Dissertation (Ph.D. Thesis)

# An Incrementally Trainable Statistical Approach to Information Extraction Based on Token Classification and Rich Context Models

Christian Siefkes

Disputationen:
16th February 2007

Primary Supervisor: Prof. Dr. Heinz F. Schweppe

Database and Information Systems Group

Fachbereich Mathematik und Informatik

Freie Universität Berlin

Supervisors:

Prof. Dr. Heinz F. Schweppe
Database and Information Systems Group
Institute for Computer Science
Freie Universität Berlin

Prof. Dr. Bernhard Thalheim
Systems for Information Management
Institute of Computer Science and Applied Mathematics
Christian-Albrechts-Universität zu Kiel

Dedicated to
the memory of my parents,

Uta Siefkes
(1940–2002)

and

Harm Siefkes
(1936–1989)

**Abstract**

Most of the information stored in digital form is hidden in natural language (NL) texts. While *information retrieval* (IR) helps to locate documents which might contain the facts needed, there is no way to answer queries. The purpose of *information extraction* (IE) is to find desired pieces of information in NL texts and store them in a form that is suitable for automatic querying and processing.

The goal of this thesis has been the development and evaluation of a trainable statistical IE approach. This approach introduces new functionality not supported by current IE systems, such as support for *incremental training* to reduce the human training effort by allowing a more interactive workflow.

The IE system introduced in this thesis is designed as a generic framework for statistical classification-based information extraction that allows modifying and exchanging all core components (such as classification algorithm, context representations, tagging strategies) independently of each other. The thesis includes a systematic analysis of switching one such component (the tagging strategies).

Several new sources of information are explored for improving extraction quality. Especially we introduce rich tree-based context representations that combine document structure and generic XML markup with more conventional linguistic and semantic sources of information. Preparing these rich context representations makes it necessary to unify various and partially conflicting sources of information (such as structural markup and linguistic annotations) in XML-style trees. For this purpose, we develop a merging algorithm that can repair nesting errors and related problems in XML-like input.

As the core of the classification-based IE approach, we introduce a generic classification algorithm (Winnow+OSB) that combines online learning with novel feature combination techniques. We show that this algorithm is not only suitable for information extraction, but also for other tasks such as text classification. Among other good results, the classifier was found to be one of the two best filters submitted for the 2005 Spam Filtering Task of the *Text REtrieval Conference (TREC)*.

The thesis includes a detailed evaluation of the resulting IE which shows that the results reached by our system are better than or competitive with those of other state-of-the-art IE systems. The evaluation includes an ablation study that measures the influence of various factors on the overall results and finds that all of them contribute to the good results of our system. It also includes an analysis of the utility of interactive incremental training that confirms that this newly introduced training regimen can be very helpful for reducing the human training effort. The quantitative evaluation is complemented with an analysis of the kinds of mistakes made during extraction and their likely causes that allows a better understanding of where and how we can expect further improvements in information extraction quality to be made and which limits might exist for information extraction systems in general.

*Wir sehen ein kompliziertes Netz von Ähnlichkeiten, die einander übergreifen und kreuzen. Ähnlichkeiten im Großen und Kleinen.*

*Ich kann diese Ähnlichkeiten nicht besser charakterisieren als durch das Wort „Familienähnlichkeiten"; denn so übergreifen und kreuzen sich die verschiedenen Ähnlichkeiten, die zwischen den Gliedern einer Familie bestehen: Wuchs, Gesichtszüge, Augenfarbe, Gang, Temperament, etc. etc. – Und ich werde sagen: die „Spiele" bilden eine Familie. [...]*

*Wie würden wir denn jemandem erklären, was ein Spiel ist? Ich glaube, wir werden ihm* Spiele *beschreiben, und wir könnten der Beschreibung hinzufügen: „das,* und *Ähnliches, nennt man ‚Spiele'". Und wissen wir selbst denn mehr? Können wir etwa nur dem Anderen nicht genau sagen, was ein Spiel ist? – Aber das ist nicht Unwissenheit. Wir kennen die Grenzen nicht, weil keine gezogen sind.*

— *Ludwig Wittgenstein, Philosophische Untersuchungen*

Leeloo: *Hello.*

Korben Dallas: *Oh, so you speak English now.*

Leeloo: *Yes. I learned.*

— *The Fifth Element (1997)*

# Contents

*Contents*

# List of Tables

# List of Figures