

An End-to-End PDF Toolchain for Marking Up Scientific Documents

Sanna Hulkkonen, Oliver Ray

University of Bristol, UK

National Institute of Informatics, Japan

sh16145.2016@my.bristol.ac.uk, csxor@bristol.ac.uk

Abstract

This paper proposes a system for making sentence-level semantic enrichment of scientific publications more user-friendly by developing an end-to-end toolchain for augmenting PDFs with automatically determined textual annotations and visual highlights. The aim is to categorise each sentence according to a given classification scheme and display the labels in a visually appealing way that preserves document structure and formatting while allowing users to work with standard PDF tools they are already accustomed to. This is in contrast to existing approaches which provide an XML representation of document content obtained by abstracting away formatting and structural details in order to focus on the raw text. In particular, we present a toolchain that automatically marks up each sentence in the body of a PDF with a Core Scientific Concept category using a classifier trained with a corpus of papers on social insect biology that we manually labelled. Preliminary testing with domain experts provides anecdotal evidence that end-users do find such automatically derived sentence classifications useful and that they prefer to work directly with marked up PDFs.

Keywords: Portable Document Format, Mark Up, Core Scientific Concepts

1. Introduction

This work is concerned with the semantic enrichment of scientific publications using sentence-level classifications like Argumentative Zoning (AZ) (Teufel, 1999) or Core Scientific Concepts (CoreSC) (Liakata et al., 2012). Its specific focus is on finding sentence labels with a combination of Natural Language Processing (NLP) and Machine Learning (ML) and projecting the labels back onto the original document using PDF manipulation tools.

Our goal is to display sentence labels in a visually intuitive way that preserves document formatting while allowing users to work with standard PDF software they are already accustomed to. This is in contrast to existing approaches which provide an XML representation of a document's content by abstracting away formatting details so the plain text of each sentence can be enclosed within semantic tags. While XML is not a suitable filetype for most end-users to work with, the sentence labels can be conveniently visualised using a text-based annotation interface like Brat (Stenetorp et al., 2012), as shown in Figure 1:

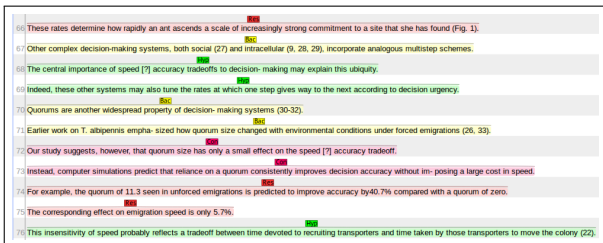


Figure 1: Typical view of sentence mark up.

Our main contribution is developing an end-to-end toolchain that inserts CoreSC annotations directly into the original PDF using classifiers trained on a corpus of papers from social insect biology that we labelled ourselves. Preliminary tests with domain experts suggest they do find such annotations useful and that they prefer to work with marked up PDFs of the sort illustrated in Figure 2:



Figure 2: Proposed view of sentence mark up.

2. System

Our initial work for generating marked up PDFs of the form shown in Figure 2 is detailed in (Hulkkonen, 2017). The motivation for that work comes from a requirement to present social insect biologists with intuitive sentence-level annotations highlighting some key aspects of a scientific investigation as formalised by the CoreSC ontology (Liakata et al., 2012). To do this, we used the Sapienia system (Liakata et al., 2012) to convert PDFs into XML and classify sentences

among 11 CoreSC categories which we used 4 colours to highlight: Motivation, Goal, Object, Hypothesis, Background, Method, Model, Experiment and Observation, Result, Conclusion. This resulted in an output of the form shown in Figure 1.

After that we used pdftotext to extract bounding boxes of the words in the classified sentences in order to project the annotations back onto the PDF as shown in Figure 2.

But, while this is mostly satisfactory, we discovered some incompatibilities between Sapienta’s XML conversions (that use the PDFX utility for textual content extraction) and our back projections (that use pdftotext for bounding box extraction) which mean that some highlighting errors are unavoidable using this method.

To get around this, we re-implemented our approach using a prototype PDF toolkit called PDFNLT that is currently being developed by the Aizawa lab at the Japanese National Institute of Informatics (NII). The advantage of this system is that, given a PDF as input, it produces as output an XHTML version of the document that includes bounding box information along with a CSV file that indicates which words belong to which sentences.

To give users the ability to manually edit sentence labels, we helped the PDFNLT team develop the web interface shown in Figure 3. In essence the left hand pane corresponds to textual view of Figure 1 (and is well suited to relabelling a selected sentence by overwriting the relevant text box) while the right hand pane corresponds to graphical view of Figure 2 (and is well suited to visualising the context of a selected sentence through highlighting).

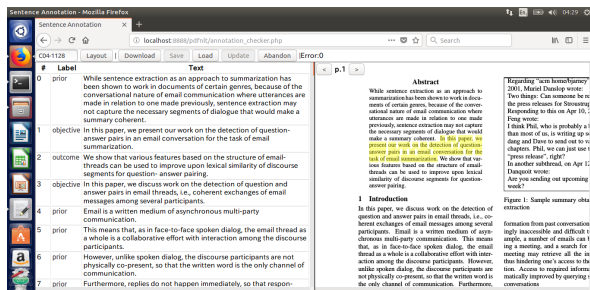


Figure 3: Developmental web interface that combines the content-based and presentation-based approaches.

3. Results

Although our visualisation tool can be used with any pre-trained sentence classifier, we also wanted to test the effect of retraining Sapienta’s CRF model. So we manually annotated 5,300 sentences from 27 papers from our ant corpus. We asked three domain experts to verify the annotations of one document; and, in this small sample, we found an inter-annotator agreement above 90%. We then used 19 papers for training and 8 for testing. These tests showed that re-training significantly improved accuracy and further gains were also achieved by debugging some of Sapienta’s XML conversion and sentence splitting code. The domain experts stated they preferred to work with highlighted PDFs rather than an XML format as they liked to see

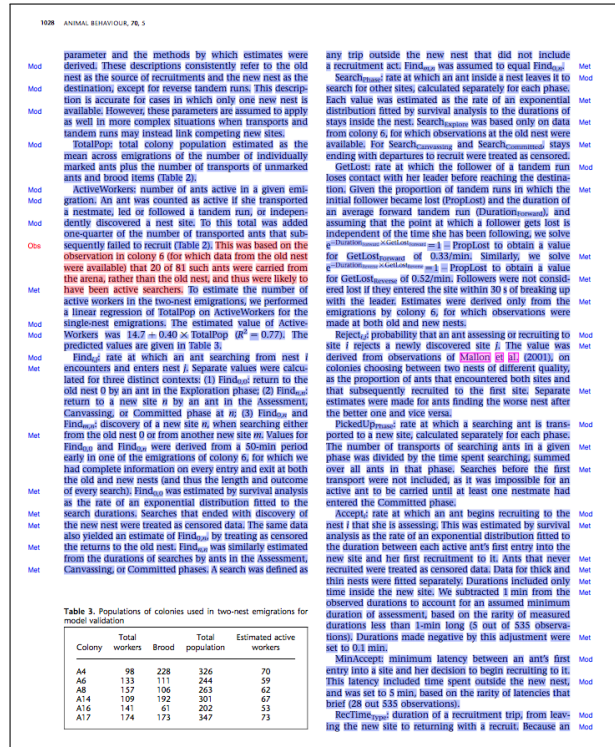


Figure 4: Example highlighting an important observation buried within a large block of text about the methodology.

each sentence in the context of the surrounding textual and graphical cues. As shown in Figure 4 above, the highlighting even led to the discovery of an important observation within a whole page of methodology that would not have been so easily found without our tool. The observation described a seemingly commonplace behaviour (social carries initiated outside the nest) of whose existence our biological collaborators were unfamiliar.

4. Acknowledgements

We thank Goran Topic and Akiko Aizawa for the development of PDFNLT; we thank Maria Liakata and James Ravenscroft for help with the deployment of Sapienta; we thank Nigel Franks, Ana Sendova-Franks and Alan Worley for their domain expert feedback; and we thank Simone Teufel and the anonymous reviewers for useful feedback.

5. Bibliographical References

Hulkkonen, S. (2017). Automatically classifying the content of research publications. Master’s thesis, Univ. of Bristol, UK.

Liakata, M. et al. (2012). Automatic recognition of conceptualization zones in scientific articles and two life science applications. *Bioinformatics*, 28(7):991–1000.

Stenetorp, P. et al. (2012). Brat: A web-based tool for nlp-assisted text annotation. In *Proc. of the Demonstrations at the 13th Conf. of the Europ. Chap. of the Assoc. for Comp. Linguistics*, EACL12, pages 102–107. Assoc. for Comp. Linguistics.

Teufel, S. (1999). *Argumentative Zoning: Information Extraction from Scientific Text*. Ph.D. thesis, Univ. of Edinburgh.