# Marking Up Cultural Materials for Time and Geography

Fredric Gey, Ryan Shaw, Ray Larson,
Michael Buckland, Barry Pateman and Dan Melia

[1] University of California, School of Information,
South Hall, Berkeley, CA 94720-4600 USA
{gey,dmelia}@berkeley.edu {ryanshaw,ray,buckland}@ischool.berkeley.edu

**Abstract.** This paper describes two projects which fundamentally change access to cultural materials by automatically marking up textual materials for time, events, geography, and people. Because of the imperfect nature of named entity recognition within natural language processing, it is also important to allow for flexible, interactive markup of these materials by editors, scholars and compilers of cultural collections. Relating people in time and place may help to explore hitherto unrecognized relationships in history and social movements.

**Keywords:** biography, document markup, named entity recognition, Irish and Celtic studies, Emma Goldman

## 1  Introduction

Cultural heritage, history, and social sciences are fundamentally about human activity. Everyone is interested in what other people do and have done. Biographies are regularly among the best-selling books. Not only History, but also Geography and most other subjects can come alive in the travelogues, journeys of discovery, and the life-stories of those involved. But mere narrative is not enough. *Understanding the context* differentiates education from memorizing. Building and supporting a community of learners needs more than facts. Understanding the *circumstances* of people's actions illuminates their lives. However, there is a significant gap in the infrastructure developed by libraries, museums, and publishers in this area. We have standards for the computerized handling people's *names*, but not for their *lives*. Take an example from American history, the US frontier hero Davy Crockett – how many people know that after serving as a Colonel under General Andrew Jackson, he was twice elected a congressman from Tennessee? These facts leads to other contextual questions: What was the population of Tennessee when Crockett was a congressman? Who were his contemporaries in Congress (e.g. Daniel Webster)? We will show mapping of biographies and historical events and their demographic contexts being developed in two closely-related projects: *Bringing Lives to Light: Biography in Context* and *Contexts and Relationships: Ireland and Irish Studies*. [1]

---

[1] http://ecai.org/imls2006, http://ecai.org/neh2007

## 2 Biography markup

Our presentation will include prototype interfaces and an "under-the-hood" look at pioneering XML event markup standard proposals which facilitate the demonstration. We have developed named entity markup grammars using the Gate technology [1] and, in particular, the Golden Gate XML editor from Karlsruhe [2] to automatically markup biographies for time and place. The work is similar in spirit to that of Smith and Crane [3]. Figure 1 below shows parts of two screens – to the right is the official USA congressional biography for Davey Crockett from the Congressional Research Service web site.

To the left in the figure is an XML version of the biographical text which has automatically marked up state and county names (counties are political sub-jurisdictions of states in the USA) using named entity recognition. Combining such markups with lookup from a gazetteer such as the Geonames server,[2] will add latitude and longitude to the marked up data, enabling locating parts of a personal biography on online maps such as Google maps or Google Earth.



Figure 1: Split screenshots of a biography page and its XML markup of geographic names

## 3    Markup for Time and Geography

In addition to working with USA congressional biographies, we have followed the timeline of the great anarchist thinker and lecturer Emma Goldman as she toured the USA from 1910 to 1916, lecturing about drama, status of women, and anarchists such as Francisco Ferrer y Guardia, the Spanish-Catalan anarchist. Developing a timeline marked up with geography using the XML Atom format[3] allows us to input the data file directly into Google Maps and provide a geographic user interface to the Goldman travels and lectures. Working from a 5-page textual narrative of the Goldman travels provided by our partner, the Emma Goldman papers project,[4] the application of a Gate grammar marked up over 900 personal names of individuals who co-lectured or were otherwise associated with Emma Goldman. Figure 2 provides a screenshot of the resulting interface.[5]



Figure 2: Automatically create a mapping interface to a biographical timeline

## 4   Markup Options for Newly Digitized Scholarly Materials

Scholarly materials are currently being digitized at a prodigious rate in the areas of digital humanities. Advanced digitization technology offers the opportunity for interactive interfaces to enable knowledgeable scholars and editors to supply their

---

[3] http://en.wikipedia.org/wiki/Atom
[4] http://sunsite.berkeley.edu/goldman/
[5] http://gray.ischool.berkeley.edu/omeka/

own markup which can then be interfaced to external related contextual resources. Our collaborative work with the Centre for Data Digitisation and Analysis at Queens University, Belfast has made available a large number of scanned documents on Irish history, culture, and scholarship. Because the digital technology can supply horizontal and vertical coordinates for each word in the scanned text, we can allow editors and scholars (and users of the materials) the opportunity to highlight segments of text and identify them as names or places (for example). Figure 3 shows a screenshot where a page from The Irish Review for March 1911 has a user-highlighted segment to identify the Irish poet Aubrey De Vere. Links are automatically made to external resources which related to this person.
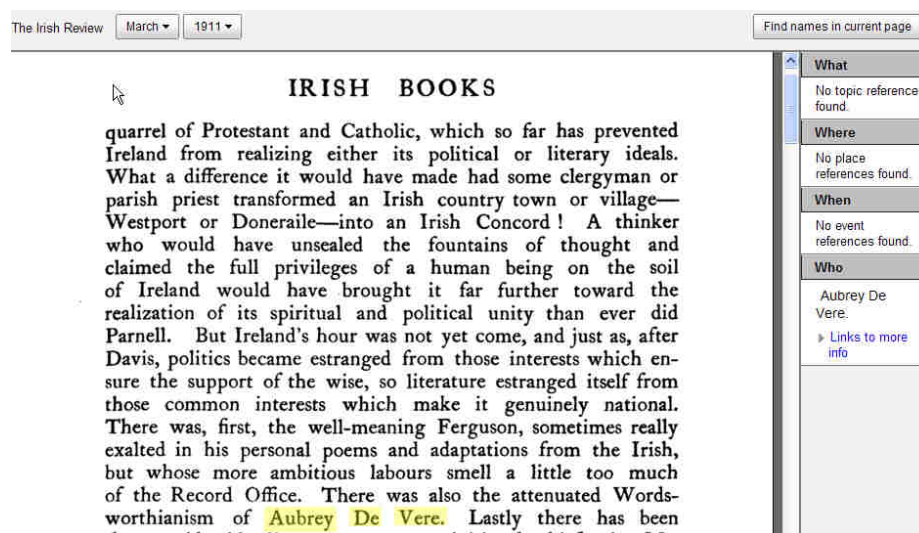


Figure 3: Manual editorial markup of scanned scholarly materials for a personal name

Note that the viewer can add markup on four facets: what, where, when and who, not just where and when (geography and time).

## References

1. Cunningham, H, Maynard, D, Bontcheva, K, Tablan, V: GATE: A Framework and Graphical Development Environment for Robust NLP Tools and Applications. *Proceedings of the 40th Anniversary Meeting of the Association for Computational Linguistics (ACL'02)*, Philadelphia, US.
2. Sautter, G, Padberg, K, Tichy, W: Empirical Evaluation of Semi-Automated XML Annotation of Text Documents with the GoldenGATE Editor. In: L Kovács, N Fuhr, C Meghini (eds.) ECDL 2007. LNCS, vol. 4675, pp. 357--367. Springer, Heidelberg (2007).
3. Smith, D, Crane, G: Disambiguating Geographic Names in a Historical Digital Library. In P Constantopoulos, I Sølvberg (eds.). ECDL 2001. LNCS, vol.2163, pp. 127 – 136. Springer, Heidelberg (2001)