# Human Language Technology for Automatic Annotation and Indexing of Digital Library Content

**Kalina Bontcheva, Hamish Cunningham**
Dept. of Computer Science
University of Sheffield
Regent Court, 211 Portobello St
Sheffield, S1 4DP, UK
`kalina,hamish@dcs.shef.ac.uk`

May 3, 2002

This demo is intended to present the domain-independent and customisable human language technology and the way it was applied for annotating $18th$ century OldBailey proceedings and indexing multimedia content. This demo is intended to accompany the paper with the same title, submitted to ECDL.

Each of these applications posed a unique challenge: the court trials required adapting the language processing components to the non-standard written conventions of old English, while the football collection presented the challenge of indexing material in multiple modalities - video, audio, semi-structured documents, and free text (newspaper reports). In both cases, we used as a basis our robust and customisable named entity recogniser, which comes as part of ANNIE- A Nearly New Information Extraction system. In OldBaileyIE we also used the graphical environment of GATE, which allows manual annotation verification and correction to be carried out on the processed texts.

First we will demonstrate the set of language technology tools which constitute ANNIE and the visual environment which is used for creating and customising new applications (see Figure 1).

Next we will demonstrate the automatic annotation of the OldBailey texts (Figure 2) and present the visual environment used for annotation error correction (Figure 3). We will also demonstrate how this environment can be used for manual annotation, independent of the language technology tools.
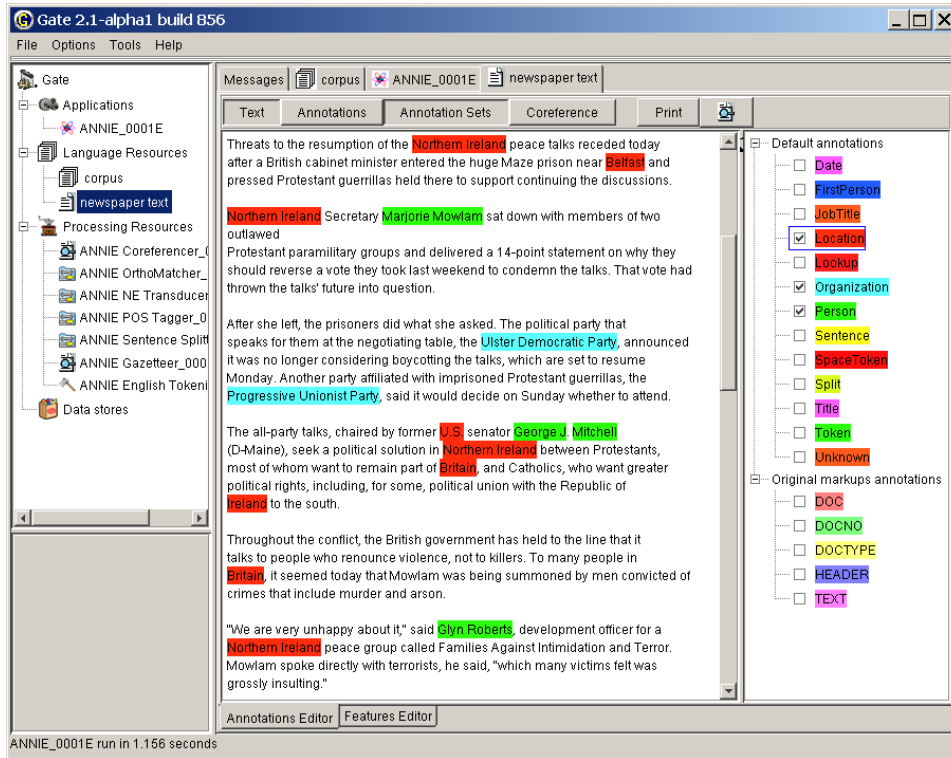
Figure 1: Named Entities recognised by ANNIE

Finally we will demonstrate how multimedia content can be indexed and searched using human language technology. The demonstrator is in the football domain, where video and textual sources about matches are indexed and can be accessed via a user-friendly interface (see Figure 4).
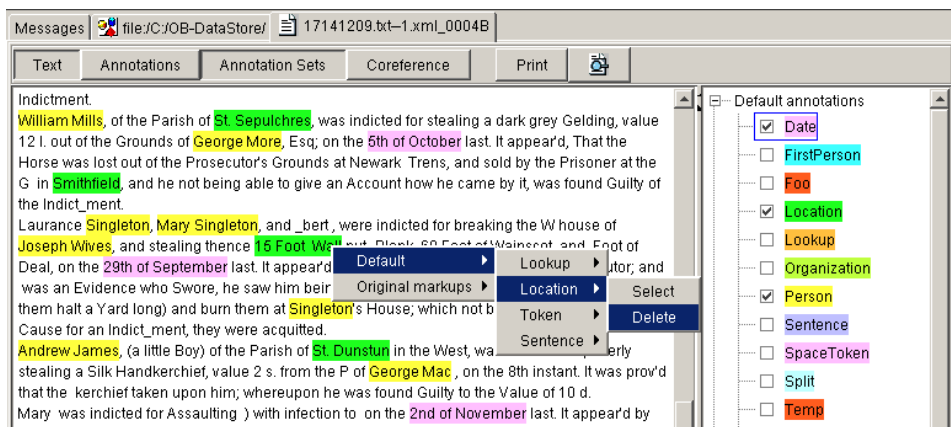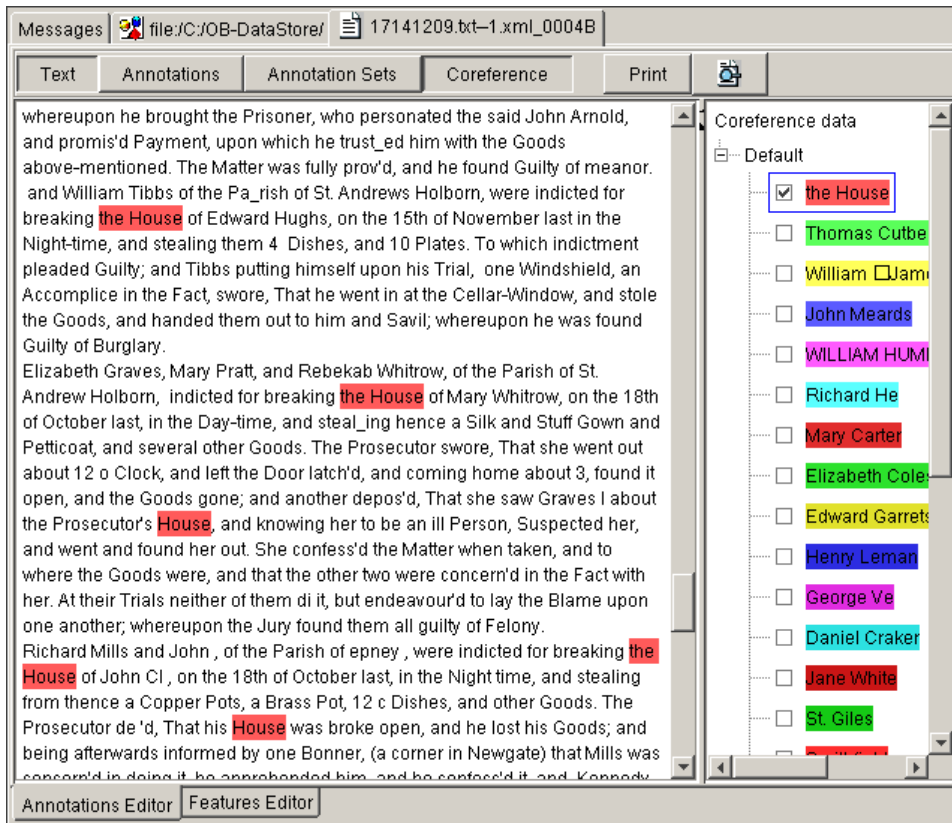
Figure 2: An Old Bailey example text from 1714
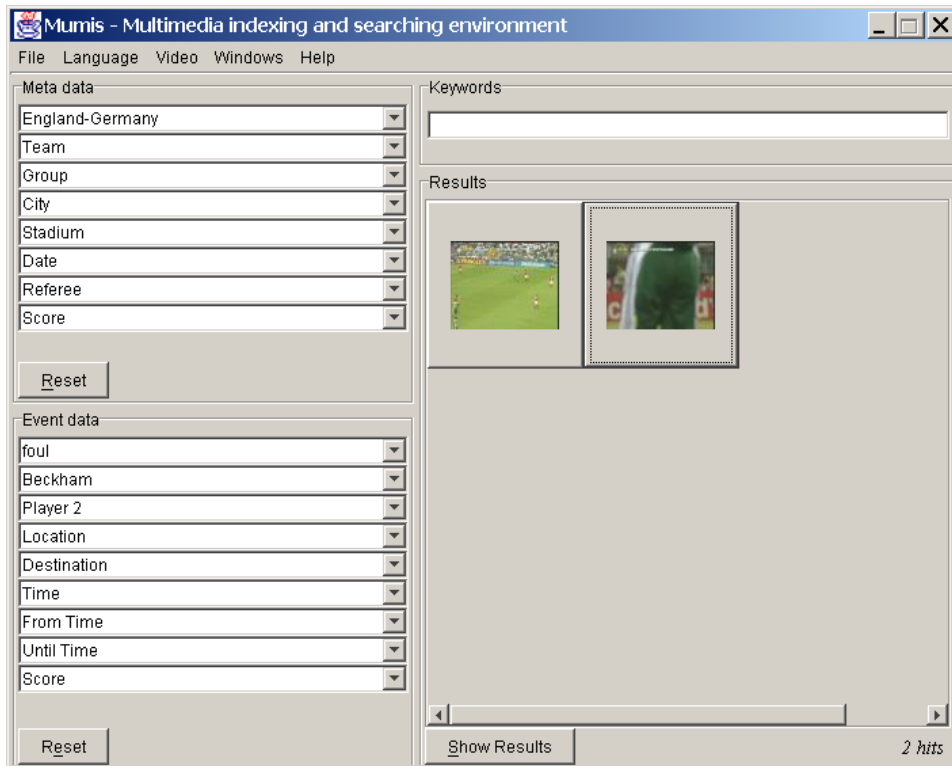
Figure 3: Viewing and deleting all annotations for an entity

Figure 4: The MUMIS interface