.; Esch, T.; Stilla, U. Deep learning earth observation classification using imagenet
512 pretrained networks. *IEEE Geoscience and Remote Sensing Letters* **2016**, *13*, 105-109.
513 <http://dx.doi.org/10.1109/lgrs.2015.2499239>
- 514 21. Romero, A.; Gatta, C.; Camps-Valls, G. Unsupervised deep feature extraction for remote sensing image
515 classification. *IEEE Transactions on Geoscience and Remote Sensing* **2016**, *54*, 1349-1362.
516 <http://dx.doi.org/10.1109/tgrs.2015.2478379>
- 517 22. Scott, G.J.; England, M.R.; Starns, W.A.; Marcum, R.A.; Davis, C.H. Training deep convolutional neural
518 networks for land-cover classification of high-resolution imagery. *IEEE Geoscience and Remote Sensing*
519 *Letters* **2017**, *14*, 549-553. <http://dx.doi.org/10.1109/lgrs.2017.2657778>
- 520 23. Zhao, W.; Jiao, L.; Ma, W.; Zhao, J.; Zhao, J.; Liu, H.; Cao, X.; Yang, S. Superpixel-based multiple local cnn
521 for panchromatic and multispectral image classification. *IEEE Transactions on Geoscience and Remote Sensing*
522 **2017**, *55*, 4141-4156. <http://dx.doi.org/10.1109/tgrs.2017.2689018>
- 523 24. Zhang, L.; Zhang, L.; Du, B. Deep learning for remote sensing data: A technical tutorial on the state of the
524 art. *IEEE Geoscience and Remote Sensing Magazine* **2016**, *4*, 22-40. <http://dx.doi.org/10.1109/mgrs.2016.2540798>
- 525 25. Zhu, X.X.; Tuia, D.; Mou, L.; Xia, G.-S.; Zhang, L.; Xu, F.; Fraundorfer, F. Deep learning in remote sensing:
526 A comprehensive review and list of resources. *IEEE Geoscience and Remote Sensing Magazine* **2017**, *5*, 8-36.
527 <http://dx.doi.org/10.1109/mgrs.2017.2762307>
- 528 26. Ball, J.E.; Anderson, D.T.; Chan, C.S. Comprehensive survey of deep learning in remote sensing: Theories,
529 tools, and challenges for the community. *Journal of Applied Remote Sensing* **2017**, *11*, 1.
530 <http://dx.doi.org/10.1117/1.jrs.11.042609>
- 531 27. Uhl, J.H.; Leyk, S.; Yao-Yi, C.; Weiwei, D.; Knoblock, C.A. Extracting human settlement footprint from
532 historical topographic map series using context-based machine learning. In *8th International Conference of*
533 *Pattern Recognition Systems (ICPRS 2017)*, Institution of Engineering and Technology: 2017.
534 <http://dx.doi.org/10.1049/cp.2017.0144>
- 535 28. Uhl, J.H.; Leyk, S.; Yao-Yi, C.; Weiwei, D.; Knoblock, C.A. Spatializing uncertainty in image segmentation
536 using weakly supervised convolutional neural networks: a case study from historical map processing
537 (under review)
- 538 29. Duan, W.; Chiang, Y.-Y.; Knoblock, C.A.; Jain, V.; Feldman, D.; Uhl, J.H.; Leyk, S. Automatic alignment of
539 geographic features in contemporary vector data and historical maps. In *Proceedings of the 1st Workshop on*
540 *Artificial Intelligence and Deep Learning for Geographic Knowledge Discovery - GeoAI '17*, ACM Press: 2017.
541 <http://dx.doi.org/10.1145/3149808.3149816>
- 542 30. Duan, W.; Chiang, Y.-Y.; Knoblock, C.A.; Uhl, J.H.; Leyk, S. Automatic generation of precisely delineated
543 geographic features from georeferenced historical maps using deep learning (under review)
- 544 31. Fishburn, K.A.; Davis, L.R.; Allord, G.J. Scanning and georeferencing historical usgs quadrangles. In *Fact*
545 *Sheet*, US Geological Survey: 2017. <http://dx.doi.org/10.3133/fs20173048>
- 546 32. U.S. Library of Congress. Available online: <http://www.loc.gov/rr/geogmap/sanborn/san6.html> (accessed
547 on 28/02/2018).
- 548 33. U.S. Library of Congress. Available online: <http://www.loc.gov/rr/geogmap/sanborn/> (accessed on
549 28/02/2018).
- 550 34. U.S. Library of Congress. Available online: [https://www.loc.gov/item/prn-17-074/sanborn-fire-insurance-](https://www.loc.gov/item/prn-17-074/sanborn-fire-insurance-maps-now-online/2017-05-25/)
551 [maps-now-online/2017-05-25/](https://www.loc.gov/item/prn-17-074/sanborn-fire-insurance-maps-now-online/2017-05-25/) accessed on 28/02/2018).
- 552 35. National Library of Scotland. Available online: <https://maps.nls.uk/os/index.html> (accessed on 28/02/2018).
- 553 36. National Library of Scotland. Available online: <http://maps.nls.uk/geo/explore> (accessed on 28/02/2018).
- 554 37. Datcu, M.; Daschiel, H.; Pelizzari, A.; Quartulli, M.; Galoppo, A.; Colapicchioni, A.; Pastori, M.; Seidel, K.;
555 Marchetti, P.G.; D'Elia, S. Information mining in remote sensing image archives: System concepts. *IEEE*
556 *Transactions on Geoscience and Remote Sensing* **2003**, *41*, 2923-2936. <http://dx.doi.org/10.1109/tgrs.2003.817197>
- 557 38. Quartulli, M.; G. Olaizola, I. A review of eo image information mining. *ISPRS Journal of Photogrammetry and*
558 *Remote Sensing* **2013**, *75*, 11-28. <http://dx.doi.org/10.1016/j.isprsjprs.2012.09.010>
- 559 39. Espinoza-Molina, D.; Alonso, K.; Datcu, M. Visual data mining for feature space exploration using in-situ
560 data. In *2016 IEEE International Geoscience and Remote Sensing Symposium (IGARSS)*, IEEE: 2016.
561 <http://dx.doi.org/10.1109/igarss.2016.7730543>
- 562 40. Espinoza Molina, D.; Datcu, M. Data mining and knowledge discovery tools for exploiting big earth-
563 observation data. *ISPRS - International Archives of the Photogrammetry, Remote Sensing and Spatial Information*
564 *Sciences* **2015**, *XL-7/W3*, 627-633. <http://dx.doi.org/10.5194/isprsarchives-xl-7-w3-627-2015>

- 565 41. Griparis, A.; Faur, D.; Datcu, M. Dimensionality reduction for visual data mining of earth observation
566 archives. *IEEE Geoscience and Remote Sensing Letters* **2016**, *13*, 1701-1705.
567 <http://dx.doi.org/10.1109/lgrs.2016.2604919>
- 568 42. Durbha, S.S.; King, R.L. Semantics-enabled framework for knowledge discovery from earth observation
569 data archives. *IEEE Transactions on Geoscience and Remote Sensing* **2005**, *43*, 2563-2572.
570 <http://dx.doi.org/10.1109/tgrs.2005.847908>
- 571 43. Kurte, K.R.; Durbha, S.S.; King, R.L.; Younan, N.H.; Vatsavai, R. Semantics-enabled framework for spatial
572 image information mining of linked earth observation data. *IEEE Journal of Selected Topics in Applied Earth
573 Observations and Remote Sensing* **2017**, *10*, 29-44. <http://dx.doi.org/10.1109/jstars.2016.2547992>
- 574 44. Silva, M.P.S.; Camara, G.; Souza, R.C.M.; Valeriano, D.M.; Escada, M.I.S. Mining patterns of change in
575 remote sensing image databases. In *Fifth IEEE International Conference on Data Mining (ICDM'05)*, IEEE.
576 <http://dx.doi.org/10.1109/icdm.2005.98>
- 577 45. Library of Congress. Available online: <http://www.loc.gov/rr/geogmap/sanborn/country.php?countryID=1>
578 (accessed on 28/02/2018).
- 579 46. ESRI Basemaps: Esri, DigitalGlobe, GeoEye, Earthstar Geographics, CNES/Airbus DS, USDA, USGS, AEX,
580 Getmapping, Aerogrid, IGN, IGP, swisstopo, and the GIS User Community
- 581 47. Huang, Z.-C.; Chan, P.P.K.; Ng, W.W.Y.; Yeung, D.S. Content-based image retrieval using color moment
582 and gabor texture feature. In *2010 International Conference on Machine Learning and Cybernetics*, IEEE: 2010.
583 <http://dx.doi.org/10.1109/icmlc.2010.5580566>
- 584 48. Van der Maaten, L; Hinton, G. Visualizing data using t-SNE. *Journal of Machine Learning Research*, **2008**, *9*,
585 2579-2605
- 586 49. Leyk, S.; Uhl, J.H.; Balk, D.; Jones, B. Assessing the accuracy of multi-temporal built-up land layers across
587 rural-urban trajectories in the united states. *Remote Sensing of Environment* **2018**, *204*, 898-917.
588 <http://dx.doi.org/10.1016/j.rse.2017.08.035>