

LDV-FORUM

Forum der Gesellschaft für Linguistische Datenverarbeitung (GLDV)

LDV-Forum 14.2 (1997)

Forum der Gesellschaft für Linguistische Datenverarbeitung e. V.

Herausgeber

Prof. Dr. Gerhard Knorz; Gesellschaft für Linguistische Datenverarbeitung e. V.

Anschrift: Fachhochschule Darmstadt, Fachbereich Information und Dokumentation, Haardtring 100, D-64289 Darmstadt; Tel: (06151) 16-8499; Fax: (06151) 16-8980; e-mail: knorz@www.iud.fh-darmstadt.de

Redaktion

Gerhard Knorz

Wissenschaftlicher Beirat

Prof. Dr. W. Hoepfner (hoepfner@unidui.uni-duisburg.de); Prof. Dr. Gerhard Knorz; Prof. Dr. Winfried Lenders (lenders@uni-bonn.de);

Editorial

Mit der Auslieferung des letzten LDV-Forum wurde ein paar Tage zugewartet, um ein entsprechendes Schreiben des Wahlvorstandes dem Heft beizulegen. Aber erst die zweite Ausgabe des LDV-Forum in diesem Jahr kann die *Wahl zu Vorstand und Beirat der GLDV* redaktionell zur Kenntnis nehmen. Zugegeben, die Kandidatenlage zur Besetzung des Vorstandes trug wenig dazu bei, die Wahl zu einem spannenden Ereignis zu machen, und natürlich erwartet niemand ernstlich, daß eine Beiratswahl mit 6 Kandidaten für 4 Positionen die Mitglieder innerlich aufgewühlt zum Stimmzettel greifen läßt. Aber 93 abgegebene Wahlentscheidungen bei insgesamt 281 erfolgreich angeschriebenen Mitgliedern ist doch wirklich etwas dünn. Was soll man daraus schließen? Immerhin weiß man nun, daß die 93 wählenden Mitglieder sich dem gedruckten Wort verpflichtet fühlen: 78 von ihnen, also ca. 85%, sprechen sich für den *Erhalt der Papierversion* des Publikationsorgans der GLDV aus (siehe dazu den entsprechenden Bericht auf Seite 67). Halt, daß ich es nicht vergesse: Dem neugewählten Vorstand alle guten Wünsche für eine erfolgreiche Amtszeit!

Prof. Dr. Ulrich Schmitz
(e-mail: ulrich.schmitz@uni-essen.de)

Erscheinungsweise

2 Hefte im Jahr, halbjährlich zum 31. Mai und 31. Oktober

Bezugsbedingungen

Für Mitglieder der GLDV ist der Bezugspreis des LDV-Forum im Jahresbeitrag mit eingeschlossen.

Jahresabonnements können zum Preis von DM 40,- (incl. Versand), Einzel-exemplare zum Preis von DM 20,- (zuzügl. Versandkosten) bestellt werden:

LDV-Forum, c/o IKS, Poppelsdorfer Allee 47, 53115 Bonn

Fachbeiträge

Unaufgefordert eingesandte Fachbeiträge werden vor Veröffentlichung von mindestens 2 ReferentInnen begutachtet. Manuskripte sollten deshalb möglichst frühzeitig eingereicht werden und bei Annahme zur Veröffentlichung in jedem Fall elektronisch und zusätzlich auf Papier übermittelt werden.

Rubriken

Die namentlich gezeichneten Beiträge geben ausschließlich die Meinung der AutorInnen wieder. Einreichungen sind – wie bei Fachbeiträgen – an den Herausgeber zu übermitteln.

Ein wenig mehr an Schwung wäre der GLDV also sichtlich zu wünschen. Vorstand und Beirat bemühen sich, das ihre zu tun, indem sie eine programmatische Diskussion auf die Tagesordnung gesetzt und diese auch mit ausreichend Zeit und Raum ausgestattet haben. Für die Mitglieder besteht neben vielen anderen Gelegenheiten die Chance, die *Konvens '98* in Bonn zu einem Aktivposten der Gesellschaft zu machen. Der Titel „*Computer, Linguistik und Phonetik zwischen Sprache und Sprechen*“ ist anregend definiert und trifft eine Thematik, die aus wissenschaftlicher und praktischer Sicht vielversprechend ist. Vergessen Sie also nicht, sich die Deadline zur Einreichung von Beiträgen dick im Terminkalender anzustreichen.

Wenn Ihre Arbeitsschwerpunkte aber doch dem KONVENS-Schwerpunkt zu fern liegen, dann ersehen Sie aus dem aktuellen Forum, daß man dort dringend auf Ihren Beitrag wartet! Es gibt auch gegenwärtig wieder/noch enttäuschend wenig „konstruktives Interesse“ der LeserInnen. Um es nochmals ganz deutlich zu machen: Dies ist eine Aufforderung! Wenn die gegenwärtige Seitenschätzung nicht ganz daneben liegt, so dürfte es sich bei der aktuellen Ausgabe um das papiersparendste Heft handeln, das ich seit Beginn des LDV-Forum verantwortet habe.

Ich weiß nicht, wie sehr der optische Eindruck des vorliegenden Heftes es gleich offenbar werden läßt, daß die beim letzten Mal gefundene, scheinbar stabile neue Lösung für die zukünftige LDV-Forum-Produktion durch einen Vorstandsbeschluß durch eine preisgünstigere Alternative ersetzt wurde. Klüger geworden, werde ich mich über die neue Situation erst äußern, wenn erkennbar ist, daß die zukünftige Realität meine Worte von heute nicht unmittelbar Lügen straft.

Redaktionsschluß

Für Heft 14 (1997) 2:
15. August 1997

Druck und Vertrieb

IKS, Poppelsdorfer Allee 47,
53115 Bonn

Herstellung

Kurt Thomas, Universität
Bonn, IKP. e-mail: thomas@
uni-bonn.de

Zu einem knappen Heft paßt wohl ein knappes Editorial. Und diese Kürze paßt überdies auch in die Zeit des Semesteranfangs, in der dieses Vorwort geschrieben wird, obwohl eigentlich gar keine Zeit verfügbar ist. Und diese fehlende Zeit ist auf Leserseite wohl auch der Grund für all das, was ich oben als fehlenden Schwung und mangelndes „konstruktives Interesse“ umschrieben habe. Wenn wir mal Zeit haben, sollten wir uns fragen, ob wir nicht doch etwas grundsätzlich falsch machen? Ich zumindest nehme mir diese Frage mal vor.

Auflage

400 Exemplare

Darmstadt, den 9.10.97

G.K.

Anzeigen

Preisliste und Informationen:
IKS e.V., Poppelsdorfer Allee
47, D-53115 Bonn, Tel.: (0228)
735621, Fax: (0228) 735639,
e-mail: iks@uni-bonn.de

Bankverbindung

IKS e.V.: PGA Köln (BLZ
370 100 50), Konto 385647-
505

GLDV-Anschrift

Prof. Dr. R. Hausser,
Universität Erlangen-Nürnberg,
Abteilung für Computerlin-
guistik, Bismarkstraße 12,
D-91054 Erlangen; e-mail: rrh@
linguistik.uni-erlangen.de

Affixerkennung in deutschen Wortformen

Ein nicht-lexikalisches Segmentierungs- verfahren nach N. D. Andreev

Oliver Cromm
Universität Göttingen
e-mail: ocromm@gwdg.de

24. März 1997

A non-lexical statistical method for affix recognition within a corpus — developed by N. D. Andreev ([1], [2]) in the early 60's — is applied to German texts. The algorithm was designed to identify inflectional as well as derivational paradigms in any language, using both statistical overrepresentation of letter sequences and combinatorial systematicity as a measure for the reliability of segmentation.

By testing this method on smaller samples, the influence of several parameters is investigated and the most suitable values are selected. The resulting algorithm is then applied to a bigger corpus, the ensuing results are analysed both quantitatively and qualitatively. These display quite good recall rates but reveal some problems due to certain characteristics of German.

1 Überblick

1.1 Grundlegender Ansatz

Das Verfahren von Andreev sucht Affixe am Beginn und Ende von Wörtern. In seiner reinsten Form, wie sie hier untersucht wird, benutzt es keinerlei syntaktische oder semantische Information, sondern berücksichtigt ausschließlich Abweichungen von der Gleichverteilung der Buchstaben, nach Harris eine der wichtigsten Grundlage einer Sprachtheorie [5]. Der Erkennungsprozeß beginnt damit, nach überrepräsentierten Buchstaben im Anfangs- und Endbereich der längeren Wörter eines Korpus zu suchen, im Deutschen z. B. *e* an vorletzter Wortposition. Deren Nachbarbuchstaben werden daraufhin untersucht, ob gewisse Kombinationen auffällig häufiger auftreten als im Falle ihrer Unabhängigkeit zu erwarten, im Beispiel etwa die Endungen *-en*, *-er*, aber auch *-ten*. Diese werden zu mutmaßlichen Affixen erklärt, sofern sie gewisse formale und statistische Bedingungen erfüllen. Diese Affixe wiederum werden, wenn möglich, zu Paradigmen von Stämmen und Affixen erweitert. Nur wenn ein Paradigma gefunden wird, werden die Morphemgrenzen akzeptiert.

Es werden Kriterien gegeben, um Muster agglutinierenden Typs zu erkennen (wobei dann auch nach weiteren Affixen am selben Wort gesucht wird) und zwischen Flexion und Derivation zu unterscheiden.

Wie alle Methoden, die auf der Kookkurrenz benachbarter Buchstaben beruhen, steht das Andreev-Verfahren in der Tradition von Harris' Ansatz zur Morphemgrenzenerkennung [4].

1.2 Einschränkungen dieses Ansatzes

Die Methode setzt Input mit gegebenen Wortgrenzen voraus. Um Sprachen ohne Markierung von Wortgrenzen oder gesprochenen Input zu behandeln, muß sie modifiziert oder mit einem eigenen Algorithmus zur Wortgrenzenerkennung kombiniert werden. Die Identifikation beschränkt sich auf identische Muster nur an Wortgrenzen. Daraus resultieren Schwierigkeiten beim Erkennen von Infixen,

mehrfachen gegenseitig abhängigen Affixen, Reduplikation, nicht-zusammenhängenden Morphemen wie z. B. bestimmten Vokalfolgen in den semitischen Sprachen oder nicht-konkatenativer Morphologie wie bei Umlautphänomenen.

Das Verfahren funktioniert am besten beim Auffinden von *item-and-arrangement*-Strukturen, wie sie für agglutinierende Konstruktionen besonders typisch sind.

2 Mögliche Anwendung

Die Arbeit des Algorithmus beschränkt sich auf die Erkennung von Flexions- und Derivationsaffixen. Darüber hinaus sind, wie bei jeder statistischen Methode, die Recall- und Precision-Raten begrenzt; daher ist das System nicht direkt verwendbar, um alle Morphemgrenzen eines bestimmten Typs aufzufinden.

Andererseits könnte es Hinweise darauf geben, ob bestimmte Modelle der menschlichen Sprachverarbeitung und des Spracherwerbs realistisch sind (vgl. etwa [3], [6]), und es könnte bei der Segmentierung und dadurch bei der Dechiffrierung einer unbekanntem Sprache Anwendung finden.

3 Das Verfahren im einzelnen

Unter Auslassung technischer Details funktioniert der Algorithmus folgendermaßen:

1. Zunächst wird fortlaufender Text in eine Liste von Wörtern und ihren Häufigkeiten transformiert.
2. Nur Wörter ab einer bestimmten Mindestlänge werden als möglicherweise affixbehaftet in Betracht gezogen.
3. Für eine gewisse Anzahl initialer und finaler Wortpositionen werden die höherfrequenten Buchstaben untersucht, beginnend bei den am stärksten überrepräsentierten.
4. Diese Buchstaben werden sukzessive mit geeigneten Nachbarn verkettet, bis die Kette den Wortrand und die erwartete minimale Affixlänge erreicht.

5. Das „Affix“ wird probeweise abgeschnitten, es wird nach Kombinationen der resultierenden „Stämme“ mit anderen „Affixen“ gesucht.
6. Die letzten Schritte werden wiederholt, bis keine weiteren Affixe akzeptiert werden können.
7. Schließlich wird getestet, ob die Morphemgrenze nach links oder rechts verschoben werden sollte und ob es sich um ein Paradigma agglutinierenden Typs handelt.

Bei jedem dieser Schritte müssen gewisse statistische Kriterien erfüllt sein, damit die Arbeit des Algorithmus fortgesetzt wird.

Wenn ein Paradigma, ein *morphologischer Typ* gefunden wurde, werden seine Mitglieder aus der Liste aller Wörter entfernt (im agglutinierenden Fall nur die affixbehafteten Formen).

In der Originalarbeit von Andreev wird die weitere Bearbeitung des Textes nur mit den Nachbarn der vorher erkannten Wörter fortgesetzt, das heißt, syntaktische Information wird eingebracht. In der vorliegenden Arbeit beginnt die nächste Runde wieder mit der reinen, nur um die bisher analysierten Formen verminderten Wortliste.

Man kann, wenn alle morphologischen Typen gefunden sind, diejenigen mit ähnlichen (syntaktischen) Eigenschaften zu Wortklassen, wie Substantiven, Adjektiven und Verben, zusammenfassen. Da hier keine syntaktische Information benutzt wird, kommt eine Zusammenfassung nicht in Betracht.

4 Einstellung der Parameterwerte

Die Methode benutzt zahlreiche Parameter, deren Werte in den vorliegenden Arbeiten von Andreev nicht näher motiviert werden. Daher wurden Testläufe mit unterschiedlichen Parameter-Einstellungen auf Beispieltexen aus der Bibel und aus Computerforen durchgeführt. Die Bibelpassagen lieferten weit bessere Resultate, da sie sprachlich einheitlicher sind und viel weniger Tippfehler und ähnliche Quellen von Rauschen enthalten.

Der Einfluß von Parametern wurde untersucht, darunter

- die Anzahl der untersuchten Wortpositionen,
- die minimale Häufigkeit von Buchstaben, die in Betracht gezogen werden,
- deren minimaler Grad der Überrepräsentation,
- die minimale Rate der Kookkurrenz benachbarter Buchstaben, damit sie als Teil eines Affixes betrachtet werden.

Es stellte sich heraus, daß einige Parameter weit großzügiger gewählt werden können als von Andreev vorgeschlagen, d. h. so, daß wesentlich mehr Buchstaben und -kombinationen in Betracht gezogen werden. Dies führt zu größerem Recall ohne großen Verlust von Precision. Daß Andreev die Anforderungen eher hoch wählte, mag darin begründet liegen, daß zu seiner Zeit die Computer-Ressourcen sehr beschränkt waren, und die großzügigere Wahl der Parameter zu weit höherem Rechenaufwand führt. Tatsächlich wurden viele Ergebnisse in Andreevs Gruppe von Hand ermittelt. Im Gegensatz dazu wurden die hier vorgestellten Ergebnisse weitgehend automatisch auf einem kleinen Rechner erreicht.

Darüber hinaus dienen viele dieser Parameter dazu, mit einer hohen Wahrscheinlichkeit sicherzustellen, daß die in der Stichprobe gefundenen Abweichungen von der Gleichverteilung nicht zufällig sind, sondern Ausdruck einer Gesetzmäßigkeit. Diese Wahrscheinlichkeit steigt entsprechend dem statistischen Gesetz der großen Zahl bei gleicher Parametereinstellung auch mit der Stichprobengröße, daher kann man in größeren Stichproben wiederum geringere Abweichungen bereits als signifikant betrachten.

5 Ergebnisse in der Praxis

Die Methode mit den ermittelten optimalen Parametereinstellungen wurde auf den gesamten Text der Bibel angewendet. Die Bibel ist ein relativ freundliches Korpus, da sie ein begrenztes Vokabular mit wenigen Fremd- und sogar Lehnwörtern benutzt. Das Korpus umfaßt insgesamt 4.652.726 Bytes.

Die morphologischen Grenzen, die der Algorithmus ermittelte, wurden mit Morphemgrenzen verglichen, die vom Autor für jeden Wort-Type, ohne Berück-

sichtigung des Kontextes, intellektuell markiert wurden. Im Falle von Homonymie wurde ein Wort als flektiert betrachtet. Diese Vorgehensweise erscheint gerechtfertigt, da einerseits die statistische Methode ebensowenig zwischen homonymen Wörtern unterscheiden kann, andererseits auf diese Weise die Recall-Werte nur schlechter werden können, nicht besser.

Detaillierte quantitative Ergebnisse enthält die Tabelle 1, einen Überblick die Tabelle 2.

Das Ziel war, Flexionsparadigmen zu finden, nicht Derivationen. Daher bezeichnet der Recall den Anteil der intellektuell markierten Flexionsaffixe, die auch vom Algorithmus markiert wurden, die Precision den Anteil der vom Computer ausgegebenen Affixe, die nach menschlichem Urteil Flexionsaffixe darstellen.

Der Recall beträgt, in Types ausgedrückt, 42,4%, die Precision 89,6%. In Tokens gerechnet haben wir einen Recall von 73,8% und eine Precision von 96,4%.

5.2 Auswertung der Ergebnisse

Die statistische Natur des Algorithmus ergibt bessere Recall- und Precision-Werte bei hochfrequenten Wörtern. Dies ist der Grund für den großen Unterschied zwischen den Ergebnissen in Types und Tokens. Man vergleiche dazu bei den Wörtern mit Flexionsaffix die mittlere Häufigkeit der vom Algorithmus richtig identifizierten (25,4) mit derjenigen der nicht identifizierten (5,6), ersichtlich aus Tabelle 1.

Einige Besonderheiten der deutschen Sprache führen zu Fehlern. Diese sollen qualitativ analysiert werden.

1. Einige hochfrequente Stämme sind zu kurz, um vom Algorithmus in Betracht gezogen zu werden. Alleine das Wort *ein-e* ist für 17% aller fälschlich als nicht affixbehaftete Form (also als mutmaßlicher Stamm) erkannten Tokens verantwortlich.
2. Besonders Verben erscheinen oft mit verschiedenen Wortbildungspräfixen. Diese Kombinationen sind so vielfältig, daß sie leicht als Flexion mißinterpretiert werden können. So findet ein Testlauf des Algorithmus als 16. Paradigma die Stämme

<i>Sorte</i>	<i>Types</i>	<i>Tokens</i>
Alle Wörter	24255	708249
zu kurz (bis 3 Bst.)	348	294570
bleiben	23907	413679
als Angeh. v. Typen		
aussortiert	7224	192242
unanalysiert	16674	221437
aussortiert	7224	192242
als affixlos	378	30480
als affixbehaftet	6846	161762
affixlos	378	30480
korrekt	211	20636
eigtl. flektiert	167	9844
affixbehaftet	6846	161762
tats. flektiert	6137	155871
Wortbildung	680	5990
Komposition	27	78
Zufall	2	21
unanalysiert	16674	221437
davon mit Affix	8170	45444
unflektiert	8504	175993

Table 1: *Detaillierte Auswertung für das Korpus Bibel*

<i>Tokens</i>	<i>insges.</i>	<i>richtig erkannt</i>	<i>nicht erkannt</i>	<i>falsch erkannt</i>
flekt.	14474	6137	8337	167
unflekt.	9424	211	8504	709
<i>Tokens</i>	<i>insges.</i>	<i>richtig erkannt</i>	<i>nicht erkannt</i>	<i>falsch erkannt</i>
flekt.	211159	155871	55288	9844
unflekt.	202520	20636	175993	5891

Tabelle 2: Gesamtergebnis für das Korpus Bibel

-fahren, -genommen, -gezogen, -werfen mit den dazugehörigen Affixen *ab-, auf-, aus-, vor-, weg-*, das als Flexion eingestuft wird, und es folgen in kurzer Folge 11 weitere Paradigmen dieses Typs.

Damit nicht genug, werden die so falsch erkannten Wörter dann auch noch aus dem Korpus entfernt, so daß ihre tatsächlichen Flexionsaffixe nicht gefunden werden können.

Bei 95,9% der fälschlich als flektiert identifizierten Wörter wurde ein Wortbildungsaffix fehlinterpretiert. Gar nicht um ein Affix handelte es sich nur bei 0,4% der Types und 0,06% der Tokens.

3. Aus dem eben genannten Grund ist die Unterscheidung von Flexion und Derivation anhand der statistischen Eigenschaften im Deutschen nicht sehr zuverlässig.
4. Im Deutschen werden über alle Wortklassen hinweg homonyme Flexionsuffixe verwendet. 13 Suffixe und 2 Zirkumfixe (die beide eines der 13 Suffixe enthalten) erfüllen unter großer Überschneidung und Synkretismus alle Funktionen der Substantiv-, Adjektiv- und Verbflexion ursprünglich deutscher Wörter. So sind die drei Affixe *-en, -e, -Ø* in Substantiv-, Adjektiv- und

Verbparadigmen gemeinsam anzutreffen. Viele Suffixe bestehen zudem aus hochfrequenten Buchstaben, die auch am Ende von Stämmen häufig sind. Das führt zu einer Vermischung solcher Fälle mit Stämmen, die nur mit unvollständigem Paradigma im Korpus vorkommen, und dadurch zu vereinzelt willkürlichen Morphemgrenzen, deren Zahl im Test allerdings gering bleibt (s. Tabelle 1).

5. Stammumlaut ist ein häufiges Phänomen in deutschen Flexionsparadigmen. Dies führt im besten Fall zu einer Aufspaltung der beteiligten Paradigmen, ansonsten zur Nichterkennung der umgelauteten Formen.

Zieht man diese Tücken in Betracht, so erscheint die Zahl der korrekt identifizierten Wortformen mit Flexionsaffix von rund drei Viertel aller Tokens beeindruckend, wenn auch nicht ausreichend für praktische Anwendungen. Überdies kann man mit wachsender Korpusgröße auch ein weiteres Ansteigen dieser Rate erwarten.

6 Erweiterungsmöglichkeiten

Interessant wäre, das Verfahren auf phonematisch statt graphematisch transkribiertem Text zu testen. Dies könnte eine größere Regelhaftigkeit aufdecken, andererseits könnten aber auch gewisse Verallgemeinerungen verlorengehen, da die deutsche Orthographie teils grammatikalisch motiviert ist.

Die Voraussetzungen des Verfahrens sind sehr restriktiv. Es könnte in Richtung einer flexibleren Mustererkennung erweitert werden, um den Suchbereich auf Infixe und nicht-zusammenhängende Morpheme zu erweitern. Möglicherweise könnten statt buchstäblich übereinstimmenden Mustern ähnliche gesucht werden, obwohl fraglich ist, ob Ähnlichkeit von Morphen angemessen definiert werden kann, zumal einzelsprachunabhängig.

Statt das Verfahren zu ändern, damit es klassischen Definitionen von Morphemgrenzen entspricht, könnte man auch eine empirischere, mehr naturwissenschaftliche Sichtweise von Sprache einnehmen und es selbst als Definition solcher Grenzen ansetzen.

Literatur

- [1] Andreev, Nikolaj D. (ed.): Statistiko-kombinatornoe modelirovanie jazykov, Moskau/Leningrad 1965
- [2] Andreev, Nikolaj D.: Statistiko-kombinatornye metody v teoretičeskom i prikladnom jazykovedenii, Leningrad 1967
- [3] Bybee, Joan L.: Morphology as Lexical Organization. In: Hammond, Michael/Noonan, Michael (ed.): Theoretical Morphology. Approaches in Modern Linguistics, San Diego 1988
- [4] Harris, Zellig: From phoneme to morpheme. *Language* 31 (2), 1955, p. 190—222
- [5] Harris, Zellig: *A Theory of Language and Information*. Oxford 1991
- [6] MacWhinney, B. (ed.): *Mechanisms of Language Acquisition*, Hillsdale, N.J. 1987

Kurze Geschichte eines Linguistik-Servers im Internet

*Elisabeth Cölfen / Ulrich Schmitz
Uni GH Essen, FB 3, 45117 Essen
elisabeth.coelfen@uni-essen.de,
ulrich.schmitz@uni-essen.de*



I. Die Begeisterung

April 1995. Es fing ganz harmlos an. Wir kamen frisch nach Essen, fanden 5000 Germanistik-Studenten vor und blieben trotzdem guten Mutes. An der „Einführung in die Linguistik“ für Erstsemester nahmen Hunderte von Studenten teil (wo anderswo 30 sitzen), auch in einem Hauptseminar können sich schon mal 150 oder 200 Teilnehmer drängeln.

Schlimmer kann's kaum werden. Wo die Menschen schon von allein einander fremd geworden sind, können neue Medien kein Unheil mehr anrichten. Vielleicht im Gegenteil? Auch sollte irgendeine kreative kleine Spielwiese in diesem Alltags-Tohuwabohu vielleicht doch ganz erfrischend wirken. Und schließlich: altherge-

brachte LDV-Kenntnisse und computerlinguistisches Knowhow lassen sich möglicherweise ja auch für sehr unmittelbar praktische Anwendungszwecke modernisieren.

Wir also los: Hier ein bißchen Geld erbettelt, da ein paar Leute von der Hochschule angestachelt, dort ein paar Nächte um die Ohren geschlagen, und flugs – im August 95 – waren wir „auf Sendung“. Zuerst waren wir selbst ganz begeistert (d.h. Hermann Cölfen und die beiden Verfasser), jedenfalls mehr als wir uns gegenseitig zugaben. Das ist schon mal eine gute Voraussetzung. Glücklicherweise gab's bald aber auch noch andere, die das irgendwie ganz toll fanden, was wir da zurechtbastelten, oder die uns zumindest wohlwollend zuschauten. Monat für Monat fieberten wir den Einschaltquoten entgegen: schon wieder 20, 30 mehr? Und so ging das denn weiter, andere schrieben und arbeiteten mit, Leser-mails kamen, erste Beschwerden (immer ein gutes Zeichen: man wird ernst genommen), und bald machte die weise alma mater einen Batzen Geld locker für einen richtig tollen high-tech-Server ganz für uns alleine. Nun konnte es so richtig los gehen. Bald wuchs uns die Sache über den Kopf, dann wuchsen wir wieder nach, gerüchteweise erfuhren wir, daß man uns draußen in der Welt für ein Riesen-Team mit einer Anzahl Vollzeitstellen hielt, und so etwas spornt natürlich an. Kritische Stimmen gehen ein. Was davon muß man ernst nehmen, was ist unbegründet? Wo müssen wir besser werden? Wie gehen wir mit den Grenzen unserer Arbeitszeit um? (Eigentlich haben wir ja doch ganz andere Aufgaben.) Unsere wilde Arbeitsweise (bürokratiefreie Zone in einem überorganisierten Apparat) widerspricht allen modernen Lehren von Management und Kommunikation. Und es klappt. Allein in den 12 Monaten von September 1996 bis August 1997 hat LINSE, der Linguistik-Server Essen, gut 1,4 Gigabyte an Daten transportiert. Lokale, regionale und internationale Kontakte sind geknüpft worden. Studenten haben eine Menge gelernt über Linguistik, Redaktion, Computer, Projektarbeit. Erstaunlich wenig ist schiefgegangen (gemessen an unseren Möglichkeiten und Erwartungen eigentlich kaum etwas), vieles glücklich gelaufen. Die LINSE ist ein Arbeitsinstrument für Linguisten. Sie hilft anderen, aber sie dokumentiert auch unsere eigene Arbeit. Ihre eher anarchische Geschichte steht ihr ins Gesicht geschrieben. Angebot, Benutzungsmöglichkeiten und tatsächliche Benutzung wachsen und wuchern. Es gibt hohe Motivation, klare Qualitätsstandards, eine strikte Redaktion und eine ästhetische Linie, aber keine systematische Planung. Viel Bewegung und wenig System. Die LINSE ist, wie das ganze Internet, ein Rhizom und kein Fertighaus. Ein Projekt und kein Ergebnis; ein Weg, aber kein Ziel.

Das schafft auch Probleme und ist nicht unbedingt jedermanns Sache. Aber unsere und offenbar auch die von vielen, die wir vorher nicht kannten. Eine Menge hat sie in Gang gesetzt, bei uns und bei anderen. Die vielen einzelnen Schritte, Lehrgelder und Seiteneffekte zu dokumentieren wäre belanglos. Zur kurzen Geschichte gehört, daß junge Studenten ihre Arbeiten erstmals zur (sogar weltweit) öffentlichen Diskussion stellen konnten und sich mächtig dafür ins Zeug legten, sie dann auch entsprechend gut zu machen. Es gehört dazu, daß die Internet-Weltausstellung 1996 sie als eines von nur drei deutschen Projekten im Bereich „Bildung und Erziehung“ in ihre Pavillons aufnahm. Sicher auch, daß ihre Link-Sammlung im Herbst 1996 die zweitvollständigste (nach Rochester) in der Linguistik-Welt war. Ob das heute noch stimmt, ist fraglich; aber in Abständen tragen wir immer nach, bauen um und holen auf. Scheinbar nebenbei fiel 1997 auch noch ein Buch ab („Linguistik im Internet. Das Buch zum Netz – mit CD-ROM“, ein Reiseführer durch die virtuelle Linguistik, erschienen im Westdeutschen Verlag). Und wir waren auf Messen, Ausstellungen, Tagungen und Symposien vertreten.

Doch wo steht LINSE heute, was bietet sie an?

II. Das Angebot

LINSE (<http://www.linse.uni-essen.de>) kommt aus der germanistischen Linguistik; deshalb sind die meisten Beiträge in deutscher Sprache verfaßt. Inhaltlich steht sie aber sämtlichen Themen im weiten Umkreis einer interdisziplinär verstandenen Sprachwissenschaft und Sprachdidaktik offen, wobei verschiedene Schwerpunkte von Neuen Medien über lokale Spezialitäten bis zu geisteswissenschaftlich orientierter Allgemeinbildung reichen. LINSE wendet sich an jede(n), die oder der sich für sprach- und medienbezogene Fragen interessiert, aber auch an drei besondere Zielgruppen, nämlich (1) professionelle Sprachwissenschaftler, (2) Linguistik-Studenten und (3) Lehrer und Schüler. Dabei erfüllt sie mehrere verschiedene Aufgaben. (1) Sie bietet völligen Anfängern (insbesondere aus dem Hochschulbereich) einen ersten Einstieg letzten Endes in das gesamte Internet. (2) Sie publiziert Aufsätze, kleine Schriften, Lernsoftware sowie Rezensionen von Büchern, CDs und Software. (3) Sie liefert Informationen, Bibliographien und Arbeitsmaterial für Forschung und Lehre. (4) Sie dient dem schnellen und direk-

ten Austausch unter Wissenschaftlern und Studenten. (5) Sie ist ein Arbeitsmittel im Studium. „Virtuelle Universität“ kann und soll die traditionelle Universität nicht ersetzen, aber stark bereichern.

Natürlich liefert LINSE ausführliche Informationen über Mitarbeiter, Lehrveranstaltungen und Aktivitäten am Ort. Der weitaus größte Teil des Materials freilich kann von überregionalem und teilweise internationalem Interesse sein. Das komplette Angebot wird einheitlich moderiert und gründlich redigiert. Wir legen Wert auf Qualität und Reichtum der Informationen, Ausnutzung der medienspezifischen Möglichkeiten, ästhetisches Design und hohe Nutzerfreundlichkeit. Klicken wir uns durch!



Unter „*Publikationen*“ finden sich Original-Veröffentlichungen aus der Essener Linguistik, darunter z. B. Aufsätze über Besonderheiten der Sprache im Internet und World Wide Web, über die Sprache von Fernsehnachrichten, über intellektualistischen Sprachstil und über „eloquent silence“.

In der Abteilung „*Rezensionen*“ werden Bücher, Software und CDs besprochen. Unter den Büchern finden sich beispielsweise Volmerts „Grundkurs Sprachwissenschaft“ (mehrfach und kontrovers), Peyer/Portmanns „Norm, Moral und Didaktik – Die Linguistik und ihre Schmuttelkinder“, Hirschs „Übersetzung und Dekonstruktion“, Uskes „Fest der Faulenzer“, Perrin/Jörgs „Netzwerkbuch Computer“, Kaisers „Literarische Spaziergänge im Internet“, Donnellys „In Your Face. The Best of Interactive Interface Design“ (mit screenshots), verschiedene Bücher zu den Themen Internet und CompuServe und nicht zuletzt die neue Brockhaus-Enzyklopädie in 24 Bänden. Außerdem stellen Erstsemester eine ganze Reihe auch älterer linguistische Fachbücher vor. An CDs werden ausführlich besprochen beispielsweise Art Spiegelmanns „Maus“, CD-ROM-Lexika, das „Oxford English Dictionary“, Ingolf Frankes „Sprachlabor“ und der „Sprachkurs Englisch

EuroPlus+ Flying Colours“. Weitere Besprechungen (stets mit screenshots) widmen sich Lern- und Edutainment-Software für Kinder („Burg Drachenstein“, „Alfons Lernsoftware Deutsch“, „Ollis Welt“, „Das Geheimnis der Arche Noah“ und „Zuppel und Guppi“). Außerdem wird der Hexaglot „Sprachen-Computer Euro-Translator“ getestet.

„*Literatur*“ versammelt (teils kommentierte) Bibliographien. Dazu zählt eine Seite, auf die an vielen anderen Orten im WWW verwiesen wird: „LZL – Literatur zur Linguistik“ ist ein kommentierter und mit Leseweg-Vorschlag versehener Kanon ausgewählter Fachliteratur für Linguistik-Studenten, der als Ganzes einen vorzüglichen Einblick in die gesamte Sprachwissenschaft und angrenzende Gebiete geben sollte. Außerdem gibt es Literaturlisten zu Internet und CompuServe sowie zu Computer Mediated Communication und eine recherchierbare Datenbank mit Büchern und Aufsätzen zum Thema Hypertext.

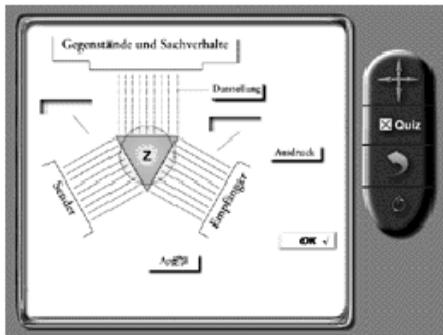
„*ESEL*“, die Essener Studienzyklopädie Linguistik, stellt Lehrmaterial und sorgfältig ausgewählte Seminar- und Examensarbeiten zur Verfügung. Dazu gehören eine Grafik zum Tempusgebrauch im Deutschen, Studienarbeiten und Aufsätze zur Entwicklung des Tempussystems vom Alt- zum Mittelhochdeutschen, zu Aspekten und Aktionsarten (besonders in slawischen Sprachen), zur Zeitbewußtseinsentwicklung bei Piaget, zur Medienkritik bei Platon und heute, zu Roland Barthes, zur e-mail- und Netzkommunikation, zum Schreibenlernen in der Grundschule (Freinet, Vereinfachte Ausgangsschrift, Computer) und andere mehr. Weiterhin gibt es hier die oben schon erwähnte wachsende Sammlung von Kurzrezensionen linguistischer Fachbücher, die Student(inn)en in ihrem ersten Semester verfaßt haben.

Unter „*Projekte*“ werden besondere Essener Aktivitäten dokumentiert, so etwa längerfristige Seminarprojekte und Lernsoftware-Entwicklung.

Das Projekt „*Kuntermund & Löwenmaul*“ fabriziert multimediale interaktive Lernsoftware für Sprache und Linguistik. Ein ausführlicher „Prospekt“, der mit 64 Seiten auch gedruckt vorliegt, informiert über theoretische Hintergründe, didaktisches Konzept und konkrete Pläne. Außerdem gibt es Kostproben aus der laufenden Arbeit, derzeit ein Lernpäckchen zu Bühler und Saussure. Wegen der großen Informationsmenge muß man hier mit längeren Ladezeiten rechnen. Ein Link weist zur teilweise notwendigen Abspielsoftware. In Zukunft möchten wir einige kleinere Lernpäckchen auch direkt zum Download zur Verfügung stellen. So kann man sich ohne jede Zusatzsoftware alles in Ruhe zu Hause ansehen.



Multimediale interaktive Lernsoftware für Sprache und Linguistik



- [Prospekt](#)
- [Lernpäckchen](#)
- [Semiotik](#)

Es werden solche „Tagungen“ dokumentiert, die das LINSE-Team oder einzelne Mitarbeiter mit organisieren. Dazu zählen etwa das Symposium Deutschdidaktik und das internationale LAUD-Symposium 1998 „Humboldt and Whorf revisited: Universal and Culture-Specific Conceptualizations in Grammar and Lexis“. „Reportagen“ berichten von besonderen Veranstaltungen im Umkreis des LINSE-Teams.

„OBST“ ist die Seite der „Osnabrücker Beiträge zur Sprachtheorie“. Sämtliche bisher erschienenen Hefte und Beihefte (also seit 1976) sind mit Titelbild, Inhaltsverzeichnis und ggf. Editorial dokumentiert. Der gesamte Bestand kann in einer Datenbank recherchiert werden. Die Reihe oder einzelne Bände können auch direkt bestellt werden.

Die Zeitschrift „Sprache und Datenverarbeitung“ ist in ähnlicher Weise vollständig (also seit 1977) und mit recherchierbarer Datenbank dokumentiert.

„Diskussion“: An dieser Stelle sind bisher drei öffentliche Diskussionsforen eingerichtet, und zwar zur Linguistik allgemein, zum Thema „Schule und Computer“ und zu neuen Medien. Hier kann jeder Fragen stellen, Thesen vertreten, Probleme wälzen und wird von irgendwoher auf der Welt Antworten bekommen.



„Schule und Computer“ ist eine eigene Abteilung für Lehrer, Lehramtsstudenten und Schüler, in der sämtliche Fragen rund um den Einsatz von Internet und Computern in der Schule behandelt werden, und zwar sowohl allgemein als auch mit einem fachdidaktischen Schwerpunkt auf dem Deutschunterricht. Sorgfältig ausgewählte und kommentierte Links führen auf wichtige Angebote auf anderen Servern; es gibt eigene Beiträge (z. B. einen Aufsatz über Computer im Schreibunterricht der Grundschule) und ein öffentliches Diskussionsforum.

Hinter „Links“ steckt eine der weltweit umfangreichsten Sammlungen linguistisch relevanter Adressen im WWW. Die Adressen sind in Rubriken sortiert und mit Stichworten kurz kommentiert. Die Liste wird ständig erweitert und aktualisiert

und schreibt somit die (allerdings sehr viel ausführlicheren) Angaben in der Monographie „Linguistik im Internet“ fort.

Die Taste „*Suchen im WWW*“ führt zu einer bequemen Sammlung von Suchmaschinen, die das gesamte World Wide Web erschließen. Hinter „Rätsel“ verbirgt sich eine alle paar Monate wechselnde Preisfrage mit Gewinnchancen; im Herbst 1997 beispielsweise ist es ein semiotisches Kreuzworträtsel zu Ecos „Der Name der Rose“.

Das rote Schildchen „*Neu*“ zeigt dem regelmäßigen Nutzer der LINSE (von denen es eine ganze Reihe gibt), welche Beiträge in letzter Zeit neu aufgenommen wurden, damit sie oder er überflüssige Klickzeit spart.

So weit, so gut, aber auch so verbesserungswürdig. Kritik und Vorschläge sind jederzeit willkommen, am einfachsten per e-mail über die auf der Homepage angegebene Adresse. Was wird in nächster Zeit hinzukommen? Weitere Aufsätze, Seminararbeiten und Rezensionen, zusätzliche Links, ein Wittgenstein-Projekt, neue Lernsoftware, die Homepage von LAUDE (Linguistic Agency Universities of Duisburg & Essen) und sicherlich die eine oder andere Überraschung, von der wir auch noch nichts wissen. Klicken Sie mal rein: <http://www.linse.uni-essen.de> – wir sind gespannt auf Ihr Urteil.

Computerlinguistik für Philologen

M. A. Moreaux / Inalco Paris

Es läßt sich nicht mehr bestreiten, daß Computerlinguistik ein weitgehend interdisziplinäres Fachgebiet darstellt. Unter den häufig aufgezählten Kenntnissen, die ein Computerlinguist beherrschen muß, werden Kenntnisse aus Linguistik und Informatik mit Recht als grundlegend betrachtet. Von entscheidender Bedeutung scheint jedoch ein drittes Gebiet, das selten oder nie sehr deutlich erwähnt wird, so, als ob man es für ganz selbstverständlich hielte: Hierbei handelt sich um alles, was Sprachfähigkeit und Sprachkenntnisse betrifft. Wir möchten deswegen ein besonderes Augenmerk auf den Stellenwert von Sprachkenntnissen in der LDV bzw. Computerlinguistik richten.

Wichtig zu bemerken ist, daß der Gegenstand der maschinellen Sprachverarbeitung nicht Sprache an sich ist, sondern eine Folge von in einer bestimmten Sprache ausgedrückten Äußerungen. Ein Sprachverarbeitungssystem muß also über ein sehr präzises und detailliertes Wissen über die Elemente und Strukturen der betroffenen Einzelsprache verfügen und dürfte um so robuster sein, desto mehr es in der Lage ist, aus seinen eigenen Sprachdaten Eigenschaften sprachlicher Einheiten zu errechnen. Daraus ergibt sich, daß eine maschinell interpretierbare Sprachdarstellung nicht nur völlig explizit formuliert werden muß, sondern auch bis ins einzelne alle Regelmäßigkeiten und Ausnahmen der zu behandelnden Phänomene beschreiben muß, auch die Regularitäten in den Irregularitäten.

Selbstverständlich müssen Sprachforschungen durch Korpora unterstützt werden. Deren Gebrauch dürfte jedoch Sprachkenntnisse nicht ersetzen, sondern von diesen Kenntnissen gesteuert werden, um eine Verfeinerung bzw. ein Ergänzen oder die Auswertung der Darstellung zu erlauben.

Maschinell interpretierbare Sprachdarstellungen sollen auch nicht bloße Aufzählungen von aus Texten und Reden direkt beobachtbaren Sprachfakten sein. Sie müssen theoretisch geleitet werden und implementierbar sein.

Die Erfordernisse der maschinellen Sprachverarbeitung zwingen dem Sprachwissenschaftler also eine strenge Methodik, eine sehr analytische Denkweise und eine ganz besondere Einstellung zur Sprache auf. Infolgedessen kann sich

eine Modellimplementierung nicht aus der Zusammenarbeit von Spezialisten mit unterschiedlichen Kompetenzen ergeben. Von einem Computerlinguisten wird unbedingt eine Mehrfachkompetenz erwartet. Er muß in der Lage sein, die vollständige Entwicklung eines computerorientierten Sprachmodells auszuführen, von der Problemstellung über den Entwurf einer linguistischen Lösung bis zur Implementierung.

Hierbei sollte man insbesondere darüber nachdenken, welcher Art die Sprachkenntnisse eines Computerlinguisten sein sollen. Die erwähnten Anforderungen setzen eine eingehende Kenntnis einer Einzelsprache voraus. Man benötigt darüber hinaus aber auch sprachwissenschaftliche Kenntnisse von der zu beschreibenden Sprache, auch wenn sie die Muttersprache ist, denn die als Muttersprachler erworbene Sprachfähigkeit scheint hier nicht ausreichend zu sein.

Äußerst wichtig ist, daß eine größere Zahl von Philologen Interesse an maschineller Sprachverarbeitung gewinnt und in die Computerlinguistik einsteigen möchte. In dieser Hinsicht ist das INALCO (Institut National des Langues et Civilisations Orientales*) in Paris eine sehr günstige Umgebung. Mit Ausnahme von westeuropäischen Sprachen können dort 80 Sprachen aus allen Erdteilen erlernt werden. Da diese Sprachen an höheren Schulen nicht unterrichtet werden, wird jeder am Institut immatrikulierte Student als Anfänger betrachtet. Je nach Sprache erstreckt sich das Anfangsstudium über zwei bis drei Studienjahre* und wird mit einem Diplom (genannt DULCO) abgeschlossen.

Das computerlinguistische Curriculum am INALCO wurde speziell für Sprachenstudierende entwickelt und setzt nach einem absolvierten sprachlichen Grundstudium von zwei bis drei Jahren ein. Der Studiengang, der somit erst Bestandteil des Hauptstudiums ist, umfaßt ein

- 1.) computerlinguistisches Grundlagenstudium, das zwei in sich abgeschlossene Teile enthält, die jeweils ein Studienjahr dauern: der erste Abschluß ist die „Licence“ und der zweite die „Maîtrise“.
- 2.) Promotionsstudium, das ebenfalls zweiteilig ist. Ein erster Teil, der ein Studienjahr dauert, führt zu dem als DEA bezeichneten

* Wortwörtlich wäre diese Bezeichnung als „Staatliches Institut für orientalische Sprachen und Kulturen“ zu übersetzen. „Orientalisch“ versteht sich hier aber viel breiter als üblich und betrifft nicht nur den geographischen Orient.

† In Frankreich dauert ein Studienjahr von Oktober bis Juni und umfaßt 27 Studienwochen.

Abschlußdiplom und ist als der erste Schritt in die Forschung zu betrachten. Der zweite Teil dehnt sich über einen Zeitraum von 3 bis 5 Jahren aus. Während dieser Zeit erarbeiten die Promovenden ihre Dissertation.

Alle Absolventen des Grundstudiums einer beliebigen Philologie* können sich im Studiengang „Computerlinguistik“ einschreiben. Diese Studenten haben die jeweils studierte Sprache schon recht gut erlernt, kennen ihre Grammatik und sind infolgedessen in der Lage, über diese Sprache als Sprachsystem nachzudenken. Bemerkenswert ist, daß sie nicht selten mehrere Sprachen studiert haben. Das kann nur von Vorteil sein, denn damit werden sie darauf vorbereitet, Funktionsunveränderlichkeit an unterschiedlichen Äußerungsformen zu erkennen. Meistens haben sie aber noch kaum Einblicke in die (allgemeine) Linguistik bekommen. Es wird von den Studenten auch nicht erwartet, daß sie Kenntnisse aus dem Bereich der Informatik besitzen.

Angesichts des Wissens und der Fähigkeiten der in die Computerlinguistik einsteigenden Studenten werden die angebotenen Lehrveranstaltungen während der beiden ersten Studienjahre auf die Grundlagenausbildung gerichtet. Im allgemeinen sind die Studenten der philologischen Fächer weder an Interdisziplinarität, noch an die mathematisch präzise Denkweise gewohnt, die Computerlinguistik erfordert. Deswegen geht es darum, Grundlagen zu vermitteln, durch die der Wissensstand der Studierenden in bezug auf interdisziplinäres Wissen erweitert und sie, wenn man es so sagen darf, in eine andere Denkweise einführt. Der Schwerpunkt liegt auf den Fächern, die von einem Philologen als neu empfunden werden und ihm die größte Mühe bereiten: alles, was sich auf Formalisierung, Algorithmenbeschreibung und Programmierung bezieht.

Die Licence- und Maîtrise-Lehrpläne sind in vier Komponenten untergliedert. Eine betrifft die Sprache, denn jeder Student, der sich für ein Computerlinguistikstudium entscheidet, muß einen Teil seiner Lehrveranstaltungen im Bereich „Sprache“ absolvieren. Jede der drei weiteren Komponenten entspricht der Vermittlung von Grundkonzepten, Methoden und Verfahren, die zu den beteiligten Disziplinen gehören und deren Kenntnis notwendig ist, um bestehende CL-Modelle verstehen und evaluieren zu können oder solche Modelle selbst zu entwick-

* Nicht nur Studierende einer „orientalischen“ Sprache, sondern auch jeder westeuropäischen Sprache (Französisch, Deutsch, Englisch, ...).

keln. Alle Lehrveranstaltungen sind Pflichtveranstaltungen und werden in Form von Vorlesungen und Übungen angeboten:

- 1.) Sprache (Licence-Studiengang: 100 Std.*; Maîtrise-Studiengang: 50 Std.): Dient der Vertiefung der Sprachkenntnisse.
- 2.) Sprachwissenschaft (Licence-Studiengang: 100 Std.; Maîtrise-Studiengang: 63 Std.): Einführung in die Grundlagen der Sprachwissenschaft. Die eingeführten Begriffe und Methodologien werden dann auf die Einheiten der verschiedenen sprachlichen Ebenen (Phonetik/Phonologie, Morphologie, Syntax, Semantik aber auch Lexikologie) angewandt. Nach einem kurzen Überblick über Anwendungsgebiete und Forschungsrichtungen der CL versucht eine der linguistischen Veranstaltungen, die Verhältnisse zwischen den verschiedenen Fächern zu skizzieren und hierbei die Studenten mit den Bedingungen einer computerorientierten Sprachbeschreibung vertraut zu machen.
- 3.) Sprachverarbeitung bzw. Computerlinguistik (Licence-Studiengang: 75 Std.; Maîtrise-Studiengang: 125 Std.): Den thematischen Schwerpunkt bildet hier die Behandlung formaler Modelle und formaler Darstellungsverfahren. Die Licence-Veranstaltungen führen in die den Philologen meistens fehlenden Grundlagen der Mathematik und der Logik ein (moderne Logik, Mengentheorie und Relationskalkül). Diese Grundkenntnisse werden dann in den Maîtrise-Veranstaltungen vertieft und durch die Darstellung der nicht-klassischen Logiken, der Theorie formaler Sprachen und der Automatentheorie erweitert. Zum Schluß werden auch Parsing-Strategien behandelt.
- 4.) Informatik und Programmierung (Licence-Studiengang: 88 Std.; Maîtrise-Studiengang: 150 Std.): Nach einer kurzen Einführung in den Aufbau und das Funktionieren eines Computers wird das Hauptgewicht auf Entwurfsprinzipien von Algorithmen und Datenstrukturen gelegt. Die Grundlagen der strukturierten Programmierung werden durch das Erlernen einer ersten

* Stunden pro Studienjahr

Programmiersprache vermittelt. Dabei ist C die hier gewählte Programmiersprache. Darüber hinaus werden dann im Rahmen des Maîtrise-Studiengangs auch die Konzepte der logik- und objektorientierten Programmierung dargestellt und Prolog und C++ erlernt. Die Studenten erhalten die Aufgabe, Basisalgorithmen der Sprachverarbeitung zu erarbeiten oder eigene Lösungen zu erstellen, die sie dann in einer Programmiersprache umsetzen müssen.

Die DEA-Veranstaltungen sind auf speziellere Kenntnisse bezogen, wie z. B. auf formale Theorien und Beschreibungsformalismen in der Computerlinguistik, die theoretisch und mathematisch komplex sind und das früher vermittelte Wissen unbedingt voraussetzen. Ein Computerlinguist muß natürlich in der Lage sein, abschätzen zu können, ob eine Grammatiktheorie und ein Beschreibungsformalismus zur Abbildung und Erklärung der zu behandelnden Phänomene besser als andere geeignet sind. Die Beteiligung von Forschern anderer französischer Hochschulen (Grenoble und Nizza), aber auch aus mehreren europäischen Ländern (Deutschland, Tschechische Republik, Großbritannien, Belgien) ermöglicht den Studenten, einen weitreichenden Überblick über zahlreiche Bereiche der Computerlinguistik zu gewinnen.

Der DEA-Lehrplan umfaßt vier Lehrveranstaltungsgruppen. Jede ist einem besonderen Thema gewidmet und wird mit einer Prüfung abgeschlossen:

- 1.) formale Beschreibungsmodelle (100 Std.)
- 2.) Methoden zur morphologischen, syntaktischen und semantischen Analyse (100 Std.)
- 3.) Anwendung von KI-Methoden und Verfahren im Bereich der Sprachverarbeitung (50 Std.)
- 4.) Anwendungsgebiete der Computerlinguistik und Systeme (125 Std.)

Der Studierende wird in eine Forschungsgruppe integriert, in der er an selbständiges wissenschaftliches Arbeiten gewöhnt wird. Zum Abschluß muß er sich inhaltlich und technisch mit einer bestimmten Problemstellung auseinandersetzen. Dabei erarbeitet er ein begrenztes aber nicht-triviales Phänomen der studierten Sprache, muß seine Ergebnisse in einer längeren schriftlichen Hausarbeit

darstellen und seine Lösung in ein Programm umsetzen. Zu den derzeit im Rahmen solcher Arbeiten behandelten Sprachen gehören Französisch, Englisch, Deutsch, Arabisch, Tschechisch, Italienisch, Malaiisch und Japanisch.

Das skizzierte Ausbildungsprofil wurde 1979 nach einem Besuch von Patrice Pognan, dem Leiter des CERTAL (Centre d'Etudes et de Recherche en grammaire et Traitement Automatique des Langues*), in Hamburg entwickelt und ist in vielerlei Hinsicht Ergebnis der Diskussionen mit Walther von Hahn. Es ist in den nachfolgenden Jahren durch zahlreiche eigene Forschungs- und Lehrerfahrungen ergänzt und den Bedürfnissen der Studenten des INALCO angepaßt worden.

*eine Forschungsgruppe, die ihren Sitz im INALCO hat

Interdisziplinäre Teamarbeit an Hochschulen

Erwartungen und Erfahrungen von Studierenden

*Studenten dreier Fachbereiche an der
Fachhochschule Darmstadt erarbeiten ein
gemeinsames Projektergebnis*

*Gerhard Knorz
Fachhochschule Darmstadt,
Fachbereich Information und Dokumentation*

1 Das interdisziplinäre Projekt „Aufbau eines WWW-Servers“

Im Sommersemester 1997 wurde an der Fachhochschule Darmstadt ein gemeinsames Projekt verschiedener Fachbereiche mit dem Ziel durchgeführt, den Prototypen eines *datenbankgestützten WWW-Informationsservers* für die Fachhochschule zu entwickeln:

- mit seiner inhaltlichen Struktur,
- seinen Navigationsmöglichkeiten,
- seiner Benutzungsoberfläche,
- seiner technischen Realisierung
- und den organisatorischen Fragen bei Aufbau und Betrieb.

Die Basisanforderungen an Informationsumfang und -strukturierung sind durch die detaillierten Vorgaben des *European Credit Transfer Systems (ECTS)* gegeben. Dieses europäische Programm will das Studieren an europäischen Hochschulen durchlässiger machen, insbesondere dadurch, daß es von den Hochschulen definierte Informationen über alle Aspekte von Studium und dessen Randbedingungen abfordert. Weitere Aspekte ergeben sich aus den Interessen der organisatorischen Einheiten der Hochschule, insbesondere der Fachbereiche und ihrer Mitglieder. Dazu gehören auch Anforderungen durch mögliche zukünftige Nutzungsoptionen wie etwa die Produktion des Vorlesungsverzeichnisses.

Das Projekt kam durch Initiativen und Lehrveranstaltungen aus drei Fachbereichen zustande:

- Fachbereich *Gestaltung* mit einem Entwurfsprojekt
- Fachbereich *Informatik* mit einem Softwareentwicklungspraktikum
- Fachbereich *Information und Dokumentation* mit einem studentischen Projekt

Das Projektziel wurde insgesamt erreicht. Zum Abschluß des Projektes wurden Konzeption und Prototyp den Verantwortlichen in der Hochschule im Rahmen einer Präsentation vorgestellt und diskutiert. Die erarbeitete Entwicklungslinie sowie die erreichten Ergebnisse wurden bestätigt und eine Weiterführung der Arbeiten verabredet.

Mit dem Ziel, Einstellungen und Erfahrungen der TeilnehmerInnen hinsichtlich der in der Hochschule nicht alltäglichen Arbeitsform eines interdisziplinären Projektes über den subjektiven Eindruck der Veranstalter hinaus zu erfassen, wurde zu Beginn und zum Abschluß des Projektes ein einfacher Fragebogen eingesetzt. Die Auswertung der Antworten und das sich daraus ergebende Bild ist das Thema dieses Beitrags.



*Abb. 1: Illustration eines der Teilergebnisse studentischer Arbeitsgruppen:
Entwurfsvariante zur Navigation im Informationssystem*

2 Innere und äußere Bedingungen der Initiierung und des Verlaufs des Projektes

Die Situation von Hochschullehre, insbesondere an Fachhochschulen bei einem Lehrdeputat von 18 Semesterwochenstunden, ist keineswegs so, daß kreative und kommunikative Freiräume Bedingungen für die Planung und Durchführung von hochaktuellen, innovativen und interdisziplinär angelegten Lehrveranstaltungen wären. So stehen der Umsetzung vielzitiert Forderungen nach stärkerer Förderung von Schlüsselqualifikationen bei Studierenden ganz konkrete Hemmnisse entgegen. Neue Formen von Lehrveranstaltungen und Zusammenarbeit über die Grenzen von Fachbereichen hinweg sind vielfach das Ergebnis von Initiativen einzelner und von glücklichem Zusammentreffen unkoordinierter Entwicklungen.

In diesem Sinn ist auch dem Projekt „Aufbau eines WWW-Servers“ keine weitsichtige Planung und fachlich ins Detail gehende Vorbereitung vorausgegangen. Vielmehr hatten sich die Veranstalter (Knorz/Krier/Pfedorf) dreier unabhängiger Projekte zunächst bilateral um Kooperationsmöglichkeiten bemüht, um dann festzustellen, daß die angestrebten Ergebnisse im Grunde nur im Gesamtzusammenhang, also (faktisch) in *einem* Projekt sinnvoll erarbeitet werden sollten. Die letztendliche Abstimmung dieses Vorhabens reichte zeitlich bis an den Anfang des Semesters und damit der Lehrveranstaltungen heran. Aus Verwaltungssicht (soweit es also etwa die Leistungsnachweise betrifft) verblieben es 3 Veranstaltungen.

Studierende verschiedener Fachbereiche in einem Projekt erfolgreich zusammenzubringen ist bereits ein Problem der terminlichen Koordination. Daß *ein* Wochentag als Projekttag definiert werden konnte, war Bedingung, aber gleichzeitig ein erster Erfolg für das Vorhaben. In der Realität sorgen dann unterschiedliche Veranstaltungsumfänge (4 – 8 SWS), unterschiedliche Veranstaltungs-Startzeiten (Datum, Uhrzeiten), konkurrierende Verpflichtungen, Wegezeiten zwischen räumlich z. T. weit entfernten Fachbereichen, die immer noch hinderliche Kluft zwischen PC- und Mac-Welt sowie unterschiedliche „Ausbildungskulturen“ für eine genügende Anzahl kleinerer und größerer Probleme.

Eine der Konsequenzen aus dem Entstehungsgeschichte des Projektes war die Tatsache, daß die Veranstalter wenig Vorwissen über und Einfluß auf die Gesamtanzahl der Teilnehmer hatten. Mit insgesamt über 30 TeilnehmerInnen (und zu Beginn über 40) waren Erwartungen und Wunschgröße deutlich übertroffen. Das Projektmanagement hatte demnach, sowohl was die Unterstützung und Koordination der Arbeitsgruppen, als auch was die Diskussion und Entscheidungsfindung im Plenum betraf, eine schwierige Aufgabe. Auch die Arbeitsgruppen selbst erwiesen sich erst nach weiterer Untergliederung als effektiv arbeitsfähig.

Das Projektmanagement lag in den Händen einer 4-köpfigen Studentengruppe, die – weil dieser Fachbereich terminlich einen Start- und Planungsvorteil verbuchen konnte – aus dem Fachbereich Information und Dokumentation entstammten. Die übrigen Arbeitsgruppen, die sich zum Teil während des Projektes mit neuen Aufgaben neu konfigurierten, waren bewußt mit Teilnehmern jeweils verschiedener Fachbereiche besetzt.

Die „einseitige“ Besetzung der Projektmanagementgruppe war sicher in manchen (vornehmlich frühen) Phasen des Projektes problematisch für die Identifikation der Teilnehmer mit den Projektzielen und der Gesamtgruppe: („Euer Projekt ...“). Eine andere Problematik stellt die Rolle der beteiligten Hochschullehrer dar, die natürlich sich in der Gefahr, Versuchung oder vor der Notwendigkeit sehen, ihre zurückgenommene Rolle als Leitungsgremium (bzw. für die Arbeitsgruppen als Berater) zu verlassen und aktiv in „das Tagesgeschehen“ einzugreifen, womit sie dann in Konsequenz die Autorität des Projektmanagements untergraben.

Als zentrales technisches Koordinations- und Informationsmedium fungierte das Internet mit e-mail und dem Web-Server „*WebSite ,Methodik‘*“, in dem das Projekt wie alle anderen Lehrveranstaltungen des Faches Informationsmethodik tagesaktuell mit Zusammenfassungen (Protokollen), Materialien und Arbeitsergebnissen präsent war (<http://www.iud.fh-darmstadt.de/iud/wwwmeth/index.htm>). Diese Infrastruktur wurde ergänzt durch einen projekteigenen Web-Server und eine erst zu Projektende ansatzweise fertiggestellte Protokoll-Datenbank. Technische Probleme bei der Internetanbindung eines entfernt untergebrachten Fachbereichs und bei der Weiterentwicklung der benutzten Werkzeuge, Schwachstellen bei Aktualisierung und Informationsaufbereitung und das erst allmählich sich entwickelnde Bewußtsein dafür, wie und wozu die gegebenen Möglichkeiten als selbstverständliche Unterstützung zu nutzen sind, haben den praktischen Wert dieser Projektinfrastruktur, besonders in der Anfangsphase, deutlich limitiert.

3 Erwartungen und Erfahrungen der TeilnehmerInnen

3.1 Die Befragung

Mit dem Ziel, Einstellungen und Erfahrungen der TeilnehmerInnen hinsichtlich der in der Hochschule nicht alltäglichen Arbeitsform eines interdisziplinären Projektes über den subjektiven Eindruck der Veranstalter hinaus zu erfassen, wurde nach der Anfangsphase und zum Abschluß des Projektes (30.4.97 und 25.6.97) ein einfacher Fragebogen mit nur 2 Fragekomplexen (insgesamt 5 Fragen) und der Gelegenheit zur freien Kommentierung ausgeteilt (siehe Abb. 2).

Der erste Fragekomplex fragte nach dem Einfluß der interdisziplinären Arbeit auf:

- die Qualität des Ergebnisses und auf
- die Arbeitszufriedenheit.

Dabei erkundete die 1. Befragung die *Erwartung* und die 2. Befragung die *Erfahrungen* der TeilnehmerInnen.

Der zweite Fragenkomplex war bei 1. und 2. Befragung unterschiedslos formuliert: *Wie schätzen Sie die Fähigkeit ein, in einem interdisziplinären Team zu arbeiten?* Drei Aspekte wurden abgefragt:

- Ist diese Fähigkeit gegenwärtig in der Berufspraxis wichtig?
- Ist diese Fähigkeit zukünftig in der Berufspraxis wichtig?
- Bedarf diese Fähigkeit einer Förderung in der Ausbildung?

Zur Beantwortung waren Skalen von +3 (sehr positiv/unterstützend) bis -3 (sehr negativ/ablehnend) vorgegeben. Mit einem Sonderzeichen konnte die Frage zurückgewiesen werden (keine Einschätzung/weiß nicht).

Wie ist Ihre Erwartung (Befragung 1) bzw. Erfahrung (Befragung 2)?

- A1: *Wie wirkt sich der interdisziplinäre Ansatz auf das Ergebnis aus?*
- A2: *Wie wirkt sich der interdisziplinäre Ansatz auf Ihre Arbeitszufriedenheit aus?*

Wie schätzen Sie die Fähigkeit, in einem interdisziplinären Team zu arbeiten, ein?

- B1: *Ist diese Fähigkeit gegenwärtig in der Berufspraxis wichtig?*
- B2: *Ist diese Fähigkeit zukünftig in der Berufspraxis wichtig?*
- B3: *Bedarf diese Fähigkeit einer Förderung in der Ausbildung?*

Abb. 2: Die Fragen des Fragebogens im Wortlaut.

Befragung 1 (28 Rückläufe)	<i>A1</i>	<i>A2</i>	<i>B1</i>	<i>B2</i>	<i>B3</i>	<i>Kommentare</i>
<i>Anzahl fehlender Antworten („weiß nicht“)</i>	2	1	0	0	2	19
Anzahl Antworten	26	27	28	28	26	9

Befragung 2 (26 Rückläufe)	<i>A1</i>	<i>A2</i>	<i>B1</i>	<i>B2</i>	<i>B3</i>	<i>Kommentare</i>
<i>Anzahl fehlender Antworten („weiß nicht“)</i>	1	0	1	2	0	16
Anzahl Antworten	25	26	25	24	26	10

Abb. 3: Quantitative Übersicht über den Rücklauf der Fragebogen

In den Balkendiagrammen zur Ergebnisdarstellung wird durchgängig die Farbe Blau der ersten und die Farbe Weinrot der zweiten Befragung zugeordnet.

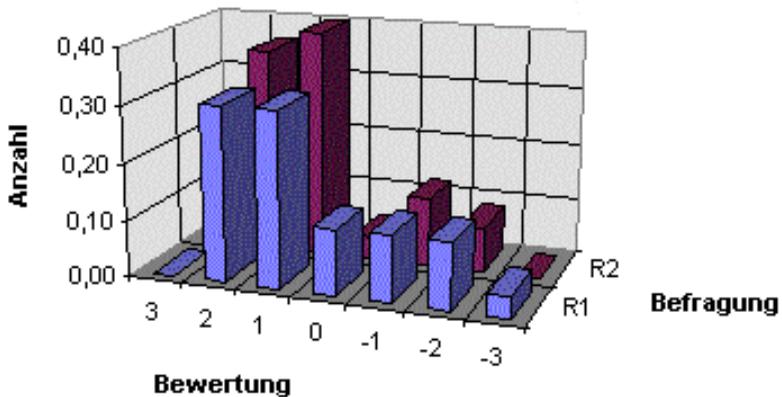


Abb. 4: Auswertung der 2 Befragungen bezüglich der Frage A1

- **Frage A1:** *Wie wirkt sich der interdisziplinäre Ansatz auf das Projektergebnis aus?* (Abb. 4)

Interpretation: Die Erfahrung hat die Erwartung hinsichtlich *Ergebnisqualität* leicht ins Positive verschoben. Insgesamt wird von einem vorsichtig bis deutlich günstigem Einfluß auf die Qualität des Ergebnisses ausgegangen.

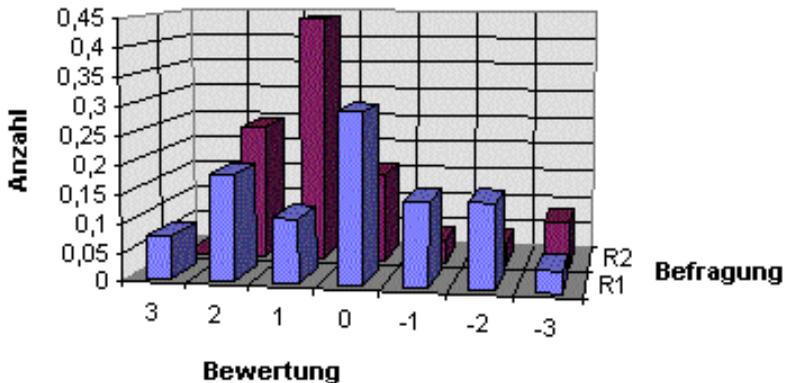


Abb. 5: Auswertung der zwei Befragungen bezüglich der Frage A2

- **Frage A2:** *Wie wirkt sich der interdisziplinäre Ansatz auf Ihre Arbeitszufriedenheit aus?* (Abb. 5)

Interpretation: Der Einfluß der Interdisziplinarität auf die Arbeitszufriedenheit wird nach der Projekterfahrung eindeutig, allerdings noch verhalten positiv eingeschätzt. Die Erwartungen dagegen waren zunächst vornehmlich indifferent bzw. ausgeglichen positiv und negativ gewesen. Hier hat also die Erfahrung eine anfängliche Unsicherheit und Skepsis widerlegt. Gleichmaßen wurden aber auch überzogene Erwartungen gedämpft.

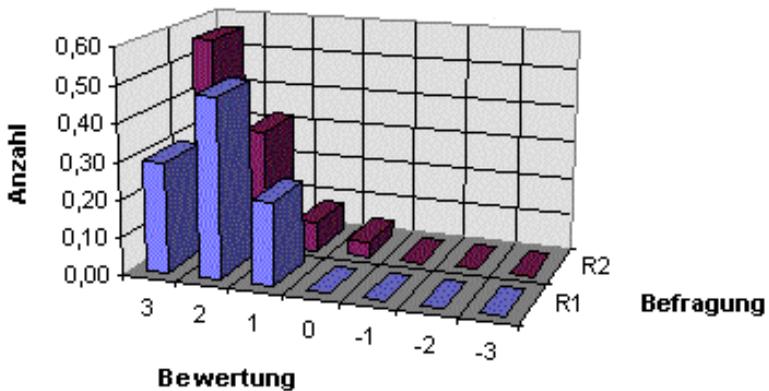


Abb. 6: Auswertung der 2 Befragungen bezüglich der Frage B1

- **Frage B1:** *Wie schätzen Sie die Fähigkeit, in einem interdisziplinären Team zu arbeiten, ein: Ist diese Fähigkeit gegenwärtig in der Berufspraxis wichtig?* (Abb. 6)

Interpretation: Arbeiten im interdisziplinären Team halten die Projektteilnehmer bereits in der *gegenwärtigen* Berufspraxis für eindeutig wichtig. Die Projekterfahrung hat diese Einschätzung so verstärkt, daß „sehr wichtig“ die dominante Angabe wird.

- **Frage B2:** *Wie schätzen Sie die Fähigkeit, in einem interdisziplinären Team zu arbeiten, ein: Ist diese Fähigkeit zukünftig in der Berufspraxis wichtig?* (Abb. 7)

Interpretation: Arbeiten im interdisziplinären Team halten die Projektteilnehmer für die *zukünftige* Berufspraxis für sehr wichtig. Die Projekterfahrung hat diese Einschätzung so verstärkt, daß „sehr wichtig“ von mehr als dreiviertel der Studierenden genannt wird.

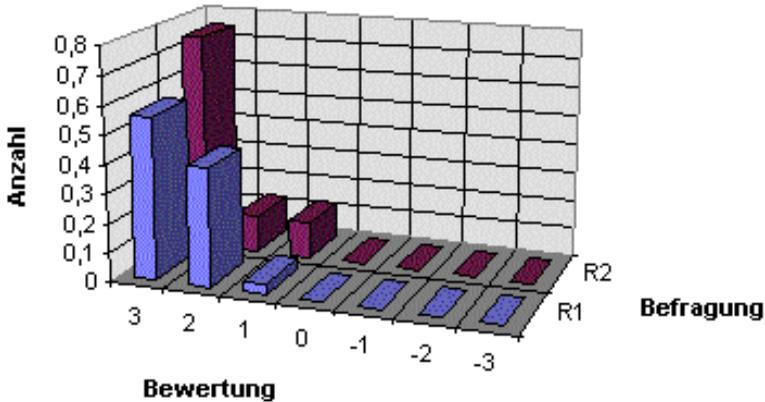


Abb. 7: Auswertung der 2 Befragungen bezüglich der Frage B2

- **Frage B3:** *Wie schätzen Sie die Fähigkeit, in einem interdisziplinären Team zu arbeiten, ein: Bedarf diese Fähigkeit einer Förderung in der Ausbildung? (Abb. 8)*

Interpretation: Daß das Arbeiten im interdisziplinären Team einer Förderung in der Ausbildung bedarf, war bereit zu Beginn des Projektes die unwidersprochene Einschätzung der TeilnehmerInnen. Nach der Projekterfahrung wird noch klarer erkannt, daß eine Förderung nicht nur wichtig, sondern „sehr wichtig“ ist.

Anmerkung: Ein Teilnehmer ist der nicht auf die Allgemeinheit übertragbaren Ansicht, daß er nach einer abgeschlossenen Lehrerausbildung das Arbeiten im Team nicht mehr üben muß.

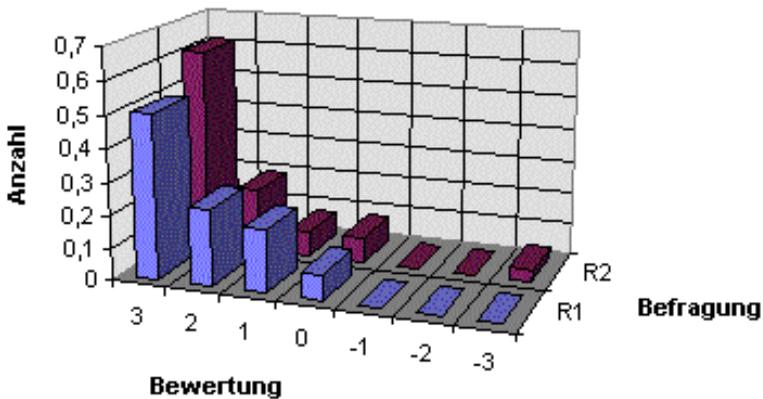


Abb. 8: Auswertung der 2 Befragungen bezüglich der Frage B3

3.3 Kommentare der Studierenden

In beiden Befragungen haben Studierende von der Möglichkeit Gebrauch gemacht, sich unabhängig von den konkreten Fragen frei zu äußern.

Im Folgenden werden ausschließlich die Argumente und Einschätzungen zusammengefaßt, die von immerhin mehr als einem Drittel aller Studierenden *nach* der Projekterfahrung abgegeben wurden. Anzahlen sind in Klammern gesetzt.

Als positive Erfahrung wird genannt, von den Spezialisten anderer Fachbereiche zu lernen (1) und mehr Offenheit und Sensibilität zu erwerben (1). Zwar brauche man am Anfang länger, aber dafür ginge es anschließend umso effizienter (1). Oder aber dasselbe als Kritik formuliert: Die Teamarbeit kam erst spät nach dem Chaos (2). Als wichtig wird eingeschätzt, Teamarbeit einzuüben (2) sowie Wissen zu integrieren und in Entscheidungen umzusetzen (1). Auch die Fortsetzung des Projektes wird als ein wichtiger Punkt genannt.

Zu besseren Arbeitsbedingungen hätten eine geringere Gesamtteilnehmeranzahl (1) mit kleineren Arbeitsgruppen (3) geführt. Ebenso auch ein besserer/gleichmäßigerer Wissenstand (über Projektziele und -kontext) der TeilnehmerInnen aus

den unterschiedlichen Fachbereichen (1) sowie eine gezieltere Zusammensetzung der Arbeitsgruppen anhand der Kompetenzen der TeilnehmerInnen (1). Auf diese Weise wäre die Koordination des Gesamtprojektes (2) und auch der Wissenstransfer innerhalb der Arbeitsgruppen (1), insbesondere in der Anfangszeit (2), nicht so problematisch gewesen.

Auch Selbstkritik innerhalb der Studentenschaft wurde geübt: Mehr Selbständigkeit wird angemahnt (1), mehr Teamgeist und Identifikation mit den Projektzielen (2), und mehr Engagement einzelner, die sich hinter den Fleißigen in einer Gruppe verstecken (1).

4 Fazit

An dem am Einzelfall erhobenen Befund läßt sich ablesen, daß Studierende in Übereinstimmung mit den aus Industrie und Wirtschaft erhobenen Forderungen das Arbeiten im interdisziplinären Team in der Berufspraxis bereits gegenwärtig und verstärkt noch zukünftig als sehr wichtig einschätzen. Gleichwohl stehen sie einer Lehrform, die diese Schlüsselqualifikation von ihnen verlangt, im konkreten Fall zunächst skeptisch bis indifferent gegenüber. Zwar erwarten sie mehrheitlich ein eventuell besseres Arbeitsergebnis, aber viele fühlen sich in ihrer eigenen Ausbildungskultur am wohlsten.

Die konkrete Erfahrung mit interdisziplinärem Arbeiten bewirkt deutlich einen Shift in den Einschätzungen: Sie verstärkt die positiven Erwartungen und stellt fest, daß die Teamarbeit sich letztlich durchaus positiv auf die Arbeitszufriedenheit auswirkt. Die Notwendigkeit zur Förderung dieser für viele neuen Arbeitsform wird deutlich klarer erkannt und auch der Bezug zur Berufspraxis wird (noch) deutlicher wahrgenommen.

Die globale Auswahl des Fragebogens ohne Differenzierung nach *Fachbereichszugehörigkeit* kann nicht wiedergeben, was die direkte Beobachtung deutlich erkennen läßt: Daß nämlich die Bereitschaft zu und die subjektive Wahrnehmung von interdisziplinärer Teamarbeit klar mit diesem Faktor zusammenhängt.

Aus Sicht der Veranstalter zeigt der Rückblick sicher vieles, was im Wiederholungsfall zu vermeiden oder zu verbessern wäre. Aber nicht nur, weil man aus Fehlern lernen kann, lohnt sich der Mut (und der Aufwand) zum Risiko: Der größte Fehler wäre es, es eine Kooperation über Fachbereichsgrenzen hinweg erst gar nicht zu versuchen!

Quellen

Projekt: Aufbau eines WWW-Servers (SS97)

Einstiegsdokument in die Planung und die Dokumentation der Durchführung der Lehrveranstaltung „Aufbau eines WWW-Servers (SS97)“ unter Einschluß aller Beiträge aller beteiligten Fachbereiche (Gestaltung, Informatik, Information und Dokumentation); Fachhochschule Darmstadt, Fachbereich Information und Dokumentation;

URL: <http://www.iud.fh-darmstadt.de/iud/wwwmeth/lv/ss97/projekt/prot1.htm>, 1997

Koblenzer Jahresbericht 1996 des Instituts für Computerlinguistik

Das Institut für Computerlinguistik im Fachbereich Informatik der Universität Koblenz-Landau, Abteilung Koblenz, gibt auch 1996 wie in den Jahren zuvor einen Jahresrückblick heraus. Neu ist allerdings, daß nicht Papier verschickt wird, sondern eine e-mail-Benachrichtigung auf die entsprechende Adresse: <http://www.uni-koblenz.de/~compling/Allgemeines/Jahresberichte/jahresbericht.html> aufmerksam macht. Und auch die LeserInnen des LDV-Forum sind spätestens jetzt informiert...

Auch eine Nachricht

Verschiedene breitgestreute Anfragen haben ergeben, daß es gegenwärtig eine Nachrichtenflaute gibt. Es scheint schlichtweg keine aktuellen Nachrichten zu geben. Wobei die logische Konsistenzprüfung dieser Nachricht, die zu behaupten scheint, daß es keine Nachrichten gibt, hier Probleme angemeldet hat SPELLING CHECK - OK; GRAMMAR CHECK - OK; LOGICAL ERROR IN PARAGRAPH NO fvlfzuln LOGICAL ERROR IN PADH;:J ihjh LOGICAL EROH hff uzffvhkv> CONTINUE (Y,N)
n CANCELED.

**Testverfahren für intelligente
Indexierungs- und
Retrievalsysteme anhand
deutschsprachiger
sozialwissenschaftlicher
Fachinformation (GIRT)**

**Bericht über einen Workshop am IZ
Sozialwissenschaften, Bonn**

**12. September 1997, 10.30 Uhr bis
17.00 Uhr**

Gerhard Knorz

Der Workshop wurde gemeinsam vom Hochschulverband Informationswissenschaft, der Fachgruppe Information Retrieval der Gesellschaft für Informatik (GI-IR) und dem IZ Sozialwissenschaften veranstaltet.

Zwischen Forschung und Praxis liegen im Bereich des Information Retrieval 20 Jahre und mehr. Wenn wir die daraus resultierende Sprachlosigkeit zwischen Forschern und Praktiker überwinden wollen, wenn wir die Anwender mit Systemen unterstützen wollen, die das gegenwärtig Gewohnte in den Schatten stellen, dann brauchen wir überzeugende Belege und erfolgreiche erste Anwendungen. Dem Workshop GIRT kann man tatsächlich zutrauen, ein erster von weiteren folgenreichen Schritten in diese Richtung zu sein.

—Gerhard Knorz

1 Die Initiative „GIRT“

Seit ca. 2 Jahren wird im Hochschulverband für Informationswissenschaft (HI) und in der GI-Fachgruppe Information Retrieval über eine Initiative diskutiert, die mit dem Workshop GIRT nun eine erste Öffentlichkeit hergestellt hat. Mit dem Wechsel von Prof. Dr. Jürgen Krause von der Informationswissenschaft der Universität Regensburg an die Spitze des IZ Sozialwissenschaften in Bonn (und gleichzeitig an die Universität Koblenz an den Lehrstuhl der Software Ergonomie im Fachbereich Informatik) hat das IZ eine radikale informationstechnologische Umorientierung vollzogen, die nun auch kürzlich durch den Wissenschaftsrat ihre Bestätigung gefunden hat. Dieses Informationszentrum, das satzungsgemäß als Mitglied der GESIS (Gesellschaft Sozialwissenschaftlicher Infrastruktureinrichtungen) die sozialwissenschaftliche Forschung dokumentiert und unterstützt, hat es sich zum Ziel gesetzt, mit einer neu eingerichteten Forschungsabteilung, unterstützt durch Drittmittelprojekte und als am Markt tätiger Datenbankanbieter die technologisch-methodische Lücke zwischen informationswissenschaftlicher Forschung und karger Praxis zu schließen und sich mit eigenen Informationsangeboten an die „Spitze des Fortschritts“ (J. Krause) zu setzen. Daß gegenwärtig im Information Retrieval vieles in Bewegung gekommen ist, ist einerseits dem Internet zu verdanken und gleichzeitig einem zweifelsfrei erfolg- und einflußreichen offenen Wettbewerb für Retrievalsysteme in den USA. Dieser Wettbewerb, TREC (Text Retrieval Conference), setzt gegenwärtig die Maßstäbe für die Effektivität von Retrievalsystemen, weit über den eigentlichen Kreis der Teilnehmer aus Forschung und Industrie hinaus. NIST, der neutrale Veranstalter von TREC, stellt Dokumentsammlungen im Gigabyte-Bereich, Retrievalfragen, Wettbewerbs-

regeln und -organisation sowie Auswertungskapazität zur Verfügung – und die bisherigen Ergebnisse haben die in der Praxis dominanten Booleschen Retrievalverfahren zugunsten rankingbasierter Systeme klar aus dem Feld geschlagen.

Aus Sicht des IZ Sozialwissenschaften – wie auch sicher vieler anderer – haben die experimentellen Ergebnisse von TREC allerdings einen entscheidenden Nachteil: Erst die neuesten Ergebnisse der allerletzten Runde (TREC6, Herbst 1997) berücksichtigen (auch) die deutsche Sprache und außerdem sind – aus der Not geboren – die Dokumentsammlungen überwiegend auf Zeitungstexte konzentriert. Inwieweit also eine an englischsprachigen Zeitungen gewonnene Erkenntnis etwa auf eine deutschsprachige sozialwissenschaftliche Datenbank anwendbar ist, für die eine thesaurusbasierte intellektuelle Indexierung vorliegt, bleibt weiterhin eine Angelegenheit von Glauben und Spekulation. Schließlich zeigen gerade auch die TREC-Ergebnisse, wie sehr Fachgebiet, Topic und Textsorte in bisher nicht vorhersehbarer Weise auf die Retrievalqualität durchschlagen.

Genau hier setzt GIRT als ein Projekt des IZ Sozialwissenschaften an, das man vereinfacht als eine Art deutsches TREC bezeichnen könnte. Das IZ stellt – wie NIST – eine Testumgebung (Dokumente aus IZ-Datenbanken, Retrievalfragen, Aufbereitungs- und Auswertungskapazität) all denen zur Verfügung, die dieses Angebot zur Evaluierung ihrer Retrievalverfahren nutzen wollen. Die Vorteile haben alle Beteiligten auf *ihrer* Seite: Forschung und Entwicklung haben eine praxisrelevante Testumgebung (die bisher im deutschen Bereich fehlte), Forschung und Praxis haben vergleichbare Testergebnisse und das IZ Sozialwissenschaften hat aus erster Hand und auf eigenen Daten Belege dafür, welche Verfahren sich für seine Zwecke als überlegen herausstellen. Insbesondere kann es für seinen Fall die jahrzehntealte Streitfrage entscheiden, ob sich der Aufwand für Thesaurusentwicklung und intellektuelle Indexierung in besserer Retrievalleistung auszahlt, oder ob die manuelle Bearbeitung durch automatische Indexierungs- und Retrievalverfahren abgelöst werden kann bzw. soll.

Der nunmehr erste Workshop im September 1997 sollte die Idee und den Kontext von GIRT sowie die Ergebnisse eines Pretests mit den Systemen *freewais-sf* und *Messenger* vorstellen. Außerdem war es Ziel, GIRT mit der neuen Initiative von TREC hinsichtlich Mehrsprachigkeit unter Einschluß von Deutsch zu koordinieren.

Im gutgefüllten großen Sitzungssaal des IZ begann pünktlich die Veranstaltung.

2 Vorträge

2.1 Ziele und Perspektiven des Projektes GIRT (Krause)

In seinem kurzen Einleitungsreferat begründete und motivierte Jürgen Krause aus seiner Sicht als Initiator von GIRT das damit verbundene wissenschaftliche und praktische Interesse des IZ Sozialwissenschaften. Evaluierung von Retrievalverfahren als ein facettenreiches und in methodischer wie auch in praktischer Hinsicht höchst anspruchsvolles wissenschaftliches Problem ist ein Anliegen, das er aus Regensburg mit nach Bonn transferiert hat. Der Wechsel von der Hochschule zu einem Institut, das selbst als Informationsanbieter auftritt und den Zugang seiner Kunden zu seinen Datenbankanhalten selbst gestalten kann, stellt die Evaluierung in einen völlig neuen Kontext: Evaluierung nicht mit dem Ergebnis letztlich doch unverbindlicher Aussagen über Effektivität und Akzeptanz, sondern als konkreter Ausgangspunkt für die Gestaltung von Informationsdiensten des IZ. Diese Instrumentalisierung von Evaluierung fokussiert die Aufmerksamkeit auf Aspekte, die aus Nutzer- und Betreibersicht von besonderer Bedeutung sind: Dies sind zum einen Stellenwert und Rolle der Indexierung, die am IZ traditionell-professionell mittels eines selbstentwickelten Thesaurus intellektuell vorgenommen wird und zum andern das Interface zum Nutzer, das unbetritten die Effektivität der Recherche wesentlich bestimmt.

2.2 Generelle Ergebnisse der TREC-Studien, einschließlich TREC-5 (Womser-Hacker)

Frau Womser-Hacker, die sich im Rahmen ihrer im Februar dieses Jahres abgeschlossenen Habilitation intensiv mit den Ergebnissen von TREC auseinandergesetzt hat, versuchte den Stellenwert von TREC herauszuarbeiten, den Fortschritt an Entwicklung und Wissen zu charakterisieren und andererseits die Lücken aufzuzeigen, die aus dem Testdesign von TREC und der praktischen Ausgestaltung dieses Retrievalwettbewerbs folgen. Besonders interessant und glücklich für die Diskussion war es, daß mit Peter Schäuble von der ETH Zürich, einem der seit Anbeginn bei TREC beteiligten Teilnehmer (und Koordinator des sich

nun entwickelnden europäischen TREC-Ablegers), ein Insider zugegen war. Auf diese Weise gewann das Thema über die Systematik hinaus ein unmittelbares und authentisches Element, das einen Workshop immer aufwertet.

Zunächst stellte Womser-Hacker Evaluierung als ein Kernthema des Information Retrieval heraus, dessen Geschichte sie in 3 Phasen mit jeweils eigenem Erkenntnisgewinn und Entwicklungsstand einteilte:

- 1.) **1955–1980** mit den Cranfield-Tests und den Experimenten mit SMART und Medlars: In dieser Phase etablierte sich das Bewußtsein, daß Standardisierung und Wissenschaft, aber auch viel harte Arbeit für die Bewertung von Retrievalverfahren notwendig sind. Man lernte, das Mögliche im Bereich von 40%–60%, was Precision und Recall als die Standardmaße betrifft, einzuschätzen, und formulierte deren inverse Beziehung als empirisches Gesetz. Als offene Frage verblieb, wie sehr die Ergebnisse vom Kontext abhängig sind und warum die Systeme gerade so gut (oder schlecht) sind, wie sie sind. Und bereits in den Anfängen öffnete sich die so viel beklagte Kluft zwischen Praxis und Forschung.
- 2.) **1980–1990** als eine Phase mit speziellem Bezug zur deutschen Szene mit Projekten wie PADOK (Regensburg), AIR (Darmstadt) und LIVE (Berlin). Standardisierung im Testdesign und Operationalisierung von Systemeigenschaften kennzeichnen den Fokus. Das Bewußtsein für statistisch und meßtheoretisch „saubere“ Bedingungen und Ergebnisse hat sich entwickelt, ebenso wie auch für die Problematik der Konzepte von Relevanz, Aboutness und Semantik. Methoden zur Schätzung schwer meßbarer Parameter wurden entwickelt, die Mehrfachnutzung von Testkollektionen gewann an Bedeutung, die Testtypen Experiment und Untersuchung differenzierten sich aus und es zeigte sich, daß mit zunehmender Komplexität der Systeme die Evaluierung immer schwieriger wird.
- 3.) **ab 1990** dominieren TREC und seine Ausdifferenzierungen (special-tracks). GIRT soll sich hier mit einordnen.

TREC (Text Retrieval Conference) hat seit dem November 1992, dem Startpunkt, mit TREC 1 eine außerordentliche Ausstrahlung erreicht und ist zweifellos für das Gebiet des Information Retrieval ein großer Gewinn. Der aktuelle Stand veröffentlichter Ergebnisse ist gegenwärtig TREC 5. TREC 6 (mit dem Schwerpunkt Crosslingual Retrieval) hat 1997 bereits stattgefunden. Interessierte können alle Veröffentlichungen einsehen unter dem URL <http://www-nlpir.nist.gov/trec/>. Plakative Basisinformation über TREC, speziell auch über das Testdesign, finden sich als Folienfolge unter <http://www.iud.fh-darmstadt.de/iud/wwwmeth/publ/slide/owfrtr1.htm>.

TREC adressiert als Zielgruppe den „dedicated searcher“, unterscheidet grundsätzlich zwischen adhoc-Abfragen und Routing sowie zwischen automatischer, manueller und interaktiver Entwicklung der Suchanfrage. Grundsätzlich wird von einem „gerankten“ Ergebnis ausgegangen. Das Evaluierungsprogramm (und damit das Erkenntnisinteresse) entwickelt sich über die Zeit, so daß sich die Testbedingungen der einzelnen Runden unterscheiden und die Ergebnisse nur sehr eingeschränkt direkt vergleichen lassen. Spezifische Probleme (Filtering, fehlerhafter Input, Interaktivität, ...) werden in special tracks behandelt. Die Testkollektionen sind heterogen, konstruiert und „groß“ (Gigabyte-Bereich), vornehmlich bestehend aus Zeitungstexten. Die Fragekollektionen umfassen jeweils 50 Suchproblemstellungen. Die Teilnehmeranzahl wächst (TREC1 25 Systeme, TREC4 35 Systeme) und umfaßt Forschungsprototypen wie auch Systeme von kommerziellen Herstellern.

Die Ergebnisse der einzelnen TREC-Runden lassen sich recht gut in eine Entwicklungslinie einpassen:

- 1.) **TREC 1:** Fast alle Teilnehmer hatten mit Effizienzproblemen zu kämpfen, aber es zeigte sich, daß system- und organisationsseitig das TREC-Konzept praktikabel war.
- 2.) **TREC 2:** Die Ergebnisse waren deutlich verbessert. Ausschlaggebender Erfolgsfaktor war die Optimierung der Termgewichtungen.
- 3.) **TREC 3:** Systeme und Experimente wurden komplexer. Hybride Ansätze führten zu Ergebnisverbesserungen.
- 4.) **TREC 4:** Erfolgsfaktor war nunmehr die Einbeziehung linguistischer Analysen (Behandlung von Phrases, Frageerweiterung).

- 5.) **TREC 5:** Verbesserung der Resultate wurden erreicht, indem man das Problem der Längennormalisierung in den Griff bekam (Berücksichtigung der Tatsache, daß das Auftreten eines Suchwortes in einer 10-Zeilen-Nachricht einen anderen Stellenwert hat als in einem 10-Seiten-Aufsatz.)

Der Shift in den Erfolgsfaktoren ist z. T. auf einen Wechsel der Testbedingungen und damit der Anforderungen zurückzuführen: Insbesondere wurden die anfangs sehr ausdifferenzierten Frageformulierungen radikal gekürzt. Dadurch wurde die Aufgabe deutlich schwieriger mit dem Effekt, daß die Retrievaleffektivität signifikant zurückging und nunmehr das Problem der Frageerweiterung in den Vordergrund rückte.

Ein interessantes und bedenkenswertes Ergebnis im Hinblick auf die Frage nach der Übertragbarkeit der Resultate lieferte ein Signifikanztest bei den TREC 3-Ergebnissen: Die Varianz durch die verschiedenen Such-Topics erwies sich als höher als die Varianz verschiedener Retrievalsysteme.

Insgesamt liefert TREC bisher folgende Aussagen:

- TREC hat die Entwicklung von Retrievalverfahren enorm befruchtet. Die heutigen Verfahren, in Konkurrenz zu den früheren auf den „alten“ Testdaten, liefern deutlich verbesserte Ergebnisse.
- Viele Systeme erreichen denselben Leistungsstandard. Unterschiedliche statistische Modelle schlagen wenig auf die Effektivität durch, genauso wie elaborierte Modelle gegenüber einfachen Ansätzen keine signifikante Verbesserung nachweisen.
- Stark linguistisch orientierte Verarbeitung rechtfertigt den Mehraufwand nicht.
- Elaborierte manuelle Frageentwicklung ist keineswegs immer besser als einfache (automatische).
- Anwendung von Relevance Feedback-Techniken zeigt eindeutige Verbesserungen.
- Data Fusion ist eindeutig positiv.

- Der Aufwand in eine verbesserte Behandlung der Anfrage schlägt besser auf Effektivität durch als eine bessere Dokumentenanalyse.
- Mit TREC 5 änderte sich die vorher bestehende Situation, daß sich die Ergebnisdokumentmengen der einzelnen Systeme (sowohl was die relevanten, als auch was die nicht-relevanten Dokumente angeht) erstaunlich wenig überschneiden.

2.3 Ergebnisse des GIRT-Pretests (Kluck)

Michael Kluck, am IZ Sozialwissenschaften für die Durchführung von GIRT zuständig, erläuterte Test-Design und Ergebnisse des ersten Tests. Ziel war es im wesentlichen gewesen, die Praktikabilität des Testkonzeptes zu erproben und Erfahrungen mit Organisation und Ablauf zu sammeln. Schließlich lassen sich im Vorfeld die für den praktischen Aufwand wichtigen Fragen nach der durchschnittlichen Dauer einer Recherche von Testpersonen, nach deren praktischen Problemen oder nach der zu erwartenden Anzahl von Trefferdokumenten nur unter großer Unsicherheit prognostizieren.

Der Pretest wurde von zwei Systemen bestritten: *freewais-sf*, eine einfache, verbreitete Suchengine im Internet mit statistischem Ranking und einfacher regelbasierter Morphologie sowie *Messenger*, das (konventionelle) Standardsystem für eine Suche nach manuell indextierten IZ-Dokumenten, in einer kategorisierten bibliographischen Datenbank mit Freitextsuche

Die Testkollektion umfaßt ca. 15.000 Dokumente der beiden sozialwissenschaftlichen Datenbanken SOLIS und FORS mit Titel und Abstract. Für eine kleine Untermenge können auch Volltexte zur Verfügung gestellt werden. 9 Anfrageprobleme waren als Standardtestfragen vorgegeben und die zugehörige Relevanzbeurteilung war durch einen IZ-Juror bereits im voraus durch exhaustive Recherchen vorgenommen worden. Grundsätzlich sollen auch Spezialfragen möglich sein, die von den Testteilnehmern eingebracht werden. Voraussetzung ist, daß sie dokumentiert und veröffentlicht sind. Die Relevanzbeurteilung findet dann im Rahmen der Testauswertung statt.

Die Testpersonen waren 8 informationswissenschaftlich vorgebildete Personen mit und ohne Erfahrung im Online-Retrieval. Jede Person bearbeitete alle 9 Fragen mit beiden Systemen (mit wechselnder Reihenfolge), wobei ein Zeitbudget von ca. 2–4 Stunden benötigt wurde. Um technische Probleme und überhaupt die Interface-Problematik aus dem Test auszuklammern, stand als technische Vermittlung ein professioneller Rechercheur als Bediener der Tastatur zur Verfügung.

Eine Festlegung, die sich als problematisch herausgestellt hat, war die Begrenzung der Antwortmengen auf 30 Treffer. Diese Restriktion sollte den Aufwand kalkulierbar machen, hat aber in unzulässiger Weise das Suchverhalten und die Ergebnisse beeinflusst.

Die Relevanz von Dokumenten wurde auf einer 4-stufigen Skala bewertet und für Precision-Recall-Auswertungen auf eine binäre Relevanzentscheidung abgebildet. Die Konsistenz von Urteilen unterschiedlicher Juroren lag bei 70–80%.

Die Ergebnisse des Pretests liegen detailliert als Precision/Recall-Diagramme vor und sollen aufgrund der geringen Anzahl von Fragen sowie der unzumutbaren Begrenzung der Antwortmengen nicht überinterpretiert werden. Dennoch ergeben sich eine Reihe interessanter Beobachtungen:

- Die Überschneidungen der Antwortmengen liegen insgesamt auf sehr niedrigem Niveau: nur 21% der relevanten Dokumente sind in mehr als einer einzigen Recherche gefunden worden, und das bei jeweils $8 * 2 = 16$ Recherchen/Frage!
- Es gibt keinen klaren Anhaltspunkt dafür, daß eines der Systeme dem anderen überlegen ist. Sollte sich dieses anhand von 9 Fragen und 72 Recherchen gewonnene Ergebnis bestätigen, so bedeutet dies, daß die manuelle Indexierung vom Kunden nicht mit Gewinn genutzt wird. Der Nutzer hat allerdings eine klare Meinung darüber, mit welchem System er besser zurechtgekommen ist und die besseren Ergebnisse erzielt hat: nämlich mit dem Booleschen System. Dieser subjektive Eindruck widerspricht der objektiven Beobachtung.
- Eine detaillierte Analyse der verwendeten Frageformulierungen ergibt, daß vielfach zentrale Aspekte der vorgegebenen Suchprobleme unberücksichtigt geblieben sind. Es drängt sich der Eindruck auf, daß ein simples Abschreiben der vorgegebenen Suchformulierung bei *freewais-sf* zu besseren Ergebnissen hätte

führen können als das Ergebnis eigenen Nachdenkens. Daß niemand überhaupt auf diese Idee gekommen ist, zeigt, daß die Testpersonen mit dem Konzept statistisch basierten Retrievals nicht vertraut waren.

- Den Testpersonen fehlte auch weitgehend das Bewußtsein, daß Boolesches Retrieval unter Verwendung von Deskriptoren die Orientierung im Vokabular nahelegt. Nur wenige Testpersonen mit Retrievalerfahrung, aber kein Laie (!) haben den bereitliegenden Thesaurus zur Suche benutzt.
- Rechercheure der IZ Sozialwissenschaften erzielen deutlich bessere Rechercheergebnisse als die Testpersonen. Hier schneidet *Messenger* als das gewohnte und die Vorgehensweise prägende Recherchesystem klar besser ab als *freewais-sf*.

2.4 Multilingualität in TREC (Schäuble)

Hintergrund für den Vortrag von Peter Schäuble von der ETH Zürich war die Neuerung im Rahmen von TREC 6 (1997), CrossLanguage Information Retrieval (CLIR) mit den Sprachen Deutsch, Englisch, Französisch und prinzipiell auch Italienisch in die Evaluierung aufzunehmen. Zur Vorbereitung dieses neuen Schwerpunktes war eine mit Europäern und US-Forschern besetzte Arbeitsgruppe ins Leben gerufen worden, deren europäischer Sprecher Peter Schäuble ist. In Weiterentwicklung dieser Aktivitäten ist für 1998 ein European TREC mit Unterstützung von CEPIS, GI, BCS und NIST geplant. Und genau dies war der konkrete Ansatzpunkt für die nachfolgende Schlußdiskussion und für das Interesse, konkrete gemeinsame Perspektiven zu entwickeln.

Wer nun bei eher strategisch/politischem Hintergrund einen tendenziell deskriptiven, farblosen Vortrag erwartet hätte, der hatte sich absolut verrechnet. Schäuble vermittelte in konzentrierter Form einen ausgezeichneten Überblick über die gegenwärtige Landschaft der Ansätze für ein sprachgrenzenüberschreitendes Information Retrieval. Die grundsätzliche Idee der einzelnen erfolgreichen Konzepte wurde gut nachvollziehbar dargestellt. Die Teilnehmer des Workshops

waren – dem an lebendiger Diskussion gezeigten Interesse zufolge – beeindruckt von der – man kann es so sagen – Raffinesse des an der ETH Zürich verfolgten Ansatzes. Und vor allen Dingen von den vorgestellten Ergebnissen!

Zunächst einmal sind die möglichen Anwendungen von Crosslingual IR eines kurzen Nachdenkens wert: Für mehrsprachige Länder (für jemanden aus Zürich ein Heimspiel!) und Organisationen liegt der Bedarf auf der Hand. Darüber hinaus bietet das Web mit seiner englischen Standardsprache das Beispiel für einen weiten Nutzerkreis mit großem passivem aber eher kleinem aktiven Wortschatz in einer Fremdsprache. In professionellem Kontext kann man aber auch an monosprachliche Nutzer denken, für die eine maschinelle Rohübersetzung eines vorher gefundenen fremdsprachlichen Dokumentes zumindest eine Relevanzentscheidung ermöglicht. Eine besonders interessante Anwendung stellt die Recherche von Bildern auf der Basis von Bildbeschreibungen dar, deren Originalsprache den Suchenden gar nicht interessieren muß.

Nachdem die Probleme „naiver Ansätze“ diese als untauglich disqualifiziert hatten, gab Schäuble einen kurzen Abriss der Geschichte einschlägiger wissenschaftlicher Arbeiten, beginnen bei Salton 1970 (wie könnte es anders sein?) und eine systematische Klassifikation prinzipiell möglicher Retrievalansätze, die Sprachgrenze zu überschreiten. (Eine entsprechende graphische Darstellung findet sich unter <http://www.iud.fh-darmstadt.de/iud/wwwmeth/publ/slide/clirva1.htm>.)

Auf zwei sehr unterschiedliche Ansätze ging Schäuble im Detail ein: zum einen auf „Latent Semantic Indexing“ (LSI) und auf den an der ETH Zürich verfolgten Ansatz, der einen Ähnlichkeitsthesaurus berechnet und für das Retrieval anwendet. Hinter LSI steckt die Idee, daß man Dokumente und Terms jeweils auf einen gemeinsamen abstrakten Vektorraum abbilden kann, der sehr viel weniger Dimensionen hat als etwa der ursprüngliche Dokumentenraum des Vektormodells. Dieser reduzierte abstrakte Raum ist, natürlich, sprachunabhängig.

Ausführlicher will ich den Ansatz unter Einbeziehung eines Ähnlichkeitsthesaurus nachzeichnen: Schäuble lag viel daran, herauszustellen, daß Dokumentenräume und Termräume sich dual verhalten. In einem Dokumentenraum beschreibt man Dokumente über die Terms (das Vokabular), die sie enthalten und in einem Termraum werden Terms über die Dokumente spezifiziert, in denen sie vorkommen. Jede wahre Aussage über Zusammenhänge im Dokumentenraum läßt sich mechanisch in eine entsprechende wahre Aussage im korrespondierendem Termraum transformieren. Mit dieser interessanten Sichtweise läßt sich neu begründen, was im Information Retrieval eine lange und keineswegs flächendeck-

kend erfolgreiche Geschichte hinter sich hat: die Assoziationsfaktoren oder anders benannt: der Ähnlichkeitsthesaurus. Terms werden demnach genau dann zueinander in Beziehung gesetzt, wenn sie häufig gemeinsam in Dokumenten auftreten.

Ein Ähnlichkeitsthesaurus kann zunächst zur Frageerweiterung eingesetzt werden: Die Suchterms werden um weitere Terms aus dem Ähnlichkeitsthesaurus angereichert. Zu einem CLIR-Ansatz wird dieses Verfahren dann, wenn man eine zweisprachige Dokumentenkollektion mit korrespondierenden Dokumenten zur Verfügung hat: Man berechnet dann die Ähnlichkeit eines Terms zu den Terms in jeweils anderssprachigen Dokumenten. So kann man etwa eine deutsche Anfrage über den Ähnlichkeitsthesaurus in eine Menge gewichteter französischer Suchterms abbilden, mit denen man dann in französischsprachigen Dokumenten recherchiert. Wichtig zu wissen ist, daß die französischen Suchterms keinesfalls Übersetzungen sein müssen bzw. sein sollten. Vielmehr geht es (nur) darum, solche französische Wörter zu finden, wie sie in Dokumenten auftreten, in deren deutschen Übersetzung das deutsche Suchwort vorkommt. Damit ist das Problem deutlich einfacher, als es z. B. in einem Übersetzungssystem gelöst werden muß.

Nun erscheinen Dokumentenkollektionen, in denen alle Dokumente in zwei Sprachen verfaßt sind, sehr selten. Die Züricher sind dennoch fündig geworden: eine erste Anwendung lieferte eine Dokumentensammlung zum Schweizer „systematischem Recht“ in Deutsch und Französisch. Anhand von 53.000 Dokumenten wurde ein Ähnlichkeitsthesaurus entwickelt, der anschließend u. a. erfolgreich für die deutschsprachige Recherche in französischsprachigen Gerichtsurteilen (Wechsel der Textsorte!) angewandt wurde. Um die detaillierten Evaluierungsergebnisse zusammenzufassen: Gegenüber einer monosprachigen einfachen Recherche *ohne* Ähnlichkeitsthesaurus bringt die Crosslingual-Recherche *mit* Ähnlichkeitsthesaurus eine deutliche Qualitätssteigerung, und letztlich kann man sagen, daß der Sprachwechsel dem Verfahren nur einen Effektivitätsverlust unter 5% kostet

Der besondere Clou ist den Zürichern allerdings damit gelungen, daß sie das Verfahren auch ohne vorliegende zweisprachige Dokumentenkollektion implementieren konnten. Sie stellten eine solche Sammlung „einfach“ selbst her: Sie filterten Agenturmeldungen eines sehr langen Zeitraums (mehr als 10 Jahre) und glichen verschiedensprachige Meldungen anhand einfacher Kriterien wie Datum, Klassifikation, gemeinsames Auftreten sprachunabhängiger Zeichenfolgen (Namen) ab. Daß die so ermittelten Dokumentenpaare nicht notwendigerweise

direkte Übersetzungen voneinander sind, muß aus praktischer und theoretischer Sicht nicht stören. Den damit berechneten Ähnlichkeitsthesaurus setzte die ETH Zürich dann in den Experimenten für TREC 6 ein. Wie sehr man umdenken muß, wenn man versucht, das Potential eines solchen Verfahrens einzuschätzen, machte Schäuble an einem Exempel plastisch deutlich: Nimmt man etwa den Schweizer Parlamentsthesaurus in die Hand, so hat man es mit 2 cm bedruckten Seiten zu tun. In gleicher Form ausgedruckt ergibt der bei TREC6 eingesetzte Ähnlichkeitsthesaurus einen 70 km hohen Papierstapel!

3 Abschlußdiskussion und Resümee

Die Abschlußdiskussion, von Gerhard Knorz (FH Darmstadt) moderiert, sollte im wesentlichen 2 Fragen klären:

1. Gibt es Interesse an und Anforderungen für das Angebot, das GIRT all denjenigen macht, die ein deutschsprachiges Retrieval anbieten und die dieses Retrieval evaluieren wollen?
2. Lassen sich die Initiativen für GIRT und für ein europäisches TREC in einen gemeinsamen Rahmen einbetten?

Der erste Punkt wurde insofern nur kurz behandelt, als übereinstimmend festgestellt wurde, daß eine Evaluierungsumgebung, wie sie GIRT anbietet, ein Desiderat für jeden darstellt, der Retrievalverfahren für deutschsprachige Texte entwickelt bzw. optimiert. Gegenüber den Planungen weitergehende Anforderungen wurden nicht formuliert.

Der zweite Punkt war im Verlauf des Workshops bereits häufiger andiskutiert worden. Im wesentlichen wurde zwischen Krause als dem Ausrichter von GIRT und Schäuble als dem europäischen Sprecher der Arbeitsgruppe für ein europäisches TREC geklärt, daß das IZ Sozialwissenschaften unbestritten und willkommen die neutrale Rolle für ein europäisches TREC spielen kann, wie sie NIST in den USA für TREC innehat. Daß die Verwaltung weiterer Dokumenten- und Fragekollektionen mit dem damit verbundenen Auswertungsaufwand nicht mit den

Mitteln des IZ getragen werden kann, wird klar akzeptiert und sollte kein Problem darstellen, da eine Projektförderung für ein solches Vorhaben sicher erreichbar scheint.

So hatte der Workshop allen Teilnehmern eine Menge zu bieten gehabt: Interessante Informationen, lebhafte Diskussionen und ein sehr vielversprechendes praktisches Ergebnis. Der Aufschwung des Gebietes „Information Retrieval“ setzt sich fort!

KONVENS 98

Computer, Linguistik und Phonetik zwischen Sprache und Sprechen

4. Konferenz zur Verarbeitung natürlicher Sprache

5.–7. Oktober 1998, Universität Bonn

Veranstaltet von:

Gesellschaft für Linguistische Datenverarbeitung (1998 federführend)

Deutsche Gesellschaft für Sprachwissenschaft

Gesellschaft für Informatik, FA 1.3.1 „Natürlichsprachliche Systeme“

Informationstechnische Gesellschaft/Deutsche Gesellschaft für Akustik

Österreichische Gesellschaft für Artificial Intelligence

Call for Papers

Thema der Tagung sind alle Bereiche der maschinellen Verarbeitung von Sprache in geschriebener und gesprochener Form.

Besondere Aufmerksamkeit soll dabei solchen Ansätzen zuteil werden, die sich mit den strukturellen und phonologisch/phonetischen Aspekten der computerunterstützten Sprachforschung befassen und dazu beitragen, eine Brücke zwischen diesen beiden Aspekten zu schlagen.

Es wird hiermit aufgefordert, Vorschläge einzureichen für:

Einzelvorträge, Workshops, Systemvorführungen und Postervorführungen.

Die eingereichten Vorschläge und Beiträge werden anonym begutachtet. Aus diesem Grund muß jedem Vorschlag ein Deckblatt beigelegt werden, auf dem Name und Institution der Verfasserin bzw. des Verfassers sowie Titel und Art des Beitrags angegeben sind.

Beiträge sollen in fünf Papier-Exemplaren, DIN A4, Times 12 pt, sowie in elektronischer Form per e-mail (vorzugsweise LaTeX oder PostScript) eingereicht werden. Der Umfang der Beiträge soll bei Vorträgen 10 Seiten, Posterbeiträgen 4 Seiten, Workshop- Anmeldungen 5 Seiten nicht überschreiten.

Workshop-Anmeldungen sollen die Bedeutung des Themas begründen und einen Überblick über die voraussichtlichen Teilnehmer und deren Beiträge enthalten. Systemvorführungen müssen kurz beschrieben werden, ferner sind die benötigten Geräte zu nennen.

Allen Beiträgen soll Zusammenfassungen auf Deutsch und Englisch von jeweils maximal 12 Zeilen beigelegt sein. Tagungssprachen sind Deutsch und Englisch.

Alle Einreichungen werden von mindestens zwei unabhängigen Gutachtern beurteilt, die vom Programmkomitee bestimmt werden. Angenommene Beiträge werden in einem Tagungsband rechtzeitig zu Beginn der Tagung gedruckt vorliegen.

Termine

- 1.3.98** Einreichen von Workshop-Vorschlägen
- 15.4.98** Einreichen von Beiträgen und Posterbeiträgen
- 15.5.98** Benachrichtigung über Annahme bzw. Ablehnungen
- 15.6.98** Abgabe der druckfertigen Vorlagen für den Tagungsband
- 15.7.98** Anmeldung von Systemvorführungen

Örtliche Organisation

Prof. Dr. Wolfgang Hess

Prof. Dr. Winfried Lenders (<http://www.ikp.uni-bonn.de/~wle>)

Dr. Thomas Portele (<http://www.ikp.uni-bonn.de/~tpo>)

Dr. Bernhard Schröder (<http://www.ikp.uni-bonn.de/~bsh>)

Programmkomitee

Dr. Ernst Buchberger, Wien (ÖGAI)

Prof. Dr. Dafydd Gibbon, Bielefeld (DGfS)

Prof. Dr. Roland Hausser, Erlangen (GLDV)

Prof. Dr. Wolfgang Hess, Bonn (ITG/DGA)

Prof. Dr. R. Hoffmann, Dresden (ITG/DGA)

Dr. Tibor Kiss, Heidelberg (DGfS)

Prof. Dr. Winfried Lenders, Bonn (GLDV)

Dr. Harald Trost (ÖGAI)

Prof. Dr. Wolfgang Hoepfner, Duisburg (GI)

Dr. Stefan Busemann, Saarbrücken (GI)

Tagungsbüro

Gisela von Neffe

Institut für Kommunikationsforschung und Phonetik

der Universität Bonn

Poppelsdorfer Allee 47

D-53115 Bonn

Internet: <http://www.ikp.uni-bonn.de/Konvens98/index.html>
(Seite entspricht diesem Artikel)

E-Mail: konvens98@uni-bonn.de
Tel.: (0228) 735638
Fax: (0228) 735639

Informationen zur online-Anmeldung finden Sie am Ende dieses Artikels.

Tagungsort

Die KONVENS 98 findet im Hauptgebäude der Universität Bonn statt. Das Hauptgebäude befindet sich im Stadtzentrum, in Geknähne zum Hauptbahnhof.

Unterkünfte werden durch die Stadt Bonn vermittelt. Wenden Sie sich an das Amt für Wirtschaftsförderung und Tourismus, Tel.: (0228) 773920; E-Mail: Amt_03@mail.bonn.de.

Voraussichtliches Zeitraster der Tagung

Mo, 5. Oktober 1998

8.00–14.00	Uhr	Registrierung
8.00–12.30	Uhr	Tutorien
14.00	Uhr	Eröffnung
14.00–18.00	Uhr	Sektionen, Workshops
19.00	Uhr	Empfang

Di, 6. Oktober 1998

9.00–11.00	Uhr	Sektionen, Workshops
11.00	Uhr	öffentlicher Vortrag
14.00–15.30	Uhr	Sektionen, Workshops
16.00–17.30	Uhr	Panel
17.30–19.00	Uhr	Mitgliederversammlungen der Fachgesellschaften

Mi, 7. Oktober 1998

9.00–11.00	Uhr	Sektionen, Workshops
11.00	Uhr	öffentlicher Vortrag
14.00–15.30	Uhr	Sektionen, Workshops
15.30–16.00	Uhr	Abschlußsitzung

Poster- und Systemvorführungen ab Di., 6. Oktober am Tagungsort (Hauptgebäude der Universität) und ggf. im Institut für Kommunikationsforschung und Phonetik (IKP).

Anmeldung

1. Teilnahmegebühr

Bei Registrierung bis zum 15.6.1998: **160,- DM**, Studenten (Ausweis mitschicken!) **80,-DM**; bei Registrierung ab dem 15.6.1998: **200,- DM**, Studenten (Ausweis mitschicken!) **100,-DM**.

Die Teilnahmegebühren sind gleichzeitig mit der Anmeldung einzuzahlen auf folgendes Konto:

Volksbank Bonn,
BLZ 38060186,
Konto-Nummer 502199013 (Winfried Lenders).

2. Online-Anmeldung

Unter dem folgenden URL können Sie sich online für die KONVENS 98 anmelden: <http://www.ikp.uni-bonn.de/Konvens98/form.html>.

Ergebnisse der Wahlen 1997 zu Vorstand und Beirat

Die folgenden Ergebnisse der Wahlen 1997 zu Vorstand und Beirat waren der letzten Ausgabe des LDV-Forum, soweit sie an GLDV-Mitglieder verschickt wurden, bereits als briefliche Benachrichtigung beigelegt. Hier nochmal – offiziell – der Text des Schreibens des Wahlvorstands vom 13. Juni 1997, nunmehr im redaktionellen Teil.

Die Auszählung der Wahlscheine für die Vorstandswahl und Beiratswahl fand am 13. Juni 1997 in Saarbrücken statt. Die Wahlbeteiligung war nicht sehr rege, es wurden nur 93 Wahlscheine eingesandt, die alle gültig waren. Der Wahlvorstand, bestehend aus Johann Haller, Gregor Thurmair und Nico Weber, hat folgende Ergebnisse ermittelt:

Ergebnisse der Vorstandswahl

Kandidat/in	Ja	Nein	Enthaltung
<i>R. Hauser (1. Vorsitz)</i>	80	10	3
<i>D. Rösner (2. Vorsitz)</i>	87	4	2
<i>U. Seewald (Schatzmeisterin)</i>	92	0	1
<i>C. Wolff (Schriftführer)</i>	86	4	3
<i>B. Schröder (Informationsreferent)</i>	90	1	2

Ergebnisse der Beiratswahl

Kandidat	Stimmen	
<i>G. Knorz</i>	72	gewählt
<i>W. Lenders</i>	65	gewählt
<i>U. Schmitz</i>	62	gewählt
<i>W. Höppner</i>	59	gewählt
<i>H. Elsen</i>	41	–
<i>G. Heyer</i>	35	–

Chronik der GLDV

Winfried Lenders

–Alle Leser dieses ersten Versuchs, die Geschichte der GLDV darzustellen, sind aufgefordert, Korrekturen und Ergänzungen mitzuteilen (Lenders@uni-Bonn.de).

Gründung

1975 Gründung der GLDV unter dem Namen LDV-Fittings in München. Seitdem ist München Sitz der Gesellschaft.

Vorstände der GLDV seit 1975

(Jeweils 1. Vorsitzender, stellvertretender. Vorsitzender, Schatzmeister, Schriftführer):

1. Vorstand 1975, München: Tillmann, Krallmann, Schweisthal, Lutz
2. Vorstand 1976, Bonn/Mannheim: Krallmann, Zimmermann, Schae-
der, Lutz
3. Vorstand 1978, Essen: wegen Beschlußfähigkeit erfolgten die
Wahlen im Herbst auf einer weiteren Sitzung in München:
Krallmann, Zimmerman, Schae-der, Endres-Niggemeyer.
1979: Jahrestagung Bonn: Rücktritt Zimmermanns,
Krause wird stellv. Vorsitzender
4. Vorstand 1976, 1980, Saarbrücken: Krause, Hellwig, Schae-der,
Schmidt
5. Vorstand 1982, Koblenz: Krause, Hellwig, Schae-der, Kroupa
6. Vorstand 1983, Trier: Krause, Hellwig, Schae-der, Kroupa
Auf der MGV in Trier am 3.3.1983 wird der Name in „*Gesellschaft
für Linguistische Datenverarbeitung*“ geändert.
7. Vorstand 1985, Hannover: Endres-Niggemeyer, Hellwig,
Schae-der, Schneider
8. Vorstand 1987, Bonn/Frankfurt: Endres-Niggemeyer, Schweisthal,
Schae-der, Schneider
9. Vorstand 1989, Ulm: Rieger, Klenk, Schae-der, Schneider.
Das Amt des Informationsreferenten wird eingeführt. Erster IR:
Haller
10. Vorstand 1991, Trier: Rieger, Klenk, Schae-der, Schneider, Haller
11. Vorstand 1993, Kiel: Lenders, Haller, Hausser, Seewald, Hitzenber-
ger
12. Vorstand 1995, Regensburg: Lenders, Haller, Hausser, Seewald,
Hitzenberger
13. Vorstand 1997, Leipzig: Hausser, Rösner, Seewald, Wolff, Schröder

Jahrestagungen

Nr	Jahr	Ort	Zeit	Organisator	Thema
1	1976	München	210./11. Dez. 1976	Bannerjee/Reinhart, Fa. Siemens	Zur Lage der LDV in der BRD
2	1978	Essen	22.-24. Feb.1978	Krallmann	Zur Lage der LDV in der BRD
3	1979	Bonn	12.-14. Dez. 1979	Lenders	Dialogsysteme und Textverarbeitung zusammen mit ALLC
4	1980	Saarbrücken	11./12. Dez. 1980	Zimmermann	LDV - Ausbildung und Berufsperspektiven
5	1982	Koblenz	17.-19. Feb. 1982	Bátori	LDV und Nachbarn
6	1983	Trier	1.-4. März 1983	Niederehe	LDV-Kolloquium
7	1984	Heidelberg	29.Feb.-2.März 1984	Hellwig	Trends Linguistischer DV in Verbindung mit dem Arbeitskreis "Mikrocomputer und Textverarbeitung"
8	1985	Hannover	5.-7.März 1985	Endres-Niggemeyer	Sprachverarbeitung in Information und Dokumentation
9	1986	Göttingen	25.-27.2.1986	Klenk, Scherber, Thaller	LDV in den Geisteswissenschaften
10	1987	Bonn	4.-6. März 1987	Willie, Tillmann	Analyse und Synthese gesprochener Sprache
11	1988	Saarbrücken	9.-11. März 1988		Computerlinguistik und ihre theoretischen Grundlagen zusammen mit DGfS
<i>- fortgesetzt auf der folgenden Seite -</i>					

<i>– fortgesetzt von der vorhergehenden Seite –</i>					
<i>Nr</i>	<i>Jahr</i>	<i>Ort</i>	<i>Zeit</i>	<i>Organisator</i>	<i>Thema</i>
12	1989	Ulm	8.-10.März 1989	Rösner	Interaktion und Kommunikation mit dem Computer
13	1990	Siegen	28.-30. März 1990	Schaeder	Lexikon und Lexikographie
14	1991	Trier	20. Sept. 1991	Rieger, Köhler	zusammen mit QUALICO
15	1993	Kiel		Pötz	Sprachtechnologie: Methoden, Werkzeuge und Perspektiven
16	1995	Regensburg		Hitzenberger	Angewandte Computerlinguistik
17	1997	Leipzig	17.-19. März 1997	Heyer	Linguistik und neue Medien

–Seit 1992 wird die GLDV-Jahrestagung alle zwei Jahre abgehalten, im Wechsel mit der Konvens.

Ja zur Sinnlichkeit – von gedrucktem Papier

Mitglieder wollen auf Printversion des LDV-Forums nicht verzichten

Zusammen mit der Wahl zum Vorstand und Beirat der Gesellschaft für Linguistische Datenverarbeitung (GLDV) wurde eine Mitgliederbefragung durchgeführt. Es ging um die Frage, ob eine Papierversion des LDV-Forum im Zeitalter des Internet noch zeitgemäß und notwendig ist. Immerhin fordern Druck und Versand des Heftes einen nicht zu vernachlässigenden Anteil der finanziellen Ressourcen des Vereins.

Die Mitglieder-Datenbank, geführt am IKP, Bonn enthielt 281 Adressen. Alle diese Mitglieder haben die Wahlunterlagen und den Stimmzettel zur Mitgliederbefragung erhalten. Nicht mehr als 93 Mitglieder haben sich dann an Wahl und Befragung beteiligt, und von diesen 93 haben sich 78 für den Erhalt der Papierversion des Forums entschieden.

Bleibt die Frage, ob die verbliebenen stummen 188 Mitglieder für Papier einfach nicht mehr ansprechbar sind (und aus diesem Grund die mit gelber Post zugesandten Wahlunterlagen nicht mehr wahrgenommen haben), oder ob einfach der Kern minimal aktiver Mitglieder so klein ist und dieser Kern für die Sinnlichkeit bedruckten Papiers in hohem Maße empfänglich ist? Wie es auch sei – das LDV-Forum mutiert nicht in eine virtuelle Existenz, sondern verbleibt bis auf weiteres real in Ihren Händen.

GK.

Geisteswissenschaftliche Hypermedia-Anwendungen

Bericht über den Workshop des Arbeitskreises Hypermedia am 20. Juni 1997 in Bonn

Mit sehr unterschiedlichen Zielsetzungen werden hypermediale Techniken im geisteswissenschaftlichen Umfeld eingesetzt. Dieser Gedanke jedenfalls legte es nahe, in einem eintägigen Workshop-Programm den Versuch zu unternehmen, zunächst zueinander bezuglos erscheinende Ansätze und Überlegungen nebeneinanderzustellen und nach Schnittstellen zu suchen. So unterschiedlich die Adressatinnen und Adressaten der angestrebten oder verfügbaren Hypermedia-Produkte und ihre Benutzungsweise sein mögen, das geisteswissenschaftliche Umfeld teilen sie; und die Produkte haben eben jene Eigenschaft gemein, die sie unter den sonstigen Medien-Produkten zu Hypermedia adelt. Gibt das bereits genügend Basis für einen Austausch zwischen den Entwicklerinnen und Entwicklern auch sehr unterschiedlicher Anwendungen?

Beteiligung und Interesse an diesem und an den vorangegangenen Workshops des Arbeitskreises Hypermedia (Leitung: *Angelika Storrer* [Mannheim], *Bernhard Schröder* [Bonn]) sprechen dafür, daß diese Frage von vielen Entwicklerinnen und Entwicklern klar bejaht wird und ein Austausch als fruchtbar empfunden wird. Der Bonner Workshop des Arbeitskreises am 20. Juni 1997 führte Beiträge zu hypermedialen Texteditionen, hypermedialen Wörterbüchern und Lernsoftware für Linguistinnen und Linguisten zusammen.

Carl-Martin Bunz (Saarbrücken) stellte das **TITUS**-Projekt vor, das die Erstellung eines *Thesaurus indogermanischer Text- und Sprachmaterialien* betreibt. Der Thesaurus soll historisch-vergleichenden Sprachforschern die Sprachdenk-

mäler der relevanten indogermanischen Sprachen in verschiedenen maschinell-verarbeitbaren Formen zugänglich machen. Dazu werden die Texte derzeit in verschiedenen Formaten bereitgestellt: ASCII- und HTML-Texte können plattformunabhängig verwendet werden, WordPerfect- und WordCruncher-für-Windows-Formate richten sich an die PC-Benutzerinnen und -Benutzer. Wo immer möglich, wird angestrebt, Abbildungen der Sprachdenkmäler, Transliterationen/Transkriptionen mit evtl. Varianten im Sinne einer kritischen Textedition zusammenzuführen. Die Integration von Abbildungen kann in HTML und WordCruncher für Windows durch die üblichen Hyperlink-Mechanismen geschehen. Die Arbeit mit Textvarianten, Transkriptions-/Transliterationsvarianten und Textübersetzungen und Kommentierungen/Annotierungen wird durch WordCruncher für Windows bei entsprechender Textvorbereitung durch die Technik des *synchronous scrolling* unterstützt: Aufeinander bezogene Texte in verschiedenen Fenstern werden automatisch parallel zum Text im aktiven Fenster positioniert.

Größere technische Probleme bereiten noch die historischen Schriftsysteme und die für ihre Transliteration/Transkription notwendigen diakritischen Erweiterungen des lateinischen Alphabets. Doch hier bietet der TITUS-Server durch die Bereitstellung von Windows-TrueType-Schriften, die für eine adäquate Darstellung sorgen, sehr weitgehende Abhilfe. Ein Grund für den Einsatz des Texterschließungssystems WordCruncher für Windows im Kontext dieser Texte ist seine Unterstützung fremdsprachiger Schriftkonventionen: Links- und rechtsläufige Schriften lassen sich in Texten mischen, für verschiedene Sprachen können auch in demselben Dokument unterschiedliche Tastaturdefinitionen wirksam sein, Sortierreihenfolgen und Akzentvarianten von Buchstaben sind sprachspezifisch zu bestimmen.

Wünschbar wäre es selbstverständlich, die Kodierung der Schriftzeichen zu standardisieren, damit die Texte unabhängig von speziellen Fonts und Systemen austauschbar bleiben. Bunz berichtete von seinen Aktivitäten bei Unicode, die auf die Integration von Schriftzeichen, die für die historisch-vergleichende Sprachforschung von Belang sind, in die Standardisierungsbemühungen zielen.

Weitere Informationen zum TITUS-Projekt bietet der TITUS-Server <http://titus.uni-frankfurt.de>. Eine ausführlichere Darstellung des TITUS-Projektes von Jost Gippert findet sich auch im LDV-Forum 12(1). Die meisten Texte des Thesaurus sind aus urheberrechtlichen Gründen derzeit nur für Teilnehmer des Projektes verfügbar.

Ebenfalls eine auf WordCruncher für Windows basierende Textedition kam in dem Beitrag von *Winfried Lenders* (Bonn) zur Sprache. Am Beispiel einer multimedialen Edition der Werke **Hartmanns von der Aue**, die Lenders in Zusammenarbeit mit Prof. Dr. Kurt Gärtner (Universität Trier) und Prof. Dr. Roy Boggs (University of South Florida) erarbeitet, stellte er den Mehrwert der hypermedialen Aufbereitung sowohl für den philologisch als auch den linguistisch interessierten Leser altdeutscher Texte heraus. Eine elektronische Ausgabe kann verschiedenartigen Bedürfnissen zugleich gerecht werden: Sie kann als Synopse verschiedener Textvarianten und unterschiedlicher Transliterationen und Übersetzungen dienen und zusätzlich qualitativ hochwertige Abbildungen der Originaldokumente anbieten. Sie kann gleichzeitig die einzelnen Varianten selbst in lesbarer Form präsentieren, als es die Buchform vermag, wenn sie zugleich den Anforderungen an eine kritische Edition genügen muß.

Gewiß gibt es stärker formatierungsorientierte Werkzeuge als WordCruncher für Windows, und diese halten für manche Layout-Fragen befriedigendere Antworten als letzteres Programm bereit. Aufgrund der oben schon im Zusammenhang mit TITUS beschriebenen Mechanismen der Hyperlinks und des *synchronous scrolling* gehören jedoch zu den Stärken von WordCruncher für Windows seine Möglichkeiten, mehrere Text- und Bildebenen zueinander in Beziehung zu setzen. Doch lassen sich diese Darstellungsmöglichkeiten auch zur Integration weiterer linguistischer Informationen nutzen: Der Hartmann-Text wird mit einer morphologischen Analyse aller Vorkommen von Wortformen parallelisiert, so daß linguistisch interessierte Benutzerinnen oder Benutzer diese Information für jede Textstelle abrufen können. Da auch die linguistische Markierung wie ein fortlaufender Text repräsentiert ist, ist diese ebenfalls mit dem WordCruncher-Search-Manager durchsuchbar. Über die vielfältigen Optionen der Kollokationssuche kann nach morphologischen oder morphosyntaktischen Konfigurationen gesucht werden.

Während es bei Textausgaben eher die metatextuellen Informationen sind, die den Einsatz hypermedialer Techniken sinnvoll machen, können es bei Wörterbüchern schon die ausdrücklichen intratextuellen Verweisstrukturen des Wörterbuchtextes selbst sein. Dies erläuterte *Ingrid Lemberg* (Heidelberg) sehr anschaulich am Beispiel des **Deutschen Rechtswörterbuchs**, dem in den 1896/97 von der Königlich Preußischen Akademie der Wissenschaften begonnenen und seit 1959 von der Heidelberger Akademie der Wissenschaften weitergeführten Wörterbuch der älteren deutschen Rechtssprache. Neben der gebundenen Ausgabe wird auch eine elektronische Bereitstellung der Informationen angestrebt. Durch den Ein-

Translingua

Language & Technology

Wir suchen ab sofort Diplom-Informatiker/innen oder Computerlinguist/inn/en für einen jeweils ca. einjährigen Einsatz bei einem Kunden in den USA als

Localizer

Notwendige Voraussetzungen

- Abschluß als Diplom-Informatiker/in oder Computerlinguist/in (durch Zeugnis nachweisbar)
- Gute Kenntnisse in Windows 95/NT auf Betriebssystemebene
- Internet-Kenntnisse (HTML, TCP/IP, CGI, etc.)
- Fähigkeit, Probleme selbständig zu lösen
- Gutes Englisch und Deutsch in Wort und Schrift

Wir bieten eine marktgerechte Vergütung und die Erledigung der Visumsangelegenheiten.

Zudem bieten wir Studenten aus dem Bereich Computerlinguistik (Schwerpunkt maschinell unterstützte Übersetzung) die Möglichkeit zu einem berufsorientierten

Praktikum

Die Mindestdauer des vergüteten Praktikums beträgt 3 Monate.

Notwendige Voraussetzungen

- Immatrikulation in einen Studiengang aus dem Bereich Computerlinguistik
- Gute Kenntnisse im Problembereich maschinelle Übersetzung (bes. Translation-Memory-Systeme)
- Gute Kenntnisse in Windows 95
- Fähigkeit, Probleme selbständig zu lösen

Bitte bewerben Sie sich mit Ihren aussagekräftigen Unterlagen bei:

Translingua
Language & Technology

z. H. Dr. Ursula Marmé
(Manager Resources and Planning)
Ludwig-Erhard-Allee 3
D-53117 Bonn
Tel.: +49.228 – 8160.117

satz hypermedialer Verweistechiken läßt sich der Nutzwert eines Wörterbuchs gegenüber der gedruckten Form deutlich erhöhen. Da sind zum einen die Verweise innerhalb von Wörterbuchartikeln auf andere Artikel durch Angabe von entsprechenden Verweiswörtern, die elektronisch als hypertextuelle Links auf die betreffenden Artikel oder Teilartikel im Wörterbuch realisiert werden können. Hier besteht der Wert der elektronischen Realisierung in der Bereitstellung von bequem zu handhabenden Navigationsmitteln für den Benutzer. In einem historischen Wörterbuch sind zum anderen zahlreiche Verweise auf Belegstellen zu finden, teils mit Zitierung des Belegs, teils ohne ein solches Zitat nur mit einer Fundstellenangabe. Hier kann die elektronische Ausgabe die Belege in größerem Umfang und mit weiteren Kontexten bereitstellen, als ein gedrucktes Lexikon dies vermag. In vielen Fällen wird auch ein Faksimile der Originalquellen wichtige außersprachliche Informationen bereithalten. Extratextuelle Links zu rechtshistorischen Bilddatenbanken erhöhen den bei einem fachsprachlichen Wörterbuch wesentlichen sachlichen Informationswert um einiges. Neben Angaben zum Wörterbuchprojekt finden sich unter dem URL <http://www.uni-heidelberg.de/institute/sonst/adw/drw/index.html> zur Demonstration auch einige Artikel in ihrer Internet-Form.

Die Vorträge des Workshop beschloß der Beitrag von *Elisabeth Cölfen* (Essen) über die am FB 3.6 der Universität Essen entwickelte und projektierte **Multimediale Interaktive Lernsoftware für Linguistik-Studentinnen und -Studenten im Grundstudium**. Von einem sehr ansprechenden und selbsterklärenden Design wird die Benutzerin oder der Benutzer in die Arbeit mit den Lernpäckchen, derzeit zu Bühler und Saussure, eingeführt. Das Äußere hat hier nicht nur werbenden Charakter, sondern soll der speziellen Benutzungssituation und dem didaktischen Zweck angepaßt sein. Als Konsequenz weicht es stark von den gewohnten graphischen Benutzeroberflächen ab: Das Interface zeigt in Form von Symbolen, die sich in ein graphisches Ganzes ohne harte Abgrenzungen integrieren, und durch Text offen, was das Programm in einer bestimmten Situation an Optionen zu bieten hat und verzichtet auf geschachtelte und damit zunächst versteckte Untermenüs, die nur mit Mausakrobatik hervorzuzaubern sind.

Die Texte der Lernpäckchen werden von Abbildungen begleitet, wo sie den Text ergänzen können oder zur assoziativen oder informativen Einordnung des Textes dienen. Eine besondere Problemstellung ergibt sich aus dem Entwurf von Navigationsmöglichkeiten, die sowohl verschiedenen Lernbedürfnissen als auch

Nachschlagebedürfnissen gerecht werden. Unter der Überschrift *Kuntermund & Löwenmaul* gibt der Linguistik-Server Essen <http://www.linse.uni-essen.de> Einblick in die aktuellen Lernpäckchen.

Den größten Teil des Nachmittags nahm der Besuch bei dem Bonner Übersetzungs-Unternehmen **Translingua** und der Schwesterfirma **tops.net** ein, zu dem *Harald Elsen* den Arbeitskreis eingeladen hatte. Der Arbeitsschwerpunkt der *Translingua Gesellschaft für Dokumentation und Software-Lokalisierung mbH* liegt in der Übersetzung und zielsprachlichen Anpassung von Software-Dokumentationen. Bei der Übersetzung von On-line-Dokumentationen kann natürlich auch deren hypertextueller Aufbau nicht außer acht gelassen werden. Breiten Raum nahm eine Diskussion über die maschinelle Unterstützung von Übersetzung ein und die Qualifikation von Absolventen computerlinguistischer Studiengänge für das Beschäftigungsfeld maschinell unterstützter Übersetzung. Auf besonderes Interesse stießen in der Gesprächsrunde auch die Verfahren zur Qualitätssicherung von Übersetzungen.

Bei *tops.net online publishing services* wurde von *Lucie Prinz* ein sehr aktueller Einblick in die kommerzielle Wirklichkeit der Internet-Angebote gewährt. Die Firma ist sowohl im Bereich der Bereitstellung von Modem- oder ISDN-Zugängen zum Internet tätig als auch beim Hosting, beim Entwurf und bei der Realisierung von Web-Sites. Die Fragen und Diskussionsbeiträge der Teilnehmerinnen und Teilnehmer kehrten immer wieder zur Sprache des Web zurück. Jedoch konnte Prinz berichten, daß eine sprachliche Kreativität, die auf eine Ersetzung der Anglizismen oder Pseudo-Anglizismen ziele, bei den Kunden nicht gefragt sei.

Der nächste Workshop des Arbeitskreises soll 1998 in Heidelberg stattfinden und unter dem Thema *Hypertextualisierung von Wörterbüchern* stehen.

Bernhard Schröder, Bonn