OXFORD

# Effective annotation for the automatic vectorization of cadastral maps

**Remi Petitpierre** [1]* , **Paul Guhennec**[2]

[1]Institute for Area and Global Studies, EPFL, Lausanne, Switzerland
[2]Digital Humanities Laboratory, EPFL, Lausanne, Switzerland
*Correspondence: Remi Petitpierre, Institute for Area and Global Studies, CDH-EPFL, CM 1 468, Station 10, 1015 Lausanne, Switzerland.
E-mail: remi.petitpierre@epfl.ch

## Abstract

The great potential brought by large-scale data in the humanities is still hindered by the time and technicality required for making documents digitally intelligible. Within urban studies, historical cadasters have been hitherto largely under-explored despite their informative value. Powerful and generic technologies, based on neural networks, to automate the vectorization of historical maps have recently become available. However, the transfer of these technologies is hampered by the scarcity of interdisciplinary exchanges and a lack of practical literature destined to humanities scholars, especially on the key step of the pipeline: the annotation. In this article, we propose a set of practical recommendations based on empirical findings on document annotation and automatic vectorization, focusing on the example case of historical cadasters. Our recommendations are generic and easily applicable, based on a solid experience on concrete and diverse projects.

## 1 Introduction

Cadastral maps are key historical sources. Both extremely detailed and very rich, they are the main witness of the evolution of the territory since the eighteenth century, and even before in some countries (Kain and Baigent, 1992; Dolej and Forejt, 2019). In theory, a systematic treatment of cadastral sources could yield the creation of large geohistorical databases and open the door to new paths of scientific analysis (Domaas *et al.*, 2003; Cousins *et al.*, 2007; Ekamper, 2010; Mou, 2012). The present challenges include the disenclavement of the city in favour of an extended analysis involving the surrounding territories and neighbouring cities, and in general the inclusion of the object studied in a macroscopic context, or a comparative angle. The shift from an idiographic analysis to a nomothetic perspective necessarily involves a process of abstraction of the complexity of the city and its development.
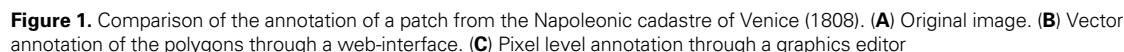
One of the most common abstractions to achieve that is the process of vectorizing data (Costes, 2016; Lelo, 2020; di Lenardo *et al.*, 2021). Vectorization is an inherently simplified representation of both the map and the geographic space, allowing to visualize and compare them. However, the process of georeferencing and then manually vectorizing the map manually is slow and expensive (Candeias *et al.*, 2003; Valent *et al.*, 2016). The manual vectorization of a city map can typically take several weeks (Costes, 2016). As a result, the size of the data suffers and the research tends to be limited to single study cases, and strictly restricted to the city boundaries, at the expense of comparative approaches.

One solution to this problem is to automate the vectorization process. Automatic vectorization of cadastral maps has triggered interest early on (Chen *et al.*, 1996; Katona and Hudra, 1999; Katona, 2000; Viglino and Pierrot-Deseilligny, 2003; Balkoca *et al.*, 2011). Several attempts have been presented in the last 50 years focusing on the automatic vectorization of topographic maps, using computer vision techniques for colour extraction or patterns recognition (Cofer and Tou, 1972; Frischknecht and Kanani, 1998; Dhar and Chanda, 2006; Pradhan *et al.*, 2010). However, the development of flexible and generic algorithms for cartography was a technical challenge until the late 2010s (Chiang *et al.*, 2014; Ignjatić *et al.*, 2018). The emergence of learning algorithms and deep neural networks allowed the development of generic vectorization models that can then specifically adapt to each case of study

by learning from a few human-annotated samples (Oliveira *et al.*, 2018; Digital Humanities Laboratory, 2021). A first attempt was presented for the Napoleonic cadastre of Venice in 2019 (Oliveira *et al.*, 2019). Similar projects have also been carried out on related problems such as the segmentation of floor plans or topographic maps (Liu *et al.*, 2017; Chiang *et al.*, 2020; Heitzler and Hurni, 2020). Most of these techniques rely on semantic segmentation, a technology in which each pixel is attributed a semantic class, which can correspond to a type of terrain or a geometric delimitation (Fig. 1). The choice of which class to attribute to each pixel is learnt by the neural network from the large number of training pixels, i.e. annotated examples. It is the requirements of this essential annotation phase that this article attempts to clarify. Vectorization occurs in a second step, with the transformation of this pixel information into vector geometries (points, lines, and polygons). However, the high level of technicality of researches on automatic vectorization and the lack of interoperability between the different tools, algorithms, GIS, or annotation softwares, make these technologies difficult to leverage for most specialists in urban planning, urban history, or geography. This is a first challenge to the development of the practice. To address this pitfall, we advocate for the development of a comprehensive and readily usable graphical interface that would implement these algorithms and include all the main steps of the automatic processing pipeline: IIIF[1] collection upload and management, georeferencing, annotation, free web training of a neural model, automatic vectorization, manual verification, correction, and rectification of the vector output, export in a common and reusable IIIF annotation geoJSON format (IIIF Maps Community Group, 2021).

The second challenge to the expansion of automatic vectorization technologies is the lack of standards, benchmarks, and generic training datasets. This makes published studies rather difficult to compare, and impacts the quality of the results due to limited domain-specific pretraining data. However, several resources were released recently. The largest one is the Historical City Maps Semantic Segmentation Dataset (Petitpierre, 2021), a five-classes ontology dataset gathering 635 annotated patches from 580 diverse historical maps. The second is constituted of five annotated maps provided for the ICDAR21 Competition on Historical Map Segmentation (Chazalon *et al.*, 2021a; Chazalon *et al.*, 2021b), and focuses on the segmentation of building blocks areas from the 'Atlas de Paris' (Service du Plan, 1894). However, subsequent publications have questioned the latter's annotation approach, advocating instead for the pixel annotation of the contours of the building blocks themselves (Chen *et al.*, 2021).

There is indeed a lack of clear documentation in the literature on the annotation process itself, and this lack of empirical knowledge concerning a key step constitutes the third large challenge to the generalization of technologies for automatic comprehension of historical maps. Not only does the annotation methodology have a determining impact on the performance itself, but it is also a time-consuming step that should be optimized. Moreover, the process of annotating for semantic segmentation is quite different from manually vectorizing cartographic documents. The purpose differs deeply: semantic segmentation (Long *et al.*, 2015), the core technology of automatic map vectorization, aims at delimiting pixel areas on the image-document (for instance a parcel), then assigning them a precise meaning, in the form of a semantic class (e.g. the class 'crops'). Manual vectorization, on the other hand, aims to translate the geometry of an historical object in its geographical context, and therefore already constitutes an interpretation of the representation, due to



**Figure 1.** Comparison of the annotation of a patch from the Napoleonic cadastre of Venice (1808). (**A**) Original image. (**B**) Vector annotation of the polygons through a web-interface. (**C**) Pixel level annotation through a graphics editor

simplification, intuition-driven extrapolation, or even modification based on specific knowledge of the historical context. The annotation, in the contrary, is based solely on the source.

For this reason, we believe that it is necessary to open a methodological debate and to found a scholarship of annotation in itself. In this article, we suggest guidelines for the annotation of cadastral documents in order to promote the learning of automatic vectorization models based on neural networks. In this perspective, we will present our empirical conclusions based on five large vectorization projects conducted in the last few years, gathering several hundred cadastral maps corresponding to different places, times, and styles. These five projects are based on: the 'Ancien Régime' cadastre of Lausanne (Melotte and Perey, 1721); the Napoleonic cadastre of Venice (Selva, 1808); and its counterpart in Lausanne, the 'Berney' cadastre (Berney, 1827). The case of the 'renovated' cadastre of Neuchatel (Offenhaüser, 1869) and the ongoing vectorization project of the 1900 Atlas of Paris (Service du Plan, 1894) will also be discussed.

## 2 Discussion

### 2.1 A priori difficulties

The success of automatic vectorization, unsurprisingly, depends both on the content and the visual qualities of the digital image to process. Degradation, poor conservation, or scanning artefacts will yield visual traces that are hard to ignore for the neural network (see Fig. 2A). The content depicted should be as visually explicit as possible and must be interpretable on the basis of figurative, morphological, or topological qualities alone. The interpretation of context-specific objects that requires knowledge of the geographic place to be identified, such as the presence of benches or barriers (see Fig. 2B), will prove difficult in a generic pipeline, just like the untangling of superimposed pieces of information (see Fig. 2C). Therefore, both scanning and digitization must be conducted with care; and, in that respect, the choice of the best preserved version of the document is decisive.
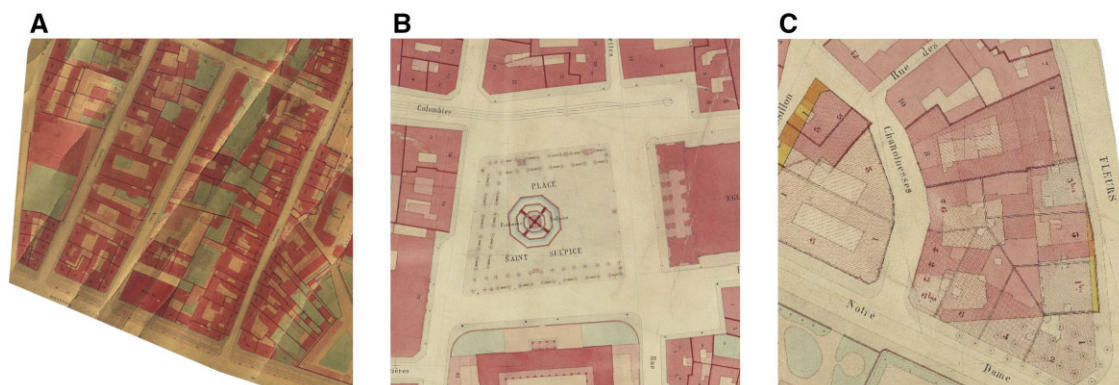
### 2.2 The need for consistent annotation

For larger projects, one might be tempted to divide the annotation work between several collaborators. Thus appears to each annotator the problem of interpreting the limits of the classification—the so-called 'annotation ontology'. For example (Figs. 3A–D and 4C and D), should the sidewalks surrounding a large private building be annotated as 'roads' or 'non-built'? Or should the thin wall surrounding a fountain be considered the edge of that fountain (and hence annotated as a line) or instead as a thin 'built' structure separate from the fountain? These edge cases are hard to gauge beforehand and will appear as soon as the annotations of different collaborators are compared.

Such inconsistencies will hinder the learning of the neural network which will be faced, in such cases, with two conflicting 'right answers' for a similar visual neighbourhood. An efficient annotation phase should take the time to establish a clear and consistent charter for all collaborators to follow, and a well-defined ontology.

### 2.3 Ontology and homogeneity of the representation

The definition of an ontology consists in classifying any element depicted in the source document into a system of semantic classes. The ontology design is decisive and must take into account the homogeneity of the representation. For instance, during the extraction of the renovated cadastre of Neuchatel, the ontology included a specialized 'stairs' class. However, the texture of the



**Figure 2.** Problems anticipated during the annotation of the 1900 Atlas of Paris. (**A**) Document degradation impacting the image quality. (**B**) Multiplication of context-specific objects (lamps, benches, barriers, pillars, etc.) and presence of non-delineating lines (tramway). (**C**) Superposition of several indistinct temporal layers drawn on the same image

**Figure 3.** Excerpt from the collaborative annotations of the renovated cadastre of Neuchatel (1872). In magenta the built, in green the non-built, in beige the road network, in purple the stairs, and in black the contours. (**A** and **B**) Conflict in the annotation of sidewalks, annotated as non-built (green) in A, and as road network (beige) in B. (**C** and **D**) Conflict in the annotation of fountain walls, annotated as building (magenta) in C, and as edges (black) in D



**Figure 4.** Excerpt from the collaborative annotations, and semantic segmentation prediction, of the renovated cadastre of Neuchatel (1872). (**A**) In purple, the specialized stairs class. (**B**) Also in purple the quay walls, falsely predicted as stairs by the neural network. (**C** and **D**) Conflict in the annotation of the walls and parcel boundaries (white), annotated as a triple line in C and as a simple line in D

stairs was very close to that of the quay walls, which might be one of the causes to the confusion observed during semantic segmentation (see Fig. 4A and B). In this case, the confusion risk is further enhanced by the underrepresentation of the stairs class in the training examples, and by the absence of quay walls.

Conventional neural networks are very sensitive to class imbalance in the training datasets, even if some specific solutions, such as focal loss and data augmentation, were developed to address this problem (Doi and Iwasaki, 2018; Li et al., 2021). When the ontology is complete and the representation of the elements in

the annotated data is balanced, however, neural networks used for semantic segmentation demonstrate a high representational flexibility, and are capable of developing advanced strategies for disambiguation, especially in some cases where colour does not allow to distinguish between several classes (such as 'non-built' and 'road network'), or in the contrary when a class (e.g. built) is represented in an expanded colour or texture palette (see Fig. 5). In fact, due to their mathematical architecture, neural networks can rely on abstract concepts such as the semantic hierarchy, the morphology, or even the topology of objects to establish a prediction (Petitpierre *et al.*, 2021). These strategies, which are often very relevant, can also lead to local errors, such as confusing a section of wall with a road, due to its elongated and rectilinear morphology (see Fig. 5). For the particular case of Venice, one of the challenges was to train the algorithm to effectively differentiate the canals from the road network, due to their morphological, figurative, and topological proximity (see Fig. 6). In most cases, however, the algorithm is able to compensate for local variations, as long as the separations between objects (name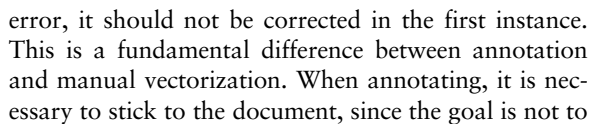ly parcel contours) are clear, by averaging the neighbourhood prediction. As an annotator, these limitations can also be corrected by adding a few well-chosen additional examples.

## 2.4 Consequences of the vectorization

The vectorization process is to be considered as a relatively literal translation of the document visual content into vector shapes. As such, the end result will be as close as possible to the source content, which can prove surprising in some cases. For example, objects whose representations lie at the boundaries of several sheets can be extracted more than once. This is the case for some of the rivers and roads of the Napoleonic cadastre of Lausanne. Similarly, the algorithm will not extrapolate missing or incomplete objects: for example, the gaps between two sheets caused by geometrical shifts in the historical measurement and drawing of the cadastral sheets will not be filled (see Fig. 7). In such instances, it may therefore be interesting to correct or rectify the vectorization result a posteriori (see Fig. 8). This step is made easier by the automatic vectorization process described here, the *raison d'être* of which is not the producing of a perfect result but instead of a first (time saving) 'draft' which the human expert can assess and correct before further analysis.

## 2.5 Formulation of recommendations

At this stage, we will try to summarize and collect our main empirical findings concerning the annotation of cadastral sources. Sharing experiences and practical knowledge should allow, in the future, to build a common and shared expertise and to democratize automatic vectorization technologies. The automatic processing of historical sources has the potential to deeply change research practices in geographical history and related sciences in the next years by promoting large-scale comparative studies and data accessibility.

When annotating a map for semantic segmentation, the document should always be considered as an image, whose pixels are being annotated. It is therefore preferable to take a certain distance from the interpretation and focus instead on the representation. First, the annotation must be concretely embodied in the image. This embodiment must be expressed directly in the visual grammar, through colour, texture, or morphology. For instance, one should not annotate an area as corresponding to a *portico* if the corresponding area is not clearly distinguished from the buildings and the road network by visual cues. The presence of a textual indication (for instance 'vineyard') is usually not sufficient for the vectorization algorithm to learn to differentiate one area from another, or to define a semantic class (see Fig. 9). Moreover, when annotating, all other historical sources must be disregarded. For example, if a feature depicted on the map is an obvious historical

**Figure 5.** Result of the semantic segmentation prediction for the Melotte cadastre of Lausanne (1721–27). In warm tones, the buildings (including the private houses in carmine, the church in purple, and the walls in orange), in beige the road network, and in green the non-built. The arrow points to an example of misclassification between a wall and a street

**Figure 6.** Venice flooded: example of confusion between the road network (transparent) and water (blue) during vectorization. Modified from Oliveira *et al.* (2019)
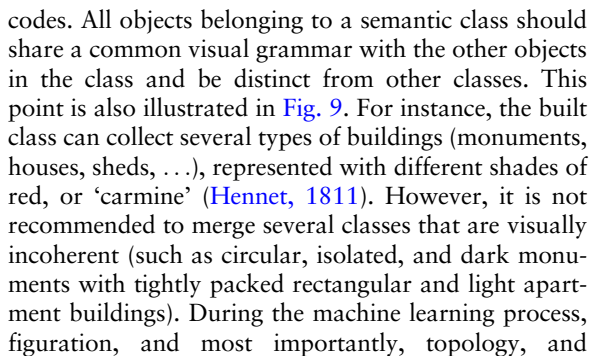


**Figure 7.** Result of the automatic vectorization of the Berney cadaster of Lausanne (1827–31). In magenta the buildings, in beige the road network, in green the non-built, and in blue the water. Yellow continuous arrows point to objects duplicated at the boundary of two cadastral sheets while dashed arrows point to gaps observed between two cadastral sheets due to sources inconsistencies.

error, it should not be corrected in the first instance. This is a fundamental difference between annotation and manual vectorization. When annotating, it is necessary to stick to the document, since the goal is not to vectorize a specific map, but rather to teach an artificial intelligence to recognize visual regularities.
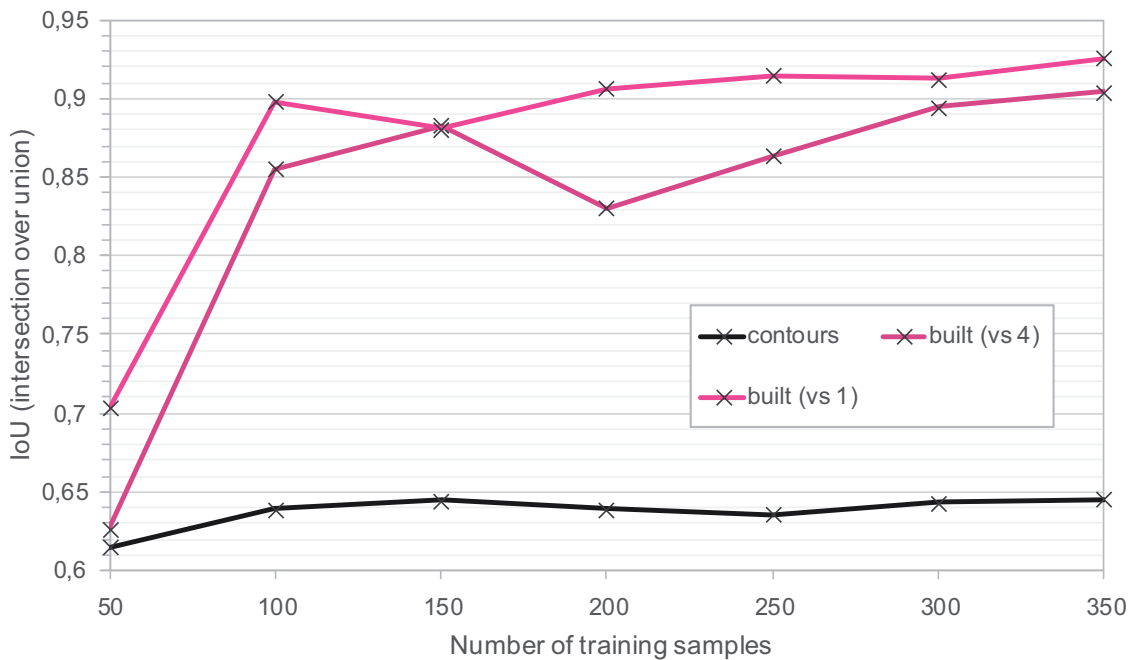
Secondly, a semantic class should be defined by congruent figurative, topological, and/or morphological

**Figure 8.** The Napoleonic cadastre of Venice (1808), after manual correction of the automatically extracted vector layers



**Figure 9.** Example of an ontology not sufficiently sticking to visual cues. The topological and morphological differences between vineyard and non-built are insufficient to discriminate between both classes during inference

codes. All objects belonging to a semantic class should share a common visual grammar with the other objects in the class and be distinct from other classes. This point is also illustrated in Fig. 9. For instance, the built class can collect several types of buildings (monuments, houses, sheds, . . .), represented with different shades of red, or 'carmine' (Hennet, 1811). However, it is not recommended to merge several classes that are visually incoherent (such as circular, isolated, and dark monuments with tightly packed rectangular and light apartment buildings). During the machine learning process, figuration, and most importantly, topology, and

morphology are decisive factors (Petitpierre et al., 2021), and the classes must therefore respect a certain homogeneity in these criteria. It is possible to derogate from this rule for the background, which represents a semantic class in itself, and which is implicit when using a vector annotator. The background includes everything that is not annotated as a specific instance, and therefore gathers 'leftovers' which are not otherwise covered by the annotation ontology.

The third recommendation is to reduce the number of classes to a minimum. The annotation ontology does not have to be exhaustive and multiplying unnecessary

**Figure 10.** Segmentation performance on the Berney cadastre of Lausanne, according to the number of training samples (768 × 768 pixels patches) and the segmentation problem. Built versus 1 refers to a setting in which the built class is segmented against background class. Built versus 4 refers to a setting in which the built class is segmented against background, non-built, water, and road network. The architecture used is state of the art (OCR–HRNetV2 W48, Yuan *et al.*, 2020)

semantic classes will likely hinder learning. Annotation should be goal oriented and purposeful. For instance, if you are conducting a study on the urban development of a city, do prefer a minimal ontology, including only the buildings (see Fig. 10) or the road network. Not only will this practice save you time, but more importantly it will greatly support the performance of the neural network, and thus the semantization of the parcels. An optimal ontology only contains two or three classes. Beyond this number, the performance will decrease progressively (Petitpierre *et al.*, 2021). Besides, the creation of specialized classes for a few occurrences can lead to an imbalance phenomenon, which is a well-known pitfall in machine learning, occurring when the number of learning examples are severely unequal (Thabtah *et al.*, 2020). The aim of automatic vectorization is to spare work time by having a computer perform the most repetitive tasks. Specific cases can be supervised later by a human expert.

Fourth, develop and apply the same rationale to all annotations. Seek impartiality and consistency in annotation. When an ambiguous case arises, it is sometimes necessary to take a stand. For instance, one might wonder whether it is better to include a paddle wheel in the 'river' class or in the built class. There is no wrong answer a priori, but once the decision is made, its application should be systematic.

Finally, the last recommendation is to adopt an iterative approach. Start with a small number of annotations, and then verify the preliminary results. Retrain the neural network with more annotations if necessary. In most cases, a dozen of cadastral sheets is already sufficient to bootstrap the model and reach an excellent performance (Petitpierre, 2020). However, the number of annotations needed depends on the annotation precision and the complexity of the task (in particular the number and the homogeneity of the classes). What is considered an acceptable performance (e.g. accuracy or intersection over union) can also vary, depending on whether one extracts linear features, such as contours for instance, or areal features, such as buildings (Fig. 10). Testing can also help to understand the weaknesses of the chosen annotation ontology and to refine the method if necessary. Trying to establish the prediction biases of the algorithm usually allows to visually identify the most complex aspects of the study case.

## 3 Conclusion

In this article, we have presented a set of guidelines, based on empirical observations and concrete examples taken from several large-scale projects. These recommendations address a gap in the scientific literature left by the lack of documentation on annotation processes

for semantic segmentation and vectorization in general, especially on the case of cadastral documents. We believe that the production of a scientific documentation on this topic is both necessary to promote the transmission of these complex technologies from computer science to other fields of research, and to investigate the intrinsic functioning of neural networks, often considered opaque.

The guidelines are summarized as follows. (1) The annotated object must be understandable and classifiable exclusively by visual cues (colour, texture, or morphology); and the hidden semantics must be disregarded. (2) The objects grouped in a semantic class should both be visually distinct from other classes and share visual characteristics that justify their grouping. (3) The number of semantic classes should be reduced to a minimum. (4) Annotations must be consistent. (5) An iterative approach should be favoured, multiplying and eventually correcting the annotations once preliminary results have been verified.

It is our hope that the observations and principles explored in this article will help foster the transfer of knowledge between computer science and geographical history, the automation of which holds immense potential. According to some estimates, cadastral records are likely representing several hundred million parcels in Europe alone (Clergeot, 2007). They represent a highly reliable and detailed geohistorical source, which allows not only to study the urban environment, but also to weave links between many sources, and thus to investigate social, historical, and economical questions. Given the vastness of the data, the time-consuming nature of manual processing, the existence of large and relatively homogeneous corpora, and the strategic importance of such processing for studies in urban history, the development of tools, and the creation of knowledge on automatic vectorization practices seems to be an essential and priority issue.

## Funding

## Notes

1. International Image Interoperability Framework.

## References

Balkoca, A., Yergök, A.I., and Yücekaya, S. (2011). Vectorization of cadastral maps using image processing algorithms. In *2011 IEEE 19th Signal Processing and Communications Applications Conference (SIU)*, 2011, pp. 900–3. IEEE. https://doi.org/10.1109/SIU.2011.5929797

Berney, A. (1827). Plan cadastral par l'arpenteur Abraham Berney. Archive of the City of Lausanne.

Candeias, T., Tomaz, F., and Shahbazkia, H. (2003). Off the shelf methods for robust Portuguese cadastral map analysis. In Sanfeliu, A. and Ruiz-Shulcloper, J. (eds), *Progress in Pattern Recognition, Speech and Image Analysis*. Berlin, Heidelberg: Springer, pp. 627–34. https://doi.org/10.1007/978-3-540-24586-5_77

Chazalon, J., Carlinet, E., Chen, Y., *et al.* (2021a). *ICDAR 2021 Competition on Historical Map Segmentation*. https://arxiv.org/abs/2105.13265

Chazalon, J., Carlinet, E., Chen, Y., *et al.* (2021b). *ICDAR 2021 Competition on Historical Map Segmentation—Dataset*, Zenodo. https://doi.org/10.5281/zenodo.4817662 (accessed 7 February 2022).

Chen, L.H., Liao, H.Y., Wang, J.Y., *et al.* (1996). An interpretation system for cadastral maps. In *Proceedings of 13th International Conference on Pattern Recognition*, August 1996, Vol. 3, pp. 711–5. IEEE. https://doi.org/10.1109/ICPR.1996.547261

Chen, Y., Carlinet, E., Chazalon, J., *et al.* (2021). Vectorization of historical maps using deep edge filtering and closed shape extraction. In Lladós, J., Lopresti, D., and Uchida, S. (eds), *Document Analysis and Recognition—ICDAR 2021*. Cham: Springer International Publishing, pp. 510–25.

Chiang, Y.-Y., Leyk, S., and Knoblock, C.A. (2014). A survey of digital map processing techniques. *ACM Computing Surveys*, 47(1): 1–44. https://doi.org/10.1145/2557423

Chiang, Y.-Y., Duan, W., Leyk, S., *et al.* (2020). *Using Historical Maps in Scientific Studies: Applications, Challenges, and Best Practices. Springer Briefs in Geography*. Cham: Springer International Publishing. https://doi.org/10.1007/978-3-319-66908-3

Clergeot, P. (2007). Cent Millions de Parcelles En France: 1807, Un Cadastre Pour l'Empire, Editions Publi-Topex.

Cofer, R.H. and Tou, J.T. (1972). Automated map reading and analysis by computer. In *Proceedings of the December 5–7, 1972, Fall Joint Computer Conference, Part I*, New York, AFIPS '72 (Fall, Part I), Association for Computing Machinery, pp. 135–45. https://doi.org/10.1145/1479992.1480010.

Costes, B. (2016). *Vers la Construction d'un Référentiel Géographique Ancien: Un Modèle de Graphe Agrégé pour Intégrer, Qualifier et Analyser des Réseaux Géohistoriques*. Ph.D. thesis, Université Paris-Est. https://tel.archives-ouvertes.fr/tel-01565850 (accessed 17 January 2022).

Cousins, S., Ohlson, H., and Eriksson, O. (2007). Effects of historical and present fragmentation on plant species diversity in semi-natural grasslands in Swedish rural landscapes. *Landscape Ecology*, 22: 723–30. https://doi.org/10.1007/s10980-006-9067-1

Dhar, D.B. and Chanda, B. (2006). Extraction and recognition of geographical features from paper maps. *International Journal of Document Analysis and Recognition (IJDAR)*, 8(4): 232–45. https://doi.org/10.1007/s10032-005-0010-9

Doi, K. and Iwasaki, A. (2018). The effect of focal loss in semantic segmentation of high resolution aerial image. In *IGARSS 2018—2018 IEEE International Geoscience and Remote*

*Sensing Symposium*, July 2018, pp. 6919–6922. IEEE. https://doi.org/10.1109/IGARSS.2018.8519409

**di Lenardo, I., Barman, R., Pardini, F.,** *et al.* (2021). Une approche computationnelle du cadastre napoléonien de Venise. *Humanités numériques*, **3**. https://doi.org/10.4000/revuehn.1786

**Digital Humanities Laboratory**. (2021). dhSegment-torch. https://github.com/dhlab-epfl/dhSegment-torch (accessed 6 February 2022).

**Dolej, M. and Forejt, M.** (2019). Franziscean cadastre in landscape structure research: a systematic review. *Quaestiones Geographicae*, **38**(1): 131–44. https://doi.org/10.2478/quageo-2019-0013

**Domaas, S.T., Hamre, L.N., and Austad, I.** (2003). Historical cadastral maps as a tool for identifying key biotopes in the cultural landscape. *WIT Transactions on Ecology and the Environment*, **64**: 913–24. https://doi.org/10.2495/ECO030162

**Ekamper, P.** (2010). Using cadastral maps in historical demographic research: some examples from the Netherlands. *The History of the Family*, **15**(1): 1–12. https://doi.org/10.1016/j.hisfam.2010.01.003

**Frischknecht, S. and Kanani, E.** (1998). Automatic interpretation of scanned topographic maps: a raster-based approach. In Tombre, K. and Chhabra, A. K. (eds), *Graphics Recognition Algorithms and Systems*, Berlin, Heidelberg: Springer Berlin Heidelberg, pp. 207–20. https://doi.org/10.1007/3-540-64381-8_50

**Heitzler, M. and Hurni, L.** (2020). Cartographic reconstruction of building footprints from historical maps: a study on the Swiss Siegfried map. *Transactions in GIS*, **24**(2): 442–61. https://doi.org/10.1111/tgis.12610

**Hennet, A.-J.-U.** (1811) *Recueil Méthodique Des Lois, Décrets, Règlemens, Instructions et Décisions Sur Le Cadastre de La France*. Paris: Imprimerie Impériale.

**Ignjatić, J., Nikolić, B., and Rikalović, A.** (2018). Deep learning for historical cadastral maps digitization: overview, challenges and potential. In *26th International Conference in Central Europe on Computer Graphics, Visualization and Computer Vision*, Václav Skala—UNION Agency, pp. 42–7. https://doi.org/10.24132/CSRN.2018.2803.6

**IIIF Maps Community Group**. (2021). *navPlace Extension*. Available online at: https://iiif.io/api/extension/navplace/ (accessed 6 February 2022).

**Kain, R.J.P. and Baigent, E.** (1992). *The Cadastral Map in the Service of the State: A History of Property Mapping*. Chicago: University of Chicago Press.

**Katona, E.** (2000). *Automatikus Térkép-Interpretáció*. Ph.D. thesis, University of Szeged.

**Katona, E. and Hudra, G.** (1999). An interpretation system for cadastral maps. In *Proceedings of the 10th International Conference on Image Analysis and Processing*, September 1999, pp. 792–7. IEEE. https://doi.org/10.1109/ICIAP.1999.797692

**Lelo, K.** (2020). Analysing spatial relationships through the urban cadastre of nineteenth-century Rome. *Urban History*, **47**(3): 467–87. https://doi.org/10.1017/S0963926820000188

**Li, K., Konstantinos, K., and Glocker, B.** (2021). Analyzing overfitting under class imbalance in neural networks for image segmentation. *IEEE Transactions on Medical Imaging*, **40**(3): 1065–77. https://doi.org/10.1109/TMI.2020.3046692

**Liu, C., Wu, J., Kohli, P.,** *et al.* (2017). Raster-to-vector: revisiting floorplan transformation. In *2017 IEEE International Conference on Computer Vision (ICCV)*, October 2017, pp. 2214–22. IEEE. https://doi.org/10.1109/ICCV.2017.241

**Long, J., Shelhamer, E., and Darrell, T.** (2015). Fully convolutional networks for semantic segmentation. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pp. 3431–40. IEEE. https://www.cv-foundation.org/openaccess/content_cvpr_2015/html/Long_Fully_Convolutional_Networks_2015_CVPR_paper.html (accessed 22 January 2022).

**Melotte, S. and Perey, C.** (1721). Plans du territoire de Lausanne, dit plan Melotte. Archive of the City of Lausanne.

**Mou, Z.** (2012). Using cadastral maps in historical GIS research: the French concession in Shanghai (1931–1941). *Annals of GIS*, **18**(2): 147–56. https://doi.org/10.1080/19475683.2012.668560

**Offenhaüser, R.** (1869). Plan cadastral du territoire de Neuchâtel, *district de Neuchâtel*. Geomatics and Land Registry Service of the Republic and Canton of Neuchâtel.

**Oliveira, S.A., Seguin, B., and Kaplan, F.** (2018). dhSegment: a generic deep-learning approach for document segmentation. In *2018 16th International Conference on Frontiers in Handwriting Recognition (ICFHR)*, Niagara Falls, August 2018, IEEE, pp. 7–12. https://doi.org/10.1109/ICFHR-2018.2018.00011

**Oliveira, S.A., di Lenardo, I., Tourenc, B.,** *et al.* (2019). A deep learning approach to cadastral computing. In *DH2019*, Utrecht, Netherlands, 8 July 2019. https://dev.clariah.nl/files/dh2019/boa/0691.html (accessed 22 January 2022).

**Petitpierre, R.** (2020). *Neural Networks for Semantic Segmentation of Historical City Maps: Cross-Cultural Performance and the Impact of Figurative Diversity*. Master's thesis, Ecole Polytechnique Fédérale de Lausanne. https://doi.org/10.13140/RG.2.2.10973.64484

**Petitpierre, R.** (2021). *Historical City Maps Semantic Segmentation Dataset*, Zenodo. https://doi.org/10.5281/zenodo.5497934 (accessed 6 February 2022).

**Petitpierre, R., Kaplan, F., and di Lenardo, I.** (2021). Generic semantic segmentation of historical maps. In *CHR 2021: Computational Humanities Research Conference* (*CEUR Workshop Proceedings*), Amsterdam, The Netherlands, 17 November 2021, pp. 228–48. http://ceur-ws.org/Vol-2989/long_paper27.pdf (accessed 6 February 2022).

**Pradhan, R., Kumar, S., Agarwal, R.,** *et al.* (2010). Contour line tracing algorithm for digital topographic maps. *International Journal of Image Processing (IJIP)*, **4**(2): 156–63. https://www.academia.edu/download/31189443/IJIP-149.pdf (accessed 7 February 2022).

**Selva, G.** (1808). *Censo Stabile, Catasto Napoleonico*. Venice State Archive.

**Service du Plan**. (1894). Atlas municipal des vingt arrondissements de Paris.

**Thabtah, F., Hammoud, S., Kamalov, F.,** *et al.* (2020). Data imbalance in classification: Experimental evaluation. *Information Sciences*, **513**: 429–41. https://doi.org/10.1016/j.ins.2019.11.004

**Valent, P., Roncák, P., Maliariková, M.,** *et al.* (2016). Utilization of historical maps in the land use change impact

studies: A case study from Myjava river basin. *Slovak Journal of Civil Engineering*, **24**(4): 15.

Viglino, J.-M. and Pierrot-Deseilligny, M. (2003). A vector approach for automatic interpretation of the French cadastral map. In *Proceedings of the Seventh International Conference on Document Analysis and Recognition*, August 2003, Vol.

**1,** pp. 304–8. IEEE. https://doi.org/10.1109/ICDAR.2003.1227678

Yuan, Y., Chen, X., and Wang, J. (2020). Object-contextual representations for semantic segmentation. In *ECCV 2020: European Conference on Computer Vision*. Cham: Springer, pp. 173–90. https://doi.org/10.1007/978-3-030-58539-6_11