

Mapping Change: A Temporal and Semantic Knowledge Base of Scottish Gazetteers (1803–1901)

Lilin Yu, Rosa Filgueira

EPCC, University of Edinburgh

L.Yu-40@sms.ed.ac.uk, r.filgueira@epcc.ed.ac.uk

Identifier: <https://rosafilgueira.github.io/MappingChange-Paper-ISWC2025/>

Abstract

We present **MappingChange**, a project that constructs a temporal and semantic knowledge base from ten 19th-century **Gazetteers of Scotland** (1803–1901), digitized as over 13,000 page-level **XML** files. These noisy, unstructured texts lack article-level markup and exhibit highly heterogeneous layouts. To segment and structure over 50,000 historical place descriptions, we employ **large language models (LLMs)** with **edition-specific prompting strategies**, tuned to handle distinct editorial conventions, abbreviations, and multi-page entries. The resulting knowledge base comprises three interlinked knowledge graphs: (1) a **basic KG**, extracted from cleaned DataFrames; (2) a **concept-enriched KG**, linking semantically similar place records across editions using **sentence embeddings**, **Wikidata**, and **DBpedia**; and (3) a **location-annotated KG**, enriched with **named entity recognition** and **geographic disambiguation**. We further align extracted entities to external sources such as Wikidata and DBpedia, enabling rich contextual integration and reuse. All are expressed in **RDF** and modeled with the updated **Heritage Textual Ontology (HTO)**, which provides a structured vocabulary for capturing textual provenance, bibliographic metadata, extraction context, and diachronic semantic alignment across editions. In addition to the knowledge graphs, we release: (a) individual DataFrames for each edition, (b) a unified cross-edition DataFrame, and (c) Elasticsearch indices. All resources are integrated into the Frances semantic web platform, enabling historical exploration through keyword and semantic search, concept-timeline navigation, and interactive geolocation visualizations.

1. Introduction

Descriptive gazetteers were central to how 19th-century Scotland documented its geography—towns, parishes, rivers, castles, lochs, and glens—embedding each place within broader historical, economic, and social narratives. These texts evolved over the century, reflecting transformations brought about by industrialization, land reform, public health, and imperial expansion. The *Gazetteers of Scotland, 1803–1901*, digitized by the National Library of Scotland (NLS), constitute one of the most comprehensive corpora for studying Scotland’s spatial knowledge in the long 19th century. The full collection comprises twenty volumes produced by different publishers and editors, and has been released as more than

13,000 high-resolution scans accompanied by ALTO XML files. These XML files encode layout and textual content extracted via Optical Character Recognition (OCR), resulting in over 1.75 million lines and 14 million words. While this makes the data technically accessible, it remains largely unsuitable for structured analysis: the texts lack article-level markup, exhibit inconsistent typographic structures, and contain significant OCR noise. Entries often begin mid-column, span multiple pages, and vary widely in format and editorial style—posing major challenges for computational processing, information retrieval, and historical reuse.

Compounding these challenges is the fact that many place names (e.g., “ABBEY” or “GREENHILL”) recur across the gazetteers, often referring to different locations. Disambiguating such entries is non-trivial, as it depends on contextual clues within each article rather than surface-level patterns. Our approach relies on LLM-based article segmentation and interpretation—capturing subtle editorial cues and semantic context to correctly associate each name with the appropriate description.

MappingChange is the first project to construct a structured, queryable, and semantically enriched temporal knowledge base from this entire collection. We extract and align over 50,000 historical place descriptions across ten gazetteer editions, using large language models (LLMs) and volume-specific prompting strategies that are carefully tuned to editorial idiosyncrasies. The result is a knowledge base composed of three interlinked knowledge graphs: a basic graph derived from structured DataFrames; a concept-enriched graph linking semantically similar entries across editions using sentence embeddings, Wikidata, and DBpedia; and a location-annotated graph generated through geographic disambiguation techniques. These graphs are serialized in RDF and modeled using the Heritage Textual Ontology (HTO), a domain ontology we developed specifically for historical and heritage corpora.

The **Heritage Textual Ontology (HTO)** (see resource [here](#)) is designed to model not just entities and attributes but also the editorial and computational processes by which each record is extracted, cleaned, and enriched. Unlike generic ontologies, HTO supports the representation of textual provenance, extraction prompts, editorial hierarchies, and diachronic linkage across editions. It enables us to track how descriptions of the same place evolve over time, with full transparency into their source structure and transformation process. Its design has been guided by real-world use cases in digital heritage, and it plays a central role in making the resulting knowledge graphs both expressive and reproducible.

The complexity and variability of these sources can be seen in Figure 1, which presents the opening pages of two editions: the 1803 *Gazetteer of Scotland* and the 1884 *Ordnance Gazetteer of Scotland*. These differences, compounded across volumes, necessitate a custom approach to segmentation, prompting, and post-processing—especially since no edition includes machine-readable metadata or reliable article delimiters. Note that

6.1 Data Access

- All cleaned DataFrames and RDF graphs are in the GitHub repository: francesNLP/MappingChange
- Scripts for reproducing those dataframes, KGs and ES are in the GitHub repository: francesNLP/MappingChange
- Zenodo DOI (to be added)

6.1 SPARQL Querying

- A Fuseki SPARQL server supports knowledge graph exploration.
- Sample queries for retrieving places, concepts, and links are included.

6.1 Frances Platform

- Users can search and explore articles via full-text or semantic search.
- Concepts are visualized through timelines and embeddings.

6.1 Notebooks

Google Colab notebooks are provided for each gazetteer to enable direct exploration and analysis.

7. Sustainability

Mapping Change is designed to support long-term historical research:

- **Archiving:** All code, data, and RDF outputs are versioned and archived on Zenodo.
- **Ontological Reuse:** The HTO ontology is maintained and extended in an open repository with permanent identifiers.
- **Frances Platform Integration:** The data is accessible through a production-ready semantic platform, ensuring ongoing usability beyond the scope of the project.
- **Extensibility:** The pipeline is modular and supports the integration of new volumes, editions, or other regional gazetteers.

8. Conclusion

Mapping Change creates a temporal, semantic infrastructure for exploring Scottish place descriptions from 1803–1901. Combining LLM-based extraction, improved ontology design, and semantic search, we deliver a reusable, interoperable dataset for historical research.

The improved HTO ontology enables robust modeling of textual provenance, record quality, and evolving concepts. The Frances platform empowers researchers to query and visualize this data across time and space.

Future work includes integrating cartographic metadata, and link it to the 100 years of the Encyclopaedia Britannica.

9. Acknowledgements

This work was supported by the Royal Society of Edinburgh (RSE Small Research Grant).

@article{semanticweb, title = {The semantic web}, author = {Berners-Lee, Tim and Hendler, James and Lassila, Ora and others}, journal = {Scientific American}, volume = {284}, number = {5}, pages = {28–37}, year = {2001}, publisher = {New York, NY, USA:} }