

Informatik-Fachberichte 243

Herausgeber: W. Brauer
im Auftrag der Gesellschaft für Informatik (GI)

Subreihe Künstliche Intelligenz

Mitherausgeber: C. Freksa
in Zusammenarbeit mit dem Fachbereich 1
„Künstliche Intelligenz“ der GI

Udo Hahn

Lexikalisch verteiltes Text-Parsing

Eine objektorientierte Spezifikation
eines Wortexpertensystems
auf der Grundlage des Aktorenmodells



Springer-Verlag

Berlin Heidelberg New York London
Paris Tokyo Hong Kong Barcelona

Autor

Udo Hahn
Universität Passau
Fakultät für Mathematik und Informatik
Postfach 2540, W-8390 Passau

CR Subject Classification (1987): I.2.7, I.2.4, H.3.1

CIP-Titelaufnahme der Deutschen Bibliothek.

Hahn, Udo:

Lexikalisch verteiltes Text-Parsing: eine objektorientierte Spezifikation eines Wortexperten-
systems auf der Grundlage des Aktorenmodells / Udo Hahn. – Berlin; Heidelberg; New York;
London; Paris; Tokyo; Hong Kong; Barcelona: Springer, 1990

(Informatik-Fachberichte; 243: Subreihe künstliche Intelligenz)

Zugl.: Konstanz, Univ., Diss., 1987

ISBN-13: 978-3-540-53230-9

e-ISBN-13: 978-3-642-76132-4

DOI: 10.1007/978-3-642-76132-4

NE: GT

Dieses Werk ist urheberrechtlich geschützt. Die dadurch begründeten Rechte, insbesondere die der Übersetzung, des Nachdrucks, des Vortrags, der Entnahme von Abbildungen und Tabellen, der Funksendung, der Mikroverfilmung oder der Vervielfältigung auf anderen Wegen und der Speicherung in Datenverarbeitungsanlagen, bleiben, bei auch nur auszugsweiser Verwertung, vorbehalten. Eine Vervielfältigung dieses Werkes oder von Teilen dieses Werkes ist auch im Einzelfall nur in den Grenzen der gesetzlichen Bestimmungen des Urheberrechtsgesetzes der Bundesrepublik Deutschland vom 9. September 1965 in der jeweils geltenden Fassung zulässig. Sie ist grundsätzlich vergütungspflichtig. Zuwiderhandlungen unterliegen den Strafbestimmungen des Urheberrechtsgesetzes.

© Springer-Verlag Berlin Heidelberg 1990

2145/3140-543210 – Gedruckt auf säurefreiem Papier

V o r w o r t

Bei der vorliegenden Arbeit handelt es sich um eine überarbeitete und aktualisierte Version meiner Dissertation, die im Februar 1987 von der Sozialwissenschaftlichen Fakultät der Universität Konstanz angenommen wurde. Zusammen mit der ebenfalls in den Informatik-Fachberichten (Bd. 198) veröffentlichten Dissertation von Ulrich Reimer dokumentiert die Arbeit den methodischen Entwicklungsstand der an Textverstehensproblemen orientierten informatischen Komponente des Faches Informationswissenschaft, in dessen institutionellem Rahmen beide Monographien entstanden sind. Der engere Projektkontext wurde durch das vom BMFT von 1981 bis 1986 geförderte Textkondensierungssystem TOPIC definiert, ein System zur flexiblen Zusammenfassung (Abstracting) deutschsprachiger Fachtexte aus dem Bereich der Informationstechnik.

Ausgangspunkt dieser Arbeit war das Problem, ein authentischen Fachtexten angemessenes linguistisches Beschreibungssystem für sprachliche Phänomene und ein Verfahren für ihre inhaltliche Analyse zu entwickeln. Zur Lösung dieses Problems ist eine Textgrammatik und ein entsprechender Text-Parser erarbeitet worden, die auf einem semantischen Beschreibungsansatz beruhen.

Eine Betrachtung der grundlegenden Merkmale semantischer Grammatiken zeigt, daß sie im wesentlichen Lexikon-Grammatiken sind; d.h. das gesamte für das Parsing relevante Wissen ist den lexikalischen Einheiten der Sprache zugeordnet. Zugleich sind semantische Grammatiken/Parser durch eine stark auf Kommunikation ausgerichtete Architektur gekennzeichnet - damit ist einerseits der parallele Zugriff auf unterschiedliche Wissensquellen beim Textverstehen gemeint, andererseits die aus der lexikalischen Verteilung erwachsende Notwendigkeit zu intensiven Interaktionen zwischen lokalen Parse-Prozessen bei der Interpretation zusammenhängender Strukturen auf der Phrasen-, Satz- und Textebene.

Im Bereich semantischer Grammatiken/Parser stellt das Wortexperten-Modell den konzeptionell innovativsten Versuch dar, die oben angesprochenen Merkmale in ein Sprachverstehensmodell zu integrieren. Gravierende methodische Mängel und nicht haltbare linguistische Hypothesen verlangten jedoch eine Reformulierung der ursprünglichen Wortexperten-Konzeption. Die wesentlichen in dieser Arbeit entwickelten Vorschläge können wie folgt zusammengefaßt werden:

Das hier vorgestellte Konzept eines lexikalisch verteilten Text-Parsers ist vollständig in einer objektorientierten Spezifikationssprache auf der Grundlage des Aktorenmodells beschrieben. Darauf aufbauend

- a) werden die deklarative Struktur der Textgrammatik und das prozedurale Verhalten des Text-Parsers einheitlich und auf einem gleichen Abstraktionsniveau behandelt,
- b) sind die Sprachprimitive der Aktoren-Spezifikationssprache im Sinne des Konzepts abstrakter Datentypen formalisiert,
- c) ist insbesondere die linguistische Spezifikation der semantischen Grammatik von der Spezifikation der Repräsentationsstrukturen des Welt- und Domänenwissens logisch strikt getrennt.

Im Mittelpunkt der linguistischen Spezifikationen stehen Beschreibungen generalisierbaren Verhaltens lexikalischer Objekte:

- a) Für das Verhalten lexikalischer Klassen repräsentative Wortexperten-Prototypen werden nach funktionalen Gesichtspunkten formuliert und lösen das alte Wortexperten-Paradigma von der Beschreibung des idiosynkratischen Verhaltens individueller Wörter ab (funktionale Wortexperten-Prototypen kann man sich etwa als hochgradig reguläre Anaphora-, Ellipsen- oder Koordinationsgrammatiken vorstellen).
- b) Prototypen setzen sich ihrerseits aus Subexperten zusammen. Ihre Spezifikation trägt zur Faktorisierung gemeinsamen Wissens in lexikalisch verteilten Grammatiken bei und reduziert erheblich die Beschreibungsredundanz in und das Wachstum von Wortexperten-Kollektionen.

Zur Bewältigung des grammatikalischen Komplexitätsproblems, dem sich alle Entwickler von natürlichsprachlichen Grammatiken mit dem Anspruch möglichst breiter linguistischer und konzeptueller Abdeckung gegenüber sehen, sind an die Körnung der Grammatikspezifikation bewußt solche Anforderungen gestellt worden, die auf einen partiellen semantischen Text-Parser für das Deutsche zielten. Das Beschreibungsproblem konnte somit darauf reduziert werden, die textuellen Bezüge von durch Nomen denotierten Konzepten auf einem "flachen" Verstehensniveau zu erfassen, wie es für die Erfüllung der Funktion des Textverstehenssystems TOPIC (Bereitstellung indikativer, d.h. themenbezogener Textzusammenfassungen), in das der Parser integriert ist, hinreichend ist.

Auf der Grundlage dieser methodisch motivierten Fortschritte läßt sich die hier gegebene Beschreibung als ein Beitrag zu einer formalisierten Theorie des semantischen Text-Parsing charakterisieren. Art und Konzeption der Beschreibungssprache garantieren zudem einen hohen Grad an Transportabilität in andere Domänen, was für semantische Grammatiken (wegen fehlender Abstraktionsmechanismen) bislang eher ungewöhnlich ist. Insgesamt stehen im Rahmen des hier entwickelten Beschreibungsansatzes konzeptuell hohe Mechanismen zur Verfügung, mit denen der Entwurf und die Wartung semantischer Grammatiken/Parser unter strenger planerischer Kontrolle bleibt.

Die Gültigkeit der textlinguistischen Beschreibungen beruht auf den Erfahrungen, die mit der implementierten Version dieser Textgrammatik im Rahmen einer insgesamt fünfjährigen Phase von Textkondensierungsexperimenten mit dem Textverstehenssystem TOPIC gesammelt wurden. Neben der Formalisierung der Textbezüge ist die hier beschriebene Textgrammatik also - wenn auch nicht unter Zugrundelegung strenger experimenteller Maßstäbe - in ihren Grundzügen validiert.

Wie immer bei längerfristigen Entwicklungsarbeiten in größeren Teams gibt es eine Reihe von Personen und sachlichen Umständen, die auf individuelle Arbeiten bestimmend einwirken. In meinem Fall gilt dies zu allererst für Prof. R. Kuhlen, der den institutionellen Rahmen für informationswissenschaftlich definierte Forschung aufgebaut und die wesentlichen Forschungsfragen thematisiert hat, die durch das TOPIC-Projekt konstruktiv ausgearbeitet wurden. Mit ausschlaggebend für die dabei erzielten Fortschritte war die Gewährung und Absicherung wissenschaftlicher Freiräume, in denen die TOPIC-Gruppe im Rahmen der Projektdefinition auf eine auch heute noch eher unüblich selbstbestimmte Weise forschen konnte. Nicht zuletzt diesen Rahmenbedingungen verdanken die beteiligten Wissenschaftler das durch eigene Publikationen erreichte Ausmaß an individueller Profilierung. Eine naturgemäß andere Form der Kooperation entwickelte sich bei der Zusammenarbeit mit meinem unmittelbaren Projektpartner Ulrich Reimer. Wir haben gemeinsam

die konzeptionellen und technischen Rahmenbedingungen des TOPIC-Systems entwickelt und waren folglich auch die "ersten" Diskussionspartner für fachliche Probleme im Zusammenhang mit unseren Dissertationen. Ulrich Thiel, Mitarbeiter in der TOPOGRAPHIC-Gruppe am Lehrstuhl, hat mich besonders in der Schlußphase der Fertigstellung dieser Arbeit durch viele ins Detail gehende Gespräche nachhaltig unterstützt. Prof. E. Pause, dem Korreferenten der Dissertation, verdanke ich wertvolle Hinweise auf die linguistische Signifikanz der hier berichteten Ergebnisse. Prof. M. Jarke ermöglichte mir schließlich die Fertigstellung der jetzt vorliegenden Version meiner ursprünglichen Dissertation in der neuen Passauer Arbeitsumgebung. Bei allen möchte ich mich aufrichtig für ihre Unterstützung bedanken. Dies gilt auch für die wissenschaftlichen Gesprächspartner auf Kongressen und Arbeitstreffen, die mein Verständnis des Wortexperten-Konzepts durch ihre Kritik und Anregungen herausgefordert und damit auch fortlaufend verschärft haben.

Passau, im Februar 1990

Udo Hahn

Inhaltsverzeichnis

1 Einführung	1
1.1 Beschreibung von Volltexten durch Textgrammatiken	1
1.2 Adäquatheitsbetrachtungen zur Wahl eines geeigneten Parsing-Ansatzes für die Analyse von Volltexten	4
1.2.1 Von syntaktischen zu semantischen Grammatiken für die Beschreibung natürlicher Sprachen	4
1.2.1.1 Parallele Modelle des Sprachverstehens	9
1.2.1.2 Lexikalisierung von Grammatiken	13
1.2.2 Von Satzgrammatiken zu Textgrammatiken	16
1.3 Prinzipien des lexikalisch verteilten Parsing	18
1.4 Beiträge dieser Arbeit zur Methodik und Empirie des lexikalisch verteilten Parsing	22
1.5 Anpassung der linguistischen Beschreibung an die Erfordernisse informationeller Sprachanalyse: partielles Text-Parsing für indikative Textzusammenfassungen	25
2 Grundlagen der Spezifikation eines lexikalisch verteilten Text-Parsers mit einer Aktorensprache	28
2.1 Methodologische Gesichtspunkte bei der Auswahl des Beschreibungsmodells für lexikalisch verteilte Grammatiken	28
2.2 Grundlagen des Aktorenmodells	32
2.3 Eine objektorientierte Spezifikationssprache zur Modellierung verteilter Systeme	36
3 Strukturelle Beschreibung der Wissensquellen des Text-Parsers	54
3.1 Das Frame-Repräsentationsmodell	55
3.1.1 Grundlegende Konstrukte des Frame-Repräsentationsmodells	56
3.1.2 Operationen im Frame-Repräsentationsmodell	61
3.1.2.1 Anfrage-Operationen	62
3.1.2.2 Änderungsoperationen	69
3.2 Die Wortexperten	77
3.3 Das Parser-Bulletin	82
3.3.1 Grundlegende formale Strukturen in Texten und Texttoken	82
3.3.2 Formale Definition des Parser-Bulletins	87
3.3.3 Spezifikation von Operationen auf dem Parser-Bulletin	89
3.3.3.1 Anfrage-Operationen	89
3.3.3.2 Änderungsoperationen	110
4 Beschreibung des Text-Parsers	112
4.1 Empirische Abdeckung der lexikalisch verteilten Textgrammatik	112
4.1.1 Die lexikalische Beschreibungsebene	112
4.1.2 Die satzbezogene Beschreibungsebene	115
4.1.3 Die textuelle Beschreibungsebene	117

4.2 Formale Spezifikation eines lexikalisch verteilten Text-Parsers	123
4.2.1 Administrative Aktoren	125
4.2.2 Semantischer Kern der lexikalisch verteilten Textgrammatik:	
Kohäsionsaktoren für lokale Verkettungsprozesse in Texten	139
4.2.2.1 Nominale Anaphora	145
4.2.2.2 Nominale lexikalische Korrespondenz	154
4.2.2.3 Nominalkomposita	178
4.2.2.4 Adjektivische lexikalische Korrespondenz	189
4.2.2.5 Texttransformation (Fokus-Berechnung)	198
4.2.3 Pragmatischer Kern der lexikalisch verteilten Textgrammatik:	
Kohärenzaktoren für globale Vertextungsmuster	199
4.2.4 Operative Aktoren: Spezifikation einzelner Lesarten	213
4.2.4.1 Kohäsionslesarten	214
4.2.4.2 Kohärenzlesarten	219
5 Fazit	220
5.1 Die Grundzüge des Modells des lexikalisch verteilten Parsing	220
5.2 Ausblick auf mögliche Erweiterungen des lexikalisch verteilten Text-Parsers	222
Appendix-1: Kontrakt für die Datenstruktur BULLETIN	224
Appendix-2: Fragment einer Frame-Wissensbasis	228
Appendix-3: Beispiel-Text und Beispiel-Parse	236
Appendix-4: Index zum formalen Apparat	241
Literaturverzeichnis	246