

## ZUR PROBLEMATIK LANG- UND KURZFRISTIGER PROJEKTE IM BEREICH DER COMPUTERLINGUISTIK

Manfred Thiel

Sonderforschungsbereich 100 /A2  
Universität des Saarlandes  
6600 Saarbrücken

Phänomene wie Theorien oder Systeme besitzen eine gewisse Lebensdauer und manche weisen auch in diesen Zeitgrenzen eine zyklische Entwicklung auf. Nicht nur auf dem Gebiet der LDV sind diese Zyklen eng mit dem Begriff des Redesign verbunden. Im folgenden soll auf die Problematik hingewiesen werden, die bei einer phasenverschobenen zyklischen Entwicklung von linguistischen Theorien, computerlinguistischen und softwaretechnologischen Systemkonzeptionen entstehen kann. Diese Überlegungen sollen also versuchen, die Probleme aufzuzeigen; es wird nicht die Absicht verfolgt, Lösungen anzubieten, höchstens eine: Erhaltung der Pluralität der Verhaltensalternativen von Projekten. Entsprechend abrupt endet der Beitrag.

Hat ein computerlinguistisches Projekt Zeit genug, um bestimmten Fragen auf den Grund zu gehen, entsteht zwangsläufig das Problem, daß es eigentlich alle 2-5 Jahre die linguistischen und/oder informatischen Grundlagen neu überdenken und Systeme neu konzipieren müßte. Tut man dies in einem langfristigen Projekt nicht, so ist es nach einem derartigen Zyklus

**einerseits** selten in der vordersten Reihe der Forschung, wenn es darum geht, z.B. neue linguistische oder informatische Konzepte zu realisieren (dies bedeutet nicht: nicht zu perzipieren!).

**andererseits** aber ist es ziemlich einfach, innerhalb von 2 oder 4 Jahren zu zeigen, daß z.B. eine linguistische Theorie in der Computerlinguistik anwendbar ist und zumindest einen Teil der Sprache beschreibt. Es bleiben aber immer Fälle offen, die man in dem relativ kurzen Zeitraum nicht erreicht. Z.B.: wieviele Jahre muß man an einem Parser arbeiten, bis er stabil und linguistisch derart mächtig ist, daß man damit einigermaßen sicher juristische Texte analysieren kann, in denen Sätze vorkommen, die länger als 100 Wörter sind? Wie verhält sich eine linguistische Theorie, wenn es darum geht, sog. "Ausnahmen" analysieren zu können, wobei "Ausnahmen" immer in Bezug auf das implementierte linguistische Modell zu verstehen ist? So können Regeln nicht bis zu einer filigranten Exaktheit im Detail restringiert werden, ohne daß ihre Allgemeingültigkeit verloren geht. Es hat sich außerdem gezeigt, daß die letzten Prozente der zu lösenden Probleme einen immer höheren Aufwand an Zeit und Arbeitseinsatz benötigen.

Daraus lassen sich folgende Thesen ableiten:

1. In unregelmäßigen Abständen entstehen neue theoretische Konzepte.
2. Ein System sollte nach einer gewissen Zeit einem Redesign oder einer Neukonzeption unterworfen werden, um den Anschluß an neue theoretische Konzepte nicht zu verlieren.
3. Nur wenn man den für ein System gewählten theoretischen Ansatz lange und intensiv genug verfolgt, können grundlegende Erkenntnisse, wie sie oben angedeutet wurden, gewonnen werden.

Der Widerspruch zwischen der 2. und 3. These könnte so mancher linguistischen und informatischen Konzeption Schwierigkeiten bereiten. Bei der Konzeption eines langfristigen Projektes muß also großer Wert auf die Planung der projektinternen Lebensdauer der gewählten theoretischen Grundlagen gelegt werden. Es ist dabei davon auszugehen, daß viele (d.h. thematisch anspruchsvolle) computerlinguistische Projekte eine verhältnismäßig langfristige Angelegenheit sein **müssen**. Andererseits darf man aber nicht vergessen, daß der Aufwand, die zu einem beliebigen Zeitpunkt eines Projektes übriggebliebenen Fragen zu lösen, überproportional steigt. So stellt sich für eine linguistische Theorie in einem sprachverarbeitenden System irgendwann das Problem, daß die Eleganz ihrer Konstruktion in Gefahr gerät, von Notlösungen und angebauten Balkönchen überwuchert zu werden. Lohnt es sich dann noch, weiter zu machen, oder soll man sich einer neuen Theorie zuwenden? Geschieht mit dieser dann das gleiche? Der Widerspruch liegt auch darin, daß diese Balkönchen gebaut bzw. ihre Ursachen erst einmal entdeckt werden müssen, um die Grenzen einer Konzeption abstecken und neue Ansätze aus diesen Erfahrungen aufbauen zu können. Ein sehr einfaches, aber doch typisches Beispiel für das oft beträchtliche und manchmal auch unübersehbare Wachstum eines Systems auch für den Fall, daß nur ein triviales Problem gelöst werden soll, ist folgendes: Zur Unterstützung der Wörterbuchkodierung wird eine Regel benutzt, die folgenden Inhalt hat:

*wenn ein Substantiv auf das Suffix 'ung' endet, dann ist sein Genus feminin.*

Sofort fallen uns die problematischen Fälle ein: "Dung", "Sprung" usw., also Beispiele, in denen 'ung' nicht Suffix ist. Über die eigentliche linguistische Regel hinaus muß sich der Computerlinguist also noch etwas einfallen lassen. Es bieten sich dabei verschiedene Alternativen an:

1. Die Graphemfolge 'ung' wird gegebenenfalls als Suffix präkodiert. Damit ist aber prinzipiell nichts gewonnen: ob man kodiert: "ist feminin" oder: "ist Suffix" läuft auf das gleiche hinaus.
2. Die Wörter, bei denen 'ung' kein Suffix ist, müssen vorher alle im Wörterbuch stehen. Dies ist eine Notlösung und verursacht ein wenig wissenschaftliche Bauchschmerzen.
3. Eine morphologische Analyse stellt dagegen eine adäquate Vorgehensweise dar.

Allerdings tritt hier der diskutierte Effekt ein: diese morphologische Analyse schafft selbst wieder neue Probleme, die, genau wie das hier diskutierte, zu ihrer Lösung wieder anderen, möglicherweise

ähnlich großen oder noch größeren Aufwand erfordern. Sollte die Computerlinguistik als Vorlage für ein Bild von Escher herangezogen werden können? Dies scheint nicht der Fall zu sein: Es entstehen zwar neue Komponenten, die hochgradig voneinander abhängig sind, aber es ist zu erwarten, daß sich der Kreis irgendwann schließt und keine neuen Komponenten hinzukommen. Allerdings besteht immer die Gefahr, daß die Veränderung einer Komponente aufgrund der Vernetzung Auswirkungen auf andere haben kann.

Folgende Schlußfolgerungen sind daraus zu ziehen:

1. Die meisten Probleme in der Computerlinguistik sind nicht trivial.
2. Computerlinguistische Systeme sind von Natur aus sehr komplex; die Aussagekraft von Systemen, die sich auf einen Ausschnitt

der gesamten Problematik beschränken, ist mit äußerster Vorsicht zu bewerten. Es scheint oft so etwas wie ein "Alles-oder-nichts-Effekt" einzutreten.

3. Es ist notwendig, zur Unterstützung der Entwicklung von sprachverarbeitenden Systemen Werkzeuge zu besitzen.

Steht ein langfristiges Projekt diese Durststrecke des Beharrens auf einem gewählten Ansatz durch, um danach sagen zu können, welche Anforderungen eine linguistische Theorie oder eine Softwarekonzeption erfüllen muß, ist dies aufgrund der breiten Variation an Grammatikalität und Akzeptanz in der natürlichen Sprache zugegebenermaßen recht undankbar. Es soll aber auch nicht die Problematik vergessen werden: Derartige "praktische" Arbeiten müssen irgendwann einmal auch in wissenschaftliche Neuansätze münden.

## Anzeige

# Sprache und Information S+I

## Thomas Herrmann Zur Gestaltung der Mensch-Computer-Interaktion: Systemerklärung als kommunikatives Problem

Ca. 280 Seiten. Kart. ca. DM 108,-. ISBN 3-484-31914-3 (Band 14)

Die Untersuchung orientiert sich an solchen Computerbenutzern, die mit dem Rechner nicht nur Routineaufgaben erledigen, sondern abwechselnden Problemstellungen gegenüberstehen, die einen kreativen Lösungsprozeß erfordern. Der Kenntnisstand solcher Benutzer ändert sich permanent während der Mensch-Computer-Interaktion. Für diese Benutzungsweise werden auch Erklärungen während des Dialogs notwendig. Unter diesem Blickwinkel wird untersucht, welche Leistung durch systemgebundene Selbsterklärung über EDV vermittelt werden kann. Das breite Spektrum möglicher Erklärungen wird hinsichtlich der Erklärungsinhalte sowie der Orientierung und der Form ihrer Darstellung differenziert. Diese Aspekte werden in einem Modell der Mensch-Computer-Interaktion zusammengefaßt, das sich am Computerbild des (Be-)Nutzers orientiert (CombiNo-Modell der M-C-I). Insbesondere enthalten Erklärungen zu EDV-Systemen auch konzeptuelles (im Gegensatz zu ereignis-orientiertem) Wissen, dessen kommunikativ-adäquate Vermittlung vom Kontext der Erklärungssituation relativ unabhängig ist und einen argumentativen Dialog benötigt – es findet sog. kruziale Kommunikation statt. Sie ist nach heutigem Stand der Sprechakt- und Argumentationstheorie nicht im Rahmen der Computerbenutzung technisch zu realisieren, wenn man ein unangemessenes Ausmaß unerkannter Mißverständnisse vermeiden will. »Kruzial« ist eine Eigenschaft, die als Abgrenzungskriterium bei der Gestaltung selbsterklärungsfähiger Systeme dienen kann.

Die Auseinandersetzung wird mit solchen Ansätzen geführt, die davon ausgehen, daß der Mensch kommunikative Erwartungen an die Mensch-Computer-Interaktion heranträgt. Dem wird ein Vorschlag gegenübergestellt, wie man Benutzer unterstützen kann, wenn sie sich bei fehlendem Wissen oder bei Störungen selbst helfen möchten: das Konzept intervenierender Benutzbarkeit.

### Rainer Kuhlen

in Zusammenarbeit mit UDO HAHN und ULRICH REIMER

## Informationslinguistik

Theoretische, experimentelle, curriculare und prognostische Aspekte einer informationswissenschaftlichen Teildisziplin

Ca. 250 Seiten. Kart. ca. DM 74,-. ISBN 3-484-31915-1 (Band 15)

Der Band informiert mit seinen 8 Beiträgen (4 davon sind in Englisch geschrieben) über die Entwicklung der Informationslinguistik, wie sie im Bereich Informationswissenschaft an der Universität Konstanz betrieben wird. Informationslinguistik ist keine neue wissenschaftliche Disziplin, sondern wird als zentrale Teildisziplin der Informationswissenschaft angesehen. Sofern diese experimentell arbeitet, befaßt sie sich weitgehend mit Fragen des Information Retrieval. Entsprechend der allgemeinen Forschungsentwicklung vom Referenz-/Dokument-Retrieval zu intelligenten Systemen hat auch die Informationslinguistik den Übergang von mehr quantitativen, strukturalistischen Verfahren zu wissensbasierten Verfahren der Künstlichen-Intelligenz-Forschung vollzogen.

Der Band enthält die folgenden Beiträge: R. KUHNEN, Informationslinguistik. – R. KUHNEN, Quantitative Probleme beim induktiven Wörterbuchaufbau. – R. KUHNEN, Linguistic aspects of information retrieval. – U. HAHN, Curriculum Informationslinguistik. – R. KUHNEN, Similarities and difference in intellectual and machine text understanding and condensation. – U. HAHN/U. REIMER, TOPIC – A text understanding and condensation system. – R. KUHNEN, Development of information linguistics. Results of an international Delphi poll. – U. HAHN/U. REIMER, »State of the art« der Volltextanalyse und Textkondensierung.

# Niemeyer