# $\begin{array}{c} {\rm Magisterarbeit} \\ {\rm im} \\ {\rm Studiengang~Computerlinguistik} \end{array}$

# Grammatik der Menschenbezeichner in biographischen Kontexten

Michaela Geierhos

März 2006

Betreuer der Arbeit: Prof. Dr. Franz Guenthner

#### Vorwort

Wer muss heute noch eine Biographie oder einen Lebenslauf veröffentlichen, wenn er oder sie eine in der Öffentlichkeit präsente Person ist? - Eigentlich betrifft das niemand dieser Leute, denn Ausschnitte ihres Lebens werden in den Printmedien von verschiedensten Blickwinkeln beleuchtet. Natürlich veröffentlichen die wenigsten Zeitungen oder Magazine lückenlose Lebensläufe prominenter Menschen.

Liest man nur einen Artikel zu der betreffenden Person, so erfährt man nur wenig über sie und bekommt auch recht einseitige Informationen. Doch lässt man Google nach diesen bekannten Leuten suchen, so bekommt man eine Vielzahl von Artikelreferenzen, welche die unterschiedlichsten Facetten und Bereiche ihres öffentlichen und privaten Lebens beleuchten. Kurz und prägnant werden einem Informationen über den Familienstand, die Familienverhältnisse, den sozialen Status, das geschätzte oder bekannte Jahreeinkommen, sowie die Vorlieben und Freizeitaktivitäten und noch vieles mehr auf dem "silbernen Tablett" serviert. Die Fülle an Informationen, die Google ihren Kunden bietet, übersteigt oft ihre anfänglichen Erwartungen. Manchmal erfährt man sogar Details aus dem Leben der Reichen und Schönen, die derjenige selbst wohl nie so veröffentlicht hätte.

Wie CNET News.com [64] am 14. Juli 2005 berichtete, erging es dem Google CEO Eric Schmidt nicht anders. Obwohl er selbst auf seiner Homepage wenig über seine Person preisgibt, findet Google nach kürzester Zeit alle wichtigen Daten, die seine Person betreffen.

#### Google CEO Eric Schmidt doesn't reveal much about himself on his home page.

But spending 30 minutes on the Google search engine lets one discover that Schmidt, 50, was worth an estimated \$1.5 billion last year. Earlier this year, he pulled in almost \$90 million from sales of Google stock and made at least another \$50 million selling shares in the past two months as the stock leaped to more than \$300 a share.

He and his wife Wendy live in the affluent town of Atherton, Calif., where, at a \$10,000-a-plate political fund-raiser five years ago, presidential candidate Al Gore and his wife Tipper danced as Elton John belted out "Bennie and the Jets".

Schmidt has also roamed the desert at the Burning Man art festival in Nevada, and is an avid amateur pilot.  $^{1}$ 

<sup>&</sup>lt;sup>1</sup>Ausschnitt aus dem Artikel "Google balances privacy, reach" von Elinor Mills [64]

Aber warum sollte man 30 Minuten bei der Suche nach "Google CEO Eric Schmidt" damit verbringen, die einzelnen Treffer nach der biographisch relevanten Information zu durchsuchen? Wäre es nicht sinnvoller, wenn bei der Suche nach Personen auch ein Fokus auf die verschiedenen Beziehungen gelegt wird, in denen ein Mensch mit anderen Menschen, mit einer Firma, mit Wohn- und Arbeitsorten oder zeitlichen Begebenheiten in Verbindung steht? Würde es eine personenbezogene Suche nicht enorm erleichtern, wenn einer der allerersten Treffer, Aufschluss über das Alter bzw. das Geburtsdatum, dann evtl. über den Familienstand, gefolgt vom Beruf oder dem momentanen Beschäftigungsverhältnis geben würde?

Damit will ich andeuten, dass eine Staffelung nach Wichtigkeit der biographischen Daten und das Ausfiltern von biographisch irrelevanter Information die Zufriedenheit des Benutzers bei der Suche deutlich erhöhen kann.

Doch bevor man eine Skala für die Relevanz von biographischen Relationen festlegen kann, muss man sich ein Bild davon machen, welche Prädikate überhaupt dafür in Frage kommen. Denn einerseits sollten es u.a. Verbrelationen sein, die sehr häufig in Biographien auftauchen und andererseits auch interessant für den Informationssuchenden sein.

Deshalb möchte ich im Rahmen dieser Magisterarbeit versuchen, eine umfassende Definition von Prädikaten zu geben, welche in biographischen Kontexten auftreten können und essentielle Informationen über die betreffende Person geben. Dabei verstehe ich unter einer "Definition von Prädikaten" nicht nur eine reine Auflistung dieser, sondern vielmehr die Erstellung einer Grammatik - eines Regelwerkes - welche analysiert und gleichzeitig vorgibt, wie sich ein bestimmtes Verb innerhalb eines Satzgefüges verhält, d.h. welche Argumente das Verb zwingend oder optional hat, oder ob es oft im Zusammenhang mit Lokativa oder Temporalia auftritt. Natürlich ist diese syntaktische, aber auch semantische, Betrachtungsweise von personenbezogenen Sätzen sprachabhängig. Aufgrund der Vielzahl an möglichen Prädikaten, werde ich mich in meiner Arbeit ausschließlich auf eine Sprache, nämlich das Englische, beschränken.

Dabei ist mir besonders wichtig, dass der Schwerpunkt dieser Arbeit nicht das Ranking von biographischen Relationen oder eine Fallstudie ist, wie eine gute, automatisch generierte Textzusammenfassung einer Biographie auszusehen hätte, sondern vielmehr möchte ich das Augenmerk auf die Analyse von biographisch relevanten Sätzen richten. Nun wird es nicht bei einer reinen syntaktischen Studie von Satzgefügen bleiben, da die natürliche Sprache sehr viele Paraphrasierungsmöglichkeiten bietet. Das macht u.a. ein semantisches Clustering von Relationstypen - die Bildung von Synonymklassen auf der Ebene der Prädikate - aber auch eine Typisierung von Satzteilen notwendig. Letzteres macht bereits der Titel dieser Magisterarbeit deutlich, denn "Menschenbezeichner" sind bereits eine eigene Klasse, worunter u.a. Eigennamen für Personen, wie z.B. "Bill Clinton", Berufsbezeichnungen, wie z.B. "Software Engineer", oder Bezeichnungen für Verwandtschaftsverhältnisse, wie z.B. "mother, aunt, grandfather", fallen.

Aber bevor ich in Details gehe, sollte das als kleiner Vorgeschmack auf diese Arbeit ausreichend sein und natürlich hoffe ich, dass ich Ihr Interesse dafür geweckt habe.

München, den 27. März 2006

## **Danksagung**

An dieser Stelle möchte ich allen danken, die mich bei der Erstellung dieser Arbeit unterstützt haben.

Mein besonderer Dank geht an Prof. Dr. Franz Guenthner, für den ich von Beginn meines Studiums an, als studentische Hilfskraft tätig war. Dadurch bekam ich schon früh in meinem Studium einen praktischen Bezug zu fachspezifischen Themen, wie der Informationsextraktion aus Internet Ressourcen, der Named-Entity-Recognition, lexikalischer Semantik, Lexikographie und lokalen Grammatiken.

Auch die Anregung seinerseits sich mit dem Verhalten von Menschenbezeichnern in biographischen Kontexten im Rahmen meiner Magisterarbeit zu beschäftigen, war im Grunde nur eine Fortsetzung meiner bisherigen Arbeit.

In zahlreichen Gesprächen konnte ich Probleme meine Magisterarbeit betreffend mit Herrn Prof. Guenthner besprechen und bekam zahlreiche hilfreiche Ratschläge und Anregungen für mein Thema. Vorallem bin ich ihm sehr dankbar dafür, dass ich jederzeit zu ihm kommen konnte, wenn ich Fragen hatte, und er stets viel Geduld bewies. Keineswegs ist zu unterschätzen wie sehr ich mit Ressourcen - sei es in Form von Eigennamenlisten für die Lexikographie oder der Bereitstellung einer Jahresausgabe Finacial Times als Trainingskorpus - zusätzlich zu meinen aus dem Internet selbst extrahierten Ressourcen unterstützt wurde.

Des Weiteren möchte ich Herrn Sebastian Nagel für seine Unterstützung bei der Satzenderkennung im Englischen danken. Durch die Erweiterung seines Tokenizer-Programms um eine umfangreiche Abkürzungsliste, war es mir möglich, eine fast fehlerfreie Satzenderkennung auf der Zeitungstextsammlung durchzuführen.

Genau wie Frau Sandra Bsiri beriet er mich im Umgang mit dem System Unitex und gab mir nützliche Tipps zur Wörterbuchkodierung, zur Verarbeitung von großen Korpora und zur Erstellung von Graphen.

Zudem war Frau Bsiri bereit, das von ihr entworfene benutzerfreundliche Web-Interface zum Testen von Unitex-Graphen auf meine Bedürfnisse zuzuschneiden, so dass ich zeitweilig ihre Weboberfläche nutzen konnte.

Durch die Vorarbeit, die Frau Dr. Friederike Mallchok bei der Erkennung von Organisationsnamen in Wirtschaftsnachrichten geleistet hat, konnte ich meine Magisterarbeit darauf aufbauen und auf ihre Erfahrungen entsprechend zurückgreifen.

Abschließend gilt mein Dank all denen, die mir Anregungen beim Verfassen dieser Arbeit gegeben oder mich auf Sekundärliteratur bezüglich dieses Themas hingewiesen haben.

## Inhaltsverzeichnis

Vo	Vorwort					
Danksagung				4		
1	NEF	NER innerhalb biographischer Relationen in Nachrichten				
	1.1	Begrif	fsklärung: Named Entity Recognition (NER)	10		
	1.2	Defini	tion: Biographische Relationen	12		
		1.2.1	Persönliche Relationen	12		
		1.2.2	Öffentliche Relationen	13		
		1.2.3	Zufällige Relationen	13		
	1.3	Einsch	nätzung der Thematik	14		
		1.3.1	Probleme und Chancen	14		
		1.3.2	Bewältigung der Aufgabe	15		
2	Loka	okale Grammatiken				
	2.1	1 Was sind lokale Grammatiken?				
	2.2	Warum werden lokale Grammatiken verwendet?				
	2.3	.3 Unitex - Ein System zur Anwendung lokaler Grammatiken				
		2.3.1	Allgemeines	19		
		2.3.2	Textvorverarbeitung	19		
			2.3.2.1 Normalisierung	20		
			2.3.2.2 Satzenderkennung und Auflösung von Kontraktionen	20		
			2.3.2.3 Tokenisierung	20		
			2.3.2.4 Lexikalische Analyse	20		
		2.3.3	DELA Wörterbücher	20		
		2.3.4	Prioritäten bei der Anwendung der Lexika	23		
		2.3.5	Mustererkennung und Konkordanzen	23		
3	Zusa	usammenfassung früherer Arbeiten				
	3.1		ce Gross	24		
		3.1.1	Zur Person	24		
		3.1.2	Bootstrapping bei der Entwicklung lokaler Grammatiken	25		
		3.1.3	Lemmatisierung zusammengesetzter Zeiten im Englischen	26		
	3.2		Senellart	29		
		3.2.1	Zur Person	29		
		3.2.2	Erkennung von Eigennamen und Berufsbezeichnungen	29		

			3.2.2.1 Motivation des Ansatzes	29				
			3.2.2.2 Vergleich mit anderen Information Retrieval Methoden . 3	30				
				32				
			3.2.2.4 Einblick in die formalen Beschreibungsmethoden 3	32				
			3.2.2.5 Schwächen des Ansatzes	34				
		3.2.3	Bootstrapping zur Erkennung von Nominalphrasen mit FSTs 3	35				
			3.2.3.1 Die Vorgehensweise	35				
			3.2.3.2 Das Ergebnis	35				
	3.3	Natha	lie Friburger	36				
		3.3.1	Zur Person	36				
		3.3.2	Erkennung von Eigennamen in Zeitungstexten	36				
			3.3.2.1 Kaskadierung von Transduktoren	36				
			3.3.2.2 Eigennamen bei der Klassifikation von Nachrichtentexten 3	39				
	3.4	Friede	rike Mallchok	10				
		3.4.1	Zur Person	10				
		3.4.2	Erkennung von Organisationsnamen in Wirtschaftsnachrichten 4	10				
			3.4.2.1 Motivation des Ansatzes	10				
			3.4.2.2 Einsatz von Ressourcen: Korpora und Lexika 4	11				
			3.4.2.3 Entwicklung lokaler Grammatiken	11				
			3.4.2.4 Bootstrapping und Akronymbildung	12				
			3.4.2.5 Fazit der Arbeit	12				
				_				
4		Beschränkungen im System 4.1 Sprachgebundenheit						
	4.1	1 0						
	4.2		1	14				
	4.3	Priori	sierung von Entitäten	15				
5	Res	sourcer	n: Grundlagen des Systems 4	ŀ6				
_	5.1		5 ,	16				
		5.1.1		16				
		5.1.2		17				
	5.2			18				
		5.2.1		19				
		5.2.2		50				
		5.2.3		51				
				52				
			0 1 0	53				
				53				
			•	54				
				55				
		5.2.4		56				
		J. <b></b> 1		57				
				, . 57				
				, . 57				
		5.2.5		, 58				

			5.2.5.1	Allgemeine Menschenbezeichnungen aus WordNet	. 58	
			5.2.5.2	Berufsbezeichnungen	. 59	
			5.2.5.3	Einwohnerbezeichnungen	. 60	
		5.2.6		der personenbezogenen Prädikate	. 62	
		5.2.7	Lexika d	er Branchen	. 63	
			5.2.7.1	Fachbereiche und Lehrfächer	. 63	
			5.2.7.2	Sektoren- und Branchenbezeichnungen	. 64	
		5.2.8	Lexika d	er Organisationsnamen		
			5.2.8.1	Allgemeine Organisationsbezeichnungen	. 65	
			5.2.8.2	Eigennamen von Organisationen	. 65	
			5.2.8.3	Organisationsspezifische Adjektive	. 67	
			5.2.8.4	Organisationsspezifische Kontexte		
		5.2.9	Lexika d	er geographischen Begriffe		
			5.2.9.1	Kontinente und Länderbezeichungen	. 69	
			5.2.9.2	Städtenamen		
			5.2.9.3	Grafschaften, Regionen, Bezirke und Départements	. 70	
			5.2.9.4	US Bundesstaaten und ihre typischen Abkürzungen		
			5.2.9.5	Geographische Adjektive	. 71	
		5.2.10	Lexika d	er Temporalia		
				Monatsnamen und -abkürzungen		
			5.2.10.2	Wochentage und ihre Abkürzungen	. 73	
				Weitere zeitbezogene Nomina		
		5.2.11	Weitere	Lexika	. 74	
	5.3			glichkeiten bei Google		
6	Grammatik der Menschenbezeichner 76					
			rsonennamen			
	0.1	6.1.1		sche Variabilität bei Personennamen		
		0.1.1	6.1.1.1	Abkürzung vs. Vollform		
			6.1.1.2			
		6.1.2		guierung von "Scheinnamen"		
		6.1.3		Eändigung des Personennamenlexikons		
	6.2			schenbezeichner		
	6.3	_	ern auflös			
7	Cua		day Oya	aniaatiananaman	86	
7			_	anisationsnamen		
	7.1			ariabilität bei Organisationsnamen		
	7.2 Abgrenzung von unechten Organisationsnamen			<u> </u>		
	7.3	vervol	ıstandıgu	ng des Organisationsnamenlexikons	. 90	
8	Grammatik der Ortsangaben 91					
	8.1	_		Relationen mit Ortsangaben		
	8.2		0	ihrer Funktion als Attribute		
		8.2.1		ne als Attribut einer Berufsbezeichnung		
		8.2.2	Toponyn	ne als Attribut eines Organisationsnamens	. 92	

9	Grammatik der Datumsangaben					
10	) Grammatik persönlicher Relationen	105				
	10.1 Die Geburt: "to be born"	. 107				
	10.2 Die Kindheit: "to be raised (up)"	. 111				
	10.3 Der Schulabschluss: "to graduate"					
	10.4 Die Heirat: "to be married"					
	10.5 Die Scheidung: "to be divorced"					
	10.6 Der Tod: "to die"					
11	Grammatik beruflicher Relationen	130				
	11.1 Der Beginn eines Beschäftigungsverhältnisses					
	11.1.1 Die Ernennung: "to be appointed as"					
	11.1.1.1 Grammatik der Berufsbezeichnungen					
	11.1.1.2 Vervollständigung des Sektorenlexikons					
	11.1.2 Die Einstellung: "to employ so."					
	11.1.3 Der Firmeneintritt: "to join"					
	11.2 Die Ausübung des Berufes					
	11.2.1 Das Beschäftigungsverhältnis: "to be employed"					
	11.2.1 Das Beschaftigungsverhaftins: "to be employed					
	11.2.3 Die Tätigkeit: "to work as"					
	11.3 Das Ende eines Arbeitsverhältnisses					
	11.3.1 Die Entlassung: "to dismiss so." bzw. "to be dismissed"					
	11.3.2 Die Nachfolge: "to be replaced as"					
	11.3.3 Die Abdankung: "to resign as"					
12	2 Auswertung der Ergebnisse	158				
	12.1 Evaluationsmaße [97]					
	12.1.1 Precision bzw. Genauigkeit					
	12.1.2 Recall bzw. Vollständigkeit					
	12.1.3 Fall-Out					
	12.2 Qualität des Systems	. 159				
13	Anwendungen					
	13.1Extraktion von Relationen zwischen Personen und Organisationen	. 160				
	13.2 Extraktion von Relationen zwischen mindestens zwei Personen	. 161				
	13.3 Extraktion von Relationen zwischen Personen und ihren Berufen	. 162				
14	Zusammenfassung und Ausblick	163				
Α	Hilfe zum Tokenizer-Programm von Sebastian Nagel	164				
В	Übersicht aller Kategorien in den Wörterbüchern	166				
_	B.1 Semantische Kategorien					
	B.2 Grammatikalische Kategorien					

С	Syntaktische Variabilität am Beispiel von "Bill Gates"	168	
D	Weitere Berufsbezeichnergraphen	171	
Ε	Graphen und Konkordanz zur Verbalphrase "to be dismissed"	174	
Literaturverzeichnis			
Abbildungsverzeichnis			
Tabellenverzeichnis			
Index		190	
Le	Lebenslauf		
Erklärung zur Magisterarbeit			

## 1 NER innerhalb biographischer Relationen in Wirtschaftsnachrichten

#### 1.1 Begriffsklärung: Named Entity Recognition (NER)

Für die Computerlinguistik hat sich die Named-Entity-Recognition inzwischen zu einem der wichtigsten Forschungsgebiete entwickelt. Wer schon einmal den Begriff "Named-Entity-Recognition" gehört hat, weiß dass er im Bereich der Informationsextraktion (IE) anzusiedeln ist. Im deutschen Sprachraum ist die NER auch unter dem Schlagwort "Eigennamenerkennung" bekannt.

Doch ist die Erkennung von Eigennamen nicht die einzige Aufgabe der Informationsextraktion, bei der versucht wird, aus Texten nicht-ambige Daten, die ein festgelegtes Format haben, zu extrahieren (Roth, 2002 [79]). Die **Eigennamenerkennung** ist nur eine von verschiedenen Teilaufgaben der IE, wobei sie eigenständig auftreten kann, oder wiederum Teil einer anderen computerlinguistischen Anwendung sein kann. Oft sind Information-Retrieval-Systeme und Systeme zur Antwort-Extraktion, Textzusammenfassung oder maschinelle Übersetzung, sowie Textmining-Programme und Suchmaschinen auf die Dienste der Name-Entity-Recognition (NER) angewiesen (Roth, 2002 [79]).

Leider gehen die Meinungen, was man genau unter Named-Entity-Recognition zu verstehen hat, auseinander. Der Streitpunkt bei der Definitionsfindung bezieht sich hierbei auf die Klärung des Begriffs der benannten Entität (Named Entity).

Ohne sich darauf festzulegen, was unter einer benannten Entität verstanden werden soll, lässt sich zunächst folgende Definition geben:

Named-Entity-Recognition bezeichnet die automatische Erkennung von Instanzen und Einheiten bestimmter Klassen in Texten.

Natürlich gibt diese Begriffserklärung keinerlei Aufschluss darüber, welche "Klassen" von Entitäten nun bei der NER erkannt werden sollen.

Für Massimiliano Ciaramita und Yasemin Altun [9] ist der Begriff der *Named Entity* recht eng gefasst, indem sie nur Personen, Organisationen und Orte berücksichtigen, definieren sie:

**Named-entity recognition (NER)** is the task of tagging words with labels such as person, organization, and location.

Jedoch werden meist weitere Entitäten, wie Datums-, Zeit-, Prozent- und Währungsangaben mit in die Erkennung von Eigennamen einbezogen. Named entity recognition (NER) (...) seeks to locate and classify atomic elements in text into predefined categories such as the names of persons, organizations, locations, expressions of times, quantities, monetary values, percentages, etc. <sup>2</sup>

Weiterhin bleibt fragwürdig, ob Datums-, Zeit-, Prozent- und Währungsangaben wirklich in die Kategorie der "benannten Entitäten" fallen, oder zwar Entitäten, aber keine Eigennamen sind.

Zur Klärung dieser Frage trägt wohl die auf der MUC-7<sup>3</sup> festgelegte Definition zur Erkennung von "Names Entities" bei.

The Named Entity task consists of three subtasks (entity names, temporal expressions, number expressions). The expressions to be annotated are 'unique identifiers' of entities (organizations, persons, locations), times (dates, times), and quantities (monetary values, percentages).<sup>4</sup>

Auf ein ähnliches Ergebnis kommt man, wenn man vom deutschen Begriff für *Named Entity* - dem "Eigennamen" - ausgeht.

Denn **Eigennamen** können danach kategorisiert werden, welche Art von Objekt sie bezeichnen:  $^5$ 

- Die häufigsten Namensträger sind Personen. Bei Personennamen kann man Vornamen und Familiennamen unterscheiden.
- Eine weitere große Gruppe bilden die Ortsnamen (Toponyme). Diese können weiter untergliedert werden in Städtenamen, Ländernamen, Flussnamen, Flurnamen usw.
- Institutionen sind typischerweise Träger von Eigennamen.
- Eine weitere große Gruppe bilden die Produktnamen.

Als **Eigennamen** werden also Bezeichner für Personen, Orte, Organisationen und Produkte betrachtet. Datums-, Zeit-, Prozent- und Währungsangaben gelten zwar als Entitäten, aber nicht als *Named Entities* - und somit deren Bezeichner auch nicht als Eigennamen (nach Roth, 2002 [79]).

Im Grunde gibt es zwei Arten von NER Systemen: Die eine Gruppe verwendet linguistische Methoden und die andere baut auf statistischen Modellen auf.

Für den hier vorgestellten Ansatz der automatischen Erkennung von Menschenbezeichnern innerhalb biographischer Relationen in Wirtschaftsnachrichten soll nur die sprachbasierte Named-Entity-Recognition von Interesse sein.

Hierbei stellen die Personen die wichtigste Entitätsart bei der Eigennamenerkennung dar. Weitere Entitäten, die typischerweise in Wirtschaftstexten auftreten, wie Organisationen, Orte und Zeitangaben werden mit in die Suche einbezogen.

<sup>&</sup>lt;sup>2</sup>http://en.wikipedia.org/wiki/Named\_entity\_recognition

<sup>&</sup>lt;sup>3</sup>Unter der MUC-7 versteht man die im Jahre 1998 zum 7. Mal durchgeführte Message Understanding Conference.

<sup>&</sup>lt;sup>4</sup>MUC-7 Named Entity Task Definition ([79] Seite 7)

 $<sup>^5\</sup>mathrm{vgl}$ . http://de.wikipedia.org/wiki/Eigenname

#### 1.2 Definition: Biographische Relationen

Bis jetzt wurde die Aufgabenstellung der automatischen Erkennung von Menschenbezeichnern innerhalb biographischer Kontexte in Wirtschaftsnachrichten nur in dem Punkt erläutert, welche Entität in den Texten gefunden werden soll.

Inzwischen ist deutlich geworden, dass sich dieser Ansatz auf die Lokalisierung von Eigennamen, insbesondere von Personennamen als "Named Entity", konzentriert. Dennoch werden nicht nur Namen für Personen betrachtet, sondern auch andere Menschenbezeichnungen, die sich auf den Beruf, die soziale Stellung oder Verwandtschaftsverhältnisse beziehen.

Nun bleibt nur noch zu klären, was unter "biographischen Relationen" verstanden werden soll. **Biographische Relationen** sind in der Regel Prädikatrelationen von Verben, die vornehmlich in biographischen Kontexten auftreten. Da in biographischen Kontexten die Lebensgeschichte von Personen beschrieben wird, betrifft es Verben, die das Ereignis der Geburt, den schulischen und beruflichen Werdegang, sowie Beziehungen zu anderen Menschen usw. beschreiben.

Diese Prädikate lassen sich aufgrund ihrer Relevanz für die Öffentlichkeit in verschiedene Kategorien unterteilen. Einerseits gibt es eine Gruppe von Verben, die in fast jeder Biographie genannt werden und andererseits gibt es Verben bzw. Relationen, die nur in Autobiographien zur Sprache kommen. Dies ermöglicht eine Aufspaltung der biographischen Relationen in die drei Unterkategorien der persönlichen, der öffentlichen und der zufälligen Relationen.

#### 1.2.1 Persönliche Relationen

Laut Duden ist eine Biographie nichts anderes als die Niederschrift einer Lebensgeschichte. Somit ist es nicht verwunderlich, wenn manche Leute vieles aus ihrem Leben zu erzählen haben. Dabei werden oft Details aus dem Gefühlsleben Preis gegeben, und es wird "aus dem Nähkästchen geplaudert", wie es in bestimmten Autobiographien der Fall ist. In der Regel werden in diesem Zusammenhang sehr intime Dinge über Personen erzählt, welche für die Öffentlichkeit eigentlich nicht bestimmt sein sollten.

In **persönlichen Relationen** sind besonders solche Verben anzutreffen, die Gefühlsregungen ausdrücken und Informationen aus dem Privatleben der Leute liefern. Doch sind es meist Relationen, die jemanden persönlich betreffen.

Natürlich stellt sich bei diesen Prädikatrelationen nun die Frage, inwiefern sie noch biographische Relevanz haben. An dieser Stelle muss man wohl einräumen, dass persönliche Relationen in ihrer ersten Bedeutung zwar die schönsten Klatschgeschichten aus dem Leben berichten, und somit sicher als biographische Relation gezählt werden können, aber in Wirtschaftsnachrichten kaum Beachtung finden. Deshalb sind sie für den hier vorgestellten Ansatz nahezu irrelevant.

Dennoch ist die Grenze zwischen persönlichen und öffentlichen Relationen manchmal fließend. Ein solcher "Grenzgänger" ist meiner Meinung nach das englische Prädikat "to be married with". Diese Beziehung zwischen zwei Leuten wird in den meisten Biographien veröffentlicht. Oft wird noch hinzugefügt, ob es eine glückliche Ehe ist oder war, und wie lange sie schon andauert oder gedauert hat. Die meisten Menschen würden sa-

gen, dass eine Eheschließung ein rechtlicher Akt ist und aufgrunddessen keine Einwände bestehen dürften, dies Außenstehenden mitzuteilen. Doch betrifft eine Ehe immer zwei Personen und ist somit etwas sehr persönliches. Damit soll nur klar gestellt sein, dass es auch Prädikatrelationen gibt, welche sehr wohl in Lebensläufen öffentlich bekannt gegeben werden dürfen, die dennoch eine starke Brücke zum Privatleben der jeweiligen Personen schlagen.

Im Zuge dieser Arbeit werden nur persönliche Relationen in ihrer zweiten Bedeutung betrachtet. Somit werden nur Prädikatrelationen untersucht, die jemanden persönlich betreffen, wie z.B. "He was born as son of a blacksmith in 1955.".

#### 1.2.2 Öffentliche Relationen

Des Weiteren gibt es eine große Anzahl an Prädikaten, welche in die Kategorie der öffentlichen Relationen fallen. Wie der Name "öffentliche Relation" schon verrät, handelt es sich hierbei um Verben, die hauptsächlich in offiziellen Lebensläufen genannt werden und sachliche Informationen aus dem Leben dieser Personen bekannt geben. In der Regel handelt es sich hierbei um Prädikatrelationen, die biographische Fakten beschreiben, welche beispielsweise für die Leser von Wirtschaftsnachrichten von Interesse sein dürften. Hierunter fallen Relationen, welche eventuell Aufschluss darüber geben, welchen Beruf die Person ausübt, oder bei welchem Unternehmen sie gerade beschäftigt ist. Zudem enthalten öffentliche Relationen Details aus dem Leben der jeweiligen Person, bei denen abgeklärt wurde, ob die betreffende Person mit der Veröffentlichung dieser Daten einverstanden war. Manchmal ist dies auch nicht der Fall, wie der Artikel "Google balances privacy, reach" von Elinor Mills [64] (siehe Vorwort) gezeigt hat. Doch die Missachtung der Privatsphäre bei der Informationssuche ist ein anderes Thema und macht die gefundenen Fakten nicht weniger offiziell.

Da die Aufgabe der Erkennung von Menschenbezeichnern innerhalb biographischer Relationen sich auf die Domäne der englischsprachigen Wirtschaftsnachrichten beschränken wird, werden öffentliche Relationen im Zentrum dieser Untersuchung stehen.

#### 1.2.3 Zufällige Relationen

Vollständigkeitshalber sollten auch die "zufälligen Relationen" angesprochen werden. Denn im Leben der Menschen gibt es enorm viele zufällige Begebenheiten. Auch über sie lassen sich zahlreiche Geschichten erzählen. Oft können zufällige Ereignisse zusammen mit persönlichen Gefühlen auftreten und dann vermischen sich wieder persönliche mit zufälligen Relationen. Leider sind zufällige Prädikatrelationen am uninteressantesten für das Auffinden von Personen in biographischen Kontexten, da sie kaum vorhersehbar sind und in einer solche Vielfalt vorkommen, dass sie schwer aufzuzählen sind. Außerdem wird einem Beinbruch, einer Verliebtheit oder einem Streit in der Familie meist wenig Beachtung von Außenstehenden geschenkt.

Biographien schreibt das Leben - welche Art von Information bzw. Relation sie enthalten, hängt allein vom Autor ab.

#### 1.3 Einschätzung der Thematik

#### 1.3.1 Probleme und Chancen

Die Aufgabe der automatischen Erkennung von Menschenbezeichnern innerhalb biographischer Relationen in englischsprachigen Wirtschaftsnachrichten wird sicher kein leichtes Unterfangen werden. Doch ist es eine Herausforderung, der man sich ohne Weiteres stellen kann, indem man sich zunächst ein Bild von der Ausgangssituation macht und sich danach die möglichen Schwierigkeiten vor Augen führt.

Einerseits ist es wichtig, vorab abzuklären, welche Entitäten, Bezeichnungen oder andere Angaben in biographischen Kontexten vorkommen.

So kann ein Personenname beispielsweise aus einem Titel oder einer Anrede gefolgt von einem Nachnamen bestehen. Des Weiteren werden in biographischen Texten häufig Beschäftigungsverhältnisse beschrieben und in diesem Zusammenhang werden sicherlich Organisationsnamen bzw. Firmennamen auftreten. Auch eine Liste an Branchen, Fachbereichen und Industriesektoren kann von Vorteil sein, wenn nur die Arbeitsdomäne einer Person genannt wird. Zudem kommen in diesen Kontexten häufig Ortsbestimmungen und Beschäftigungszeiträume vor.

So wäre es nur sinnvoll, Wörterbücher für Titel und Anredemöglichkeiten zu erstellen, sowie Vor- und Nachnamen aufzulisten, aber auch vollständige Personennamen zu archivieren. Außerdem lassen sich weitere benannte Entitäten wie Toponyme und Organisationen ebenfalls mit der Hilfe von Lexika in den Griff bekommen. Mit anderen Kategorien von Bezeichnern kann ähnlich verfahren werden, so dass Hyperonymierelationen in Form von Wörterbüchern kodiert werden und somit semantische Klassen gebildet werden.

Andere linguistische Phänomene lassen sich dagegen schlecht mittels Lexika beschreiben, dafür scheinen sie recht gut über lokale Grammatiken dargestellt zu werden.

• Darunter fällt z.B. die **syntaktische Variabilität**. Gerade wenn man an die Beschreibung von Datumsangaben oder Personennamen denkt, gibt es eine Reihe an syntaktischen Möglichkeiten, wie diese ausgedrückt werden können.

on February 20, 2004 on 3 June 1994 on Tuesday 6th April 2005 12-Feb-2006 in March 1960

Bill Gates Mr. Gates William Henry Gates III William Gates

• Ein weiteres klassisches Problem ist die Unterscheidung zwischen der Firma, ihrer Marke und ihrem Produkt. Dafür wäre "Apple" ein Paradebeispiel, denn allein der

Kontext, in dem dieser Begriff fällt, könnte für die Auflösung dieser Ambiguität, sorgen. Solche **Disambiguierungen** können mit der Hilfe von lokalen Grammatiken relativ leicht und zugleich recht anschaulich vorgenommen werden.

• Zudem tragen sie nicht nur zur Bedeutungsunterscheidung innerhalb von Named Entities bei, sondern auch zwischen Eigennamen und allgemeinen Bezeichnungen. Beispielsweise gibt es einige Nachnamen, welche gleichzeitig auch in ihrer Funktion als Nomen eine Pflanze wie "Bush", eine Berufsbezeichnung wie "Miller" oder eine Farbe wie "Blue" sein können. Nur eine detaillierte Beschreibung des Kontextes durch eine lokale Grammatik kann verhindern, dass z.B. "The Burning Bush" sicher keine Person bezeichnet, aber "Bush jr." auf jeden Fall einen Menschen referenziert.

#### 1.3.2 Bewältigung der Aufgabe

Wie aus dem letzten Abschnitt hervorgeht, gibt es einiges zu bedenken, wenn man sich an die Aufgabe heranwagt, eine Grammatik für Menschenbezeichner in biographischen Kontexten zu entwickeln.

Aufgrunddessen werden die folgenden Kapitel einen Einblick in die Herangehensweise an dieses Thema geben und dabei die entsprechenden Lösungsansätze präsentieren.

Zunächst werden Begrifflichkeiten, Bedeutung und Funktionalität von lokalen Grammatiken in Kapitel 2 erläutert. Im Anschluss daran wird noch im selben Kapitel auf die Arbeitsweise mit dem System Unitex eingegangen, um deutlich zu machen, wie mit lokalen Grammatiken gearbeitet werden kann.

Nachdem die Grundlagen zu Grammatiken gelegt wurden, können in Kapitel 3 interessante Ansätze weiterer Linguisten vorgestellt werden, die große Fortschritte auf dem Gebiet der automatischen Erkennung von Eigennamen mittels lokaler Grammatiken erzielt haben und deren Arbeiten meinen Ansatz zur Erkennung von Menschenbezeichnern in biographischen Kontexten geprägt haben.

In Kapitel 4 werden alle gewollten Beschränkungen für meinen Ansatz zur Erkennung von Menschenbezeichnern in biographischen Kontexten beschrieben. Es werden Erklärungen gegeben, warum man sich auf die Korpusdomäne der Wirtschaftsnachrichten festgelegt hat und wieso die Personen gegenüber anderer Entitäten im Vordergrund stehen.

Daraufhin werden Einzelheiten zu den im System verwendeten Ressourcen Preis gegeben. Dabei werden in Kapitel 5 die verschiedenen Korpora und alle selbst erstellten Lexika angesprochen, mit deren Hilfe dieser Ansatz zur Erkennung von Eigennamen umgesetzt wurde.

Die folgenden Kapitel stellen die entwickelten lokalen Grammatiken für Entitäten wie Personen, Organisationen, Toponyme und Datumsangaben vor und geben Einblick in die Grammatiken der persönlichen, sowie beruflichen Relationen.

Zuletzt wird die Qualität des Systems gemessen, indem die Ergebnisse der lokalen Grammatiken auf einem Testkorpus evaluiert werden. Außerdem wird aufgezeigt, wie persönliche und berufliche Relationen aus dem Text extrahiert werden können, um die syntaktische und semantische Vielfalt dieser Prädikatrelationen zu veranschaulichen.

#### 2 Lokale Grammatiken

#### 2.1 Was sind lokale Grammatiken?

Lokale Grammatiken kann man als "Landkarten der Sprache" bezeichnen (Mallchok, 2004 [55]), die einerseits Sequenzen von Wörtern, welche semantische Einheiten bilden, und andererseits syntaktische Strukturen beschreiben.

Überdies geben sie noch Aufschluss über die morphosyntaktischen Eigenschaften, der darin beschriebenen Elemente, welche syntaktisch (Fairon, 2000 [23]) oder semantisch (Constant, 2000 [10]) geprägt sein können.

Des Weiteren können sie in den verschiedensten Varianten für automatische Sprachverarbeitung auf Textkorpora nützlich sein. Besonders auf dem Gebiet der lexikalischen Disambiguierung werden lokale Grammatiken verstärkt eingesetzt (nach Blanc & Dister, 2004 [5]).

Da Wortformen isoliert gesehen oft ambig sind, kann ein Teil von ihnen aber durch die Analyse des Kontextes disambiguiert werden. Der für die Disambiguierung relevante Kontext wird durch eine lokale Grammatik (Gross, 1997 [43]) beschrieben, die durch einen endlichen Automaten bzw. einen Transduktor repräsentiert wird. Lokale Grammatiken werden nicht nur für die Disambiguierung, sondern auch für andere Aufgaben genutzt, wie die Erkennung von Mehrwortlexemen und Komposita, die Repräsentation orthographischer Varianten im Lexikon, sowie die Überprüfung der Kongruenz oder Identifikation von Zeitangaben und anderen Entitäten (vgl. Blank, 1997 [6]).

Endliche Automaten bzw. Transduktoren beschreiben komplexe linguistische Strukturen, die so nicht in einer Lexikongrammatik oder in elektronischen Wörterbüchern formalisiert werden könnten. Eigentlich sind Transduktoren endliche Automaten, die zusätzlich eine Ausgabe erzeugen, wenn die in der Definition des Automaten spezifizierte(n) Sequenz(en) erkannt wurde(n). Der "Eingabeteil" des Transduktors dient dazu, spezifische Sequenzen im Text zu erkennen. Der "Ausgabeteil" führt einerseits Substitutionen im Text aus, versieht andererseits identifizierte Sequenzen mit zusätzlichen Informationen (z.B. einer Wortklasse) oder fügt linguistische Markierungen (z.B. die Annotation von Phrasen) in den Text ein (nach Blank, 1997 [6]).

In der Regel werden lokale Grammatiken in Form von Graphen (Silberztein, 1993 [84]) visualisiert. Die Kombination von parametrisierten Graphen mit einer Lexikongrammatik kann beispielsweise äußerst effektiv bei der syntaktischen Analyse einfacher Sätze sein (Paumier, 2001 [73] und Laporte, 2005 [53]).

Graphen sind sehr geeignete Repräsentationen für lokale Grammatiken, denn es gibt diverse Graphikprogramme, mit denen sich diese Graphen leicht erstellen, erweitern oder abändern lassen. Die beiden Systeme INTEX und Unitex bieten u.a. solche Zeichenprogramme für Automaten an.

Jeder Graph besteht aus einem Anfangszustand, der durch einen Rechtspfeil symbolisiert wird. Dieser Rechtspfeil geht von keinem Zustand aus, sondern führt lediglich zu einem der nächsten Zustände im Graphen. Außerdem enthält jeder endliche Graph einen Endzustand, welcher meist durch einen doppelt umrandeten Kreis dargestellt wird. Die Graphen werden von links nach rechts interpretiert und so werden die möglichen Pfade "abgelaufen" und ihre Muster im Text gesucht. Bei den Systemen INTEX und Unitex steht jeder Zustand bzw. jeder Knoten für Wörter (mit oder ohne ihrer morphologischen Informationen) oder für Klassen aller Flexionsformen von Wörtern, wenn diese in spitzen Klammern notiert wurden. Somit werden die Eingabesequenzen des Transduktors nicht an den Ubergängen zu den Zuständen genannt, sondern in den Zuständen selbst. Natürlich sind auch wie bei endlichen Automaten  $\epsilon$ -Transitionen erlaubt. Alle Transitionen werden durch Verbindungslinien zwischen den einzelnen Zuständen dargestellt. Das leere Wort wird als <E> in den Knoten angegeben. Es wird sogar gestattet Subgraphen innerhalb eines Automaten aufzurufen, was die Übersichtlichkeit der Graphen erhöht. Diese Subgraphen werden grau unterlegt, so dass eine Unterscheidung zwischen einem einfachen Zustand und einem Zustand, der einen weiteren Graphen aufruft, möglich

Die eben beschriebenen Graphen sind auch als gerichtete azyklische Graphen bekannt, da sie keinerlei Zyklen enthalten. Im englischen Sprachraum werden sie als "Directed Acyclic Graphs" bezeichnet und werden deshalb im deutschen Sprachraum häufig nur DAGs genannt. Mathematisch gesehen repräsentiert ein DAG eine Halbordnung.

#### 2.2 Warum werden lokale Grammatiken verwendet?

Die meisten Versuche linguistische Theorien oder Grammatiken zu entwickeln, welche umfassend und stark verallgemeinert beschreiben wollen, wie eine Sprache aufgebaut ist und wie Syntax, Morphologie und Semantik zusammenwirken, waren wenig befriedigend. Denn Ziel solch einer Grammatik sollte es immer sein, alle Sätze, die in einer Sprache möglich sind, abzudecken, und kein Satz, der mit dieser Grammatik gebildet werden konnte, durfte grammatikalisch oder semantisch unstimmig sein.

Anfangs ging man an dieses Problem so heran, dass jede explizite Komponente im Satz durch ihre jeweilige grammatikalische Kategorie ersetzt wurde. Noam Chomsky fasste 1957 diese Grammatiken unter dem Begriff "Kontextfreie Grammatik" zusammen, musste aber einräumen, dass es immer noch einige Unzulänglichkeiten in Bezug auf die formale Repräsentation natürlicher Sprache gab. Diese Grammatiken beschrieben in der Regel nur einfache Sätze und gingen kaum die Abhängigkeiten der einzelnen Satzteile untereinander ein (Gross, 1997 [43]).

Dagegen waren die späteren Ansätze von Zellig Sabbetai Harris und Noam Chomsky

schon spezieller, da sie inzwischen Bildungsregeln für die einfachen Sätze definierten und diese dann untereinander kombiniert wurden, so dass komplexe Sätze geformt wurden. Im Grunde war es damals schon ein kleiner Schritt in Richtung Diskursanalyse, den die beiden vollzogen. Denn sie legten Regeln fest, welche die Satzstellung innerhalb der einfachen Sätze variierten und einfache Sätze zu komplexen Satzgefügen verbanden.

Irgendwann stellte sich dann heraus, dass diese theoretische Sichtweise der natürlichen Sprache, die immer komplexer werdenden Beschreibungsformalismen und die vielen Ausnahmen, welche sich in die Bildungsregeln eingeschlichen hatten, nicht mehr zu handhaben waren. Daraufhin besonnen sich viele Linguisten darauf das Phänomen "Sprache" anders zu erforschen. In ihrer Herangehensweise verhielten sie sich ähnlich wie Naturwissenschaftler, denn man muss keine Sätze erfinden - es gibt sie schon und man muss das Vorhandene erst einmal untersuchen, bevor neues automatisch generiert werden kann. Laut Gross findet man eine Grammatik im Text und muss sich nicht erst eine ausdenken.

Deshalb sollte man als Linguist keine Theorie in die Welt setzen, bevor man nicht Korpusmaterial gesammelt hat und seinen Ansatz auf realem Text verifiziert hat. Denn indem Satzkorpora gebildet werden, deren syntaktische und semantische Struktur analysiert wird, entstehen indirekt schon Regeln zur Beschreibung der Sprache.

Des Weiteren war Zellig S. Harris davon überzeugt, dass die Untersuchung von Subsprachen in Verbindung mit lokalen Grammatiken besonders vielversprechend sein dürfte, weil Subsprachen

- thematisch begrenzt sind,
- lexikalischen, syntaktischen und semantischen Restriktionen unterliegen,
- in ihren grammatikalischen Eigenschaften nicht der Allgemeinsprache gleichen,
- gewisse lexikalische Strukturen relativ häufig wiederholen
- in sich strukturiert sind und
- eine gewisse Symbolik verwenden.

So können Elemente der Sprache, die in lokalen Grammatiken erfasst werden, als kleine, aber aussagekräftige Subsprachen gesehen werden und Beschreibungsversuche von Subsprachen würden in ihrer Repräsentation erweiterten lokalen Grammatiken entsprechen.

Die Einschränkung der Sprache auf eine bestimmte Bezugsdomäne wie z.B. auf Wirtschaftsnachrichten und die damit verbundene Verwendung von themenspezifischen Fachvokabular rechtfertigen gewiss den Einsatz von lokalen Grammatiken. Aufgrunddessen sind lokale Grammatiken zur syntaktischen und semantischen Analyse von Menschenbezeichnern innerhalb biographischer Relationen sicherlich die richtige Entscheidung.

## 2.3 Unitex - Ein System zur Anwendung lokaler Grammatiken

#### 2.3.1 Allgemeines

Unitex ist ein Korpusverarbeitungssystem, welches es ermöglicht, mit elektronischen Ressourcen wie z.B. elektronischen Lexika umzugehen und lokale Grammatiken zu entwickeln und anzuwenden. Dabei wird auf drei Ebenen der Sprache - der Morphologie, dem Lexikon und der Syntax - gearbeitet.

Die Hauptfunktionen von Unitex sind u.a

- das Erzeugen, sowie die Anwendung und Verarbeitung elektronischer Wörterbücher,
- die Benutzung von regulären Ausdrücken zum Pattern Matching,
- die Interpretation rekursiver Transitionsnetze zum Pattern Matching,
- die Anwendung von lokalen Grammatiken und Lexikongrammatiken und
- die Auflösung von Ambiguitäten über den Text-Automaten.

Das Konzept für das System Unitex wurde am LADL (*Laboratoire d'Automatique Documentaire und Linguistique*) unter der Leitung von Prof. Maurice Gross entwickelt, und das dazugehörige Programm wurde am Institut Gaspard-Monge (IGM) der Université de Marne la Vallée von Sébastien Paumier implementiert.

Derzeit werden für Unitex Lexika in 14 verschiedene Sprachen (Deutsch, Englisch, Finnisch, Französisch, Griechisch, Italienisch, Koreanisch, Norwegisch, Polnisch, Portugiesisch, Brasilianisches Portugiesisch, Russisch, Spanisch und Thai) angeboten.

Da Unitex im Gegensatz zu INTEX frei verfügbar ist und unter der GNU GPL (GNU General Public License) steht, kann es im Grunde jeder benutzen. Außerdem stellt es ganz ähnliche Funktionen wie INTEX zur Verfügung und ist auf allen gängigen Betriebssystemen (Windows, Linux, MacOS) lauffähig.

Vorallem bietet Unitex eine komfortable und intuitiv bedienbare Oberfläche zur Entwicklung von Grammatiken. Dabei handelt es sich um eine Java-Oberfläche, von der aus diverse C++-Programme gesteuert werden.

#### 2.3.2 Textvorverarbeitung

Unitex arbeitet mit der Kodierung "UTF-16 Little Endian" und unterstützt somit den Unicode 3.0 Standard. Dadurch wird selbst die Verarbeitung asiatischer Sprachen ermöglicht. Zur Konvertierung der Texte empfiehlt sich das Programm Convert von Unitex. Nachdem Unitex mit der gewählten Sprache gestartet worden ist, kann man einen Text mit der Kodierung UTF-16 LE öffnen. Dabei wird gefragt, wie der Text vorverarbeitet werden soll. Die Textvorverarbeitung von Unitex setzt sich aus den Schritten

 $<sup>^6 {\</sup>tt http://www-igm.univ-mlv.fr/\~unitex/download.html}$ 

Normalisierung, Satzenderkennung, Auflösung von Kontraktionen, Tokenisierung und lexikalische Analyse des Korpuses zusammen.

#### 2.3.2.1 Normalisierung

Es ist Aufgabe des Programms *Normalize* die Normalisierung des Textes vorzunehmen, indem Folgen von Leerzeichen bzw. Zeilenumbrüchen durch ein Zeichen ersetzt werden. Gleichzeitig wird die interne Syntax von eventuell lexikalisch annotierten Token überprüft.

#### 2.3.2.2 Satzenderkennung und Auflösung von Kontraktionen

Unitex bietet eine sprachspezifische Satzenderkennung mittels lokaler Grammatiken in Form von Graphen an. Des Weiteren werden Kontraktionen wie z.B. "I'm" zu "I am" oder "you're" zu "you are" aufgelöst und verschiedene Arten von Anführungszeichen vereinheitlicht.

#### 2.3.2.3 Tokenisierung

Hierfür ist das Programm *Tokenize* von Unitex zuständig. Die Tokenisierung wird aufgrund des Alphabets der jeweiligen Sprache vorgenommen. Die daraus resultierende Tokenliste wird für spätere Zwecke im Arbeitsverzeichnis des aktuellen Textes gespeichert.

#### 2.3.2.4 Lexikalische Analyse

Bei der lexikalischen Analyse werden alle Standardwörterbücher der jeweiligen Sprache und eventuell noch eigene Lexika auf die Tokenliste angewandt. Dabei kommt das Programm *Dico* zum Einsatz, welches alle Token mit der entsprechenden grammatikalischen oder semantischen Information aus den Lexika versieht. Alle Lexika, welche vom System Unitex verwendet werden sollen, müssen formal dem Standard der DELA Wörterbücher entsprechen.

#### 2.3.3 DELA Wörterbücher [29]

Das klassische Wörterbuch ist eine Sammlung von Wörtern oder einer Kategorie von Wörtern einer Sprache, die in der Regel in alphabetischer Ordnung mit Erläuterungen in derselben Sprache oder einer Übersetzung derer in eine andere Sprache aufgelistet sind (Lexis, 1975). Dagegen ist das elektronische Wörterbuch eine formale Repräsentation eines Lexikons, welche jeder Flexionsform ihr Lemma, genauso wie die entsprechende grammatikalische, Flexions- und eventuelle semantische Information zuweist (nach Sébastien Paumier)<sup>7</sup>.

Überdies hinaus wird von einem elektronischen Wörterbuch gefordert, dass es formal

<sup>&</sup>lt;sup>7</sup>Übersetzung aus dem Französischen

http://wwwigm.univmlv.fr/~paumier/DEA/Cours%206%20%20Dictionnaires%20electroniques.pdf

und vollständig ist, so dass es sich maschinell verarbeiten lässt und es von Programmen automatisch verändert werden kann. Theoretisch müsste es 100% des Lexikons abdecken, was allerdings kaum realisierbar ist.

DELA ist ein elektronisches Wörterbuchsystem und steht für "Dictionnaires électroniques du LADL" (Laboratoire d'Automatique Documentaire et Linguistique). In den 60er Jahren wurde es von Prof. Maurice Gross ins Leben gerufen, und war zunächst unter dem Namen "Lexikon Grammatik" bekannt. Das DELA ist eine formale Repräsentation der jeweiligen Sprache; das heißt, Spracheigenschaften werden strukturiert abgespeichert, wobei sowohl Vokabular als auch Morphologie berücksichtigt werden.

Die DELA-Wörterbuchfamilie gliedert sich in folgende Teillexika:

- DELAS "mots simples": Wörterbuch für die einfachen Wörter
- DELAC "mots composés": Wörterbuch für die komplexen Wörter
- DELAF "formes fléchies": Wörterbuch der einfachen Wörter, deren Flexionsmerkmale kodiert sind.
- DELACF "mots composés avec les formes fléchies": Wörterbuch der komplexen Wörter, deren Flexionsmerkmale kodiert sind.

Dabei werden als einfache Wörter ("mots simples") Sequenzen zusammenhängender Buchstaben eines Alphabets einer bestimmten Sprache verstanden, wie z.B. angry,.A oder acually,.ADV oder bodies,body.N:p.

Dagegen sind komplexe Wörter ("mots composés") Sequenzen zusammengesetzter lexikalischer Einheiten wie einfache Wörter, Trennzeichen oder Ziffern.

Beispiele aus dem Französischen wären hierfür coup de chance, .N+NDN:ms (Glückstreffer) oder coup de pied, .N+NDN:ms (Fußtritt) oder das ambige coup de foudre, .N+NDN:ms (Liebe auf den ersten Blick /Blitzschlag).

Die eben genannten Beispiele deuteten bereits an, dass hinter einem Eintrag im DELAF eine gewisse Symbolik steht.

So besteht ein Lexikoneintrag im DELAF aus 5 verschiedenen Feldern:

- 1. Flektierte Form des Wortes
- 2. Lemma des Wortes (Kanonische Form)
- 3. Charakteristische Informationen zur Lemmaform
- 4. Grammatikalische Eigenschaften der flektierten Form
- 5. Optionale Ergänzungen für den menschlichen Betrachter

Analog dazu wird ein Eintrag im DELACF gebildet. Dabei sollte man noch anmerken, dass das zweite Feld (die Lemmaform) immer dann leer ist, wenn sie mit der flektierten Form identisch ist. Dafür wird das vierte Feld (die grammatikalische Information für die flektierte Form) nicht belegt, wenn das Wort eindeutig ist, und es nicht variiert

werden kann. Außerdem wird das fünfte und letzte Feld (die Zusatzinformation) nur besetzt, wenn die flektierte Form - das Ausgangswort - ein Kompositum ist. Genau die gleichen Regeln gelten für Lexikoneinträge im DELAS und DELAC, nur dass hier die Flexionsinformation entfällt.

An einem konkreten Beispiel würde dies nun folgendes bedeuten:

bodies, body.N:p

₩

bodies : flektierte Form

body : Lemmaform

N : grammatikalische Information (Nomen)

p : grammatikalische Eigenschaft der flektierten Form (Plural)

Bei der Erstellung eigener Lexika sollte darauf geachtet werden, dass Mehrwortlexeme direkt im Lexikon kodiert werden, weil sonst Fehler bei der Tokenisierung gemacht werden. Wenn man nur ein Teilformenlexikon verwenden würde, könnte beispielsweise "grand-mère" nicht als ein Wort erkannt werden. Oft besteht auch die Möglichkeit Mehrwortlexeme wie "grand-mère" (Großmutter) anstatt des Bindestrichs mit einem Leerzeichen dazwischen zu schreiben. Dafür wäre dann grand=mères,grand=mère.N:fp der entsprechende Lexikoneintrag, denn das '=' ist ein Metazeichen, was für einen Bindestrich '-' und für ein Leerzeichen ' ' steht.

Je nachdem wie ausführlich die Kodierung eines Lexikons mit diversen grammatikalischen oder semantischen Angaben vorgenommen wurde, spricht man von 3 Stufen der Lexikonkodierung:

- **DELAF-S** (*"short"*): Es werden minimale Angaben zur grammatikalischen Analyse der einzelnen Formen gemacht. Das heißt, dass lediglich Informationen zur jeweiligen Wortart und zur Flexion kodiert werden. Hier wird ausschließlich auf die *Grammatik* Bezug genommen.
- **DELAF-M** ("medium"): Die Lexikoneinträge werden um semantische Informationen zu den Nomina erweitert. Dabei wird spezifiziert, welche Eigenschaften das Nomen hat, z.B. ob es ein Menschenbezeichner Hum, ein Konkreta Conc oder ein Tier Anl etc. ist. Außerdem werden Determinativa DET und Pronomina PRO durch weitere Unterkategorien versehen. Auf diese Weise wird die Semantik miteinbezogen.
- **DELAF-L** (*"large"*): Hierbei werden die Wörterbucheinträge um die Lexikon-Grammatik der LADL ergänzt, so dass die syntaktischen Eigenschaften der Verben im Französischen markiert werden (Berücksichtigung der *Syntax*).

Wie ausführlich nun ein Lexikoneintrag erstellt wird, hängt ganz von seiner späteren Funktion ab und über welche Art von Informationen er später angesprochen werden soll. Das heißt nichts anderes, als dass beispielsweise Nomina, welche die semantische Funktion eines Menschenbezeichners haben, auch als solche markiert werden sollten.

Legt man allerdings nachher Wert auf Kongruenzeigenschaften, so sollte man auf keinen Fall die grammatikalische Information außer Acht lassen.

#### 2.3.4 Prioritäten bei der Anwendung der Lexika

Unitex unterscheidet drei Prioritäten bei der Anwendung der Lexika, falls der Dateiname eines Lexikons (ohne die Endung .bin) auf '-' bzw. '+' endet:

- 1. \*-.bin (höchste Priorität diese Lexika werden vorrangig behandelt)
- 2. \*.bin (durchschnittliche Priorität diese Lexika werden zweitrangig behandelt)
- 3. \*+.bin (niedrigste Priorität diese Lexika werden zuletzt auf den Text angewandt)

Token, die einem der Lexika einer Prioritätsebene gefunden wurden, werden in keinem Lexikon mit untergeordneter Priorität mehr nachgeschlagen. So lassen sich z.B. bestimmte Lesarten für ein Token erzwingen, da das höher priorisierte Lexikon wie ein Filter andere Bedeutungen aussiebt. Innerhalb einer Prioritätsebene werden alle Lexika gleichrangig behandelt, d.h. verschiedene Lesarten eines Tokens aus unterschiedlichen Lexika werden ins Textlexikon geschrieben.<sup>8</sup>

#### 2.3.5 Mustererkennung und Konkordanzen

Wie bereits erwähnt, werden lokale Grammatiken im System Unitex als Graphen (DAGs) repräsentiert. Möchte man nun eine lokale Grammatik auf einem Korpus testen, so wählt man den entsprechenden Graphen aus, und das Programm *Locate* wendet diesen Graphen auf den Text an und erstellt den Index für eine Konkordanz. Dabei bietet *Locate* dem Benutzer verschiedene Arten der Textsuche an, bei der

- die kürzesten Treffer,
- die längsten Treffer oder
- alle Treffer

ausgegeben werden.

Außerdem lässt sich das Verhalten des Graphen steuern, falls es sich um einen Transduktor handelt. Es gibt folgende Möglichkeiten:

- Die Ausgabe des Transduktors bleibt unberücksichtigt.
- Die Ausgabe des Transduktors wird links vom Treffer eingefügt.
- Die gefundene Sequenz wird durch die Ausgabe des Transduktors ersetzt.

Für das Anfertigen einer Konkordanz ist das Programm Concord zuständig. Es gibt einerseits die Konkordanz in verschiedenen Formaten aus (HTML, Text) und andererseits lässt sich die Länge des Kontextes und die Sortierweise der Treffer spezifizieren.

<sup>&</sup>lt;sup>8</sup>vgl. http://www.cis.uni-muenchen.de/~wastl/lg/introUnitex.pdf

## 3 Zusammenfassung früherer Arbeiten

In diesem Kapitel möchte ich auf Personen eingehen, die große Fortschritte auf dem Gebiet der automatischen Erkennung von Eigennamen mittels lokaler Grammatiken erzielt haben und deren Arbeiten meinen Ansatz zur Erkennung von Menschenbezeichnern in biographischen Kontexten geprägt haben.

#### 3.1 Maurice Gross

Maurice Gross est décédé à Paris le 8 décembre 2001 des suites d'un cancer, insoupçonné jusqu'à peu avant sa mort.<sup>9</sup>

Né le 21 juillet 1934 à Sedan, Maurice Gross (...) dépouille d'innombrables grammaires, pour le français et l'anglais - et l'allemand, sa spécialité; pour le russe, en traduction anglaise. (...) Il rêve d'inventorier la langue et de mettre de l'ordre dans les grammaires en trouvant des règles. (...) Il prend la tête du Laboratoire d'automatique documentaire et linguistique (LADL-CNRS), à Paris-VII, qui multiplie dépouillements informatisés, dictionnaires automatiques, correcteurs de fautes, etc. (...). Son ouvre était encore en développement. Elle s'est brusquement arrêtée avec la mort. 10

#### 3.1.1 Zur Person

Maurice Gross war wohl einer der größten Linguisten unserer Zeit. Er war ein Mann voller Tatendrang und Ideenreichtum, den letztendlich sein Tod am 8. Dezember 2001 auf tragische Weise daran hinderte, sein Werk fortzusetzen.

Sein Traum war es die Sprache zu katalogisieren, er wollte ihr "Sprachinventar" wiedergeben und eine Ordnung in sie hineinbringen, indem er die Sprache anhand von Texten studierte, ihre syntaktischen Muster analysierte und daraus Regeln zur Bildung von Phrasen ableitete - im Grunde Grammatiken entwickelte.

Er war ein recht vielseitiger Mann, den viele Gebiete der Linguistik begeisterten. Doch man könnte sicher auch behaupten, dass die lokalen Grammatiken und die damit verbundene elektronische Lexikographie sein "Steckenpferd" waren. Denn in seinen zahlreichen Publikationen, nahmen sie wohl den meisten Platz und größten Stellenwert ein, da er damit jedes Detail der französischen oder englischen Grammatik untersuchen und seine Ergebnisse anschaulich präsentieren konnte.

<sup>9</sup>http://www-igm.univ-mlv.fr/~laporte/artJCChev.htm

<sup>10</sup>http://www-igm.univ-mlv.fr/~laporte/inmemoriam.doc

#### 3.1.2 Bootstrapping bei der Entwicklung lokaler Grammatiken

Bereits in Kapitel 2 wurden ein paar Arbeiten von Maurice Gross im Zusammenhang mit den lokalen Grammatiken genannt. Ein wichtiger Artikel für die praktische Arbeit mit lokalen Grammatiken ist noch "A Bootstrap Method for Constructing Local Grammars" [45], welcher auf keinen Fall in diesem Zusammenhang ungenannt bleiben sollte.

Maurice Gross stellt hierbei einen Ansatz vor, lokale Grammatiken oder elektronische Wörterbücher um einen Schlüsselbegriff oder um zusammenhängende Begriffe, die eine semantische Einheit bilden, herum zu entwickeln. Dabei wird zunächst vom vorhandenen Lexikoninventar ausgegangen und jeder Eintrag als solch ein Schlüsselbegriff gesehen. Mit Hilfe einer Suchfunktion auf dem Text kann der jeweilige Kontext zu den Ausgangsbegriffen ermittelt werden.

Zu jedem neuen Vorkommen im Korpus wird der Kontext entsprechend seiner lexikalischen Funktion manuell ausgewertet. So können relativ schnell und auf einfache Art und Weise Mehrwortlexeme gefunden werden, die das Schlüsselwort in irgendeiner Form enthalten. Als nächstes empfiehlt es sich die unmittelbaren Kontexte des Ausgangswortes zu schematisieren, um später gezielt alle möglichen Äußerungen, welche diesen Begriff enthalten, abdecken zu können.

Im konkreten Fall heißt das, dass beispielsweise das Wort "health" der Schlüsselbegriff war und in der Konkordanz lässt sich das Muster health and <N> identifizieren. Natürlich werden noch viele andere linguistische Schemata erkennbar sein, doch geht man jedes einzeln durch. Ausgehend von diesem regulären Ausdruck kann man eine Grammatik in Form eines Finite State Graphen (DAG) entwerfen, mit der alle Nomina gefunden werden, welche in diesem speziellen Kontext von "health" auftreten. Daraus ergibt sich dann wieder folgender Graph:

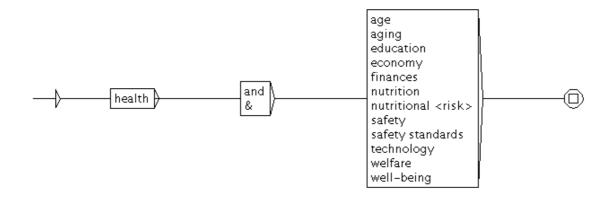


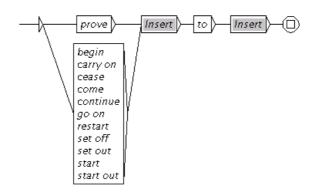
Abbildung 3.1: HealthAndN.grf aus Gross, 1999 [45]

Dieser Graph lässt sich erneut in einen anderen Graphen einbinden, welcher den linken Kontext von "health" spezifiziert. Diese Methode lässt sich nun beliebig oft wiederholen, bis alle möglichen Kontexte des Ausgangswortes ermittelt und beschrieben wurden. Diese Methode zur Entwicklung lokaler Grammatiken wird als **Bootstrapping** bezeichnet und liefert systematisch und schnell gute Ergebnisse bei der Grammatikerstellung.

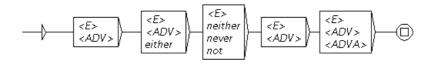
#### 3.1.3 Lemmatisierung zusammengesetzter Zeiten im Englischen

Des Weiteren hat Maurice Gross in den Jahren 1998/1999 ein sehr umfangreiches Graphenpaket zur Lemmatisierung zusammengesetzter Zeiten im Englischen entwickelt, welches er in seinem Aufsatz "Lemmatization of compound tenses in English" [44] ausführlich beschreibt.

Seine Graphen sollen später zur Erkennung personenbezogener Prädikate in dem hier vorgestellten Ansatz eingesetzt werden. Da sie mit wenigen Ausnahmen alle Verben des Englischen in verschiedenen Zeitformen finden können, sind sie eine Bereicherung für jede Arbeit. Stellvertretend für alle Graphen veranschaulicht der Graph aus Abbildung 3.2 die Struktur einer möglichen Verbalphrase des Englischen und berücksichtigt dabei auch mögliche Satzeinschübe (siehe Abbildung 3.3).



**Abbildung 3.2:** VModToV.grf aus Gross, 1998-1999 [44]



**Abbildung 3.3:** *Insert.grf* aus Gross, 1998-1999 [44]

Die Abbildungen 3.4 und 3.5 visualisieren das Zusammenspiel der einzelnen Graphen. Dabei wird deutlich, welcher Graph welchen Graphen aufruft, und es wird somit gezeigt, wie die Graphen voneinander abhängen. Der Ausgangspunkt ist der Graph VAUX, der so zu sagen alle Fäden bei der Erkennung der Verben "in der Hand hält".

Wie jeder dieser 80 Graphen genau aufgebaut ist, soll hierbei nicht von Interesse sein, weil nachher (in Abschnitt 5.2.6) nur mit dem Ergebnis, das sie erzielen, weitergearbeitet wird.

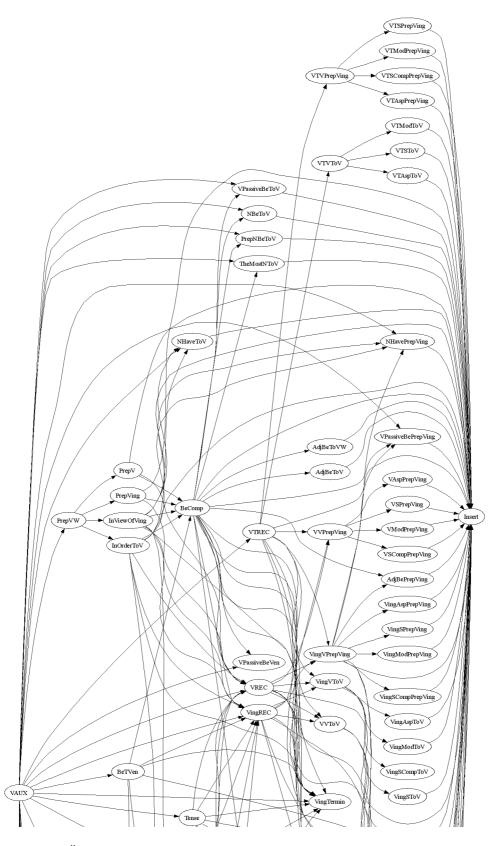


Abbildung 3.4: Übersicht der Lemmatisierungsgraphen aus Gross, 1998-1999 [44] (Teil 1)

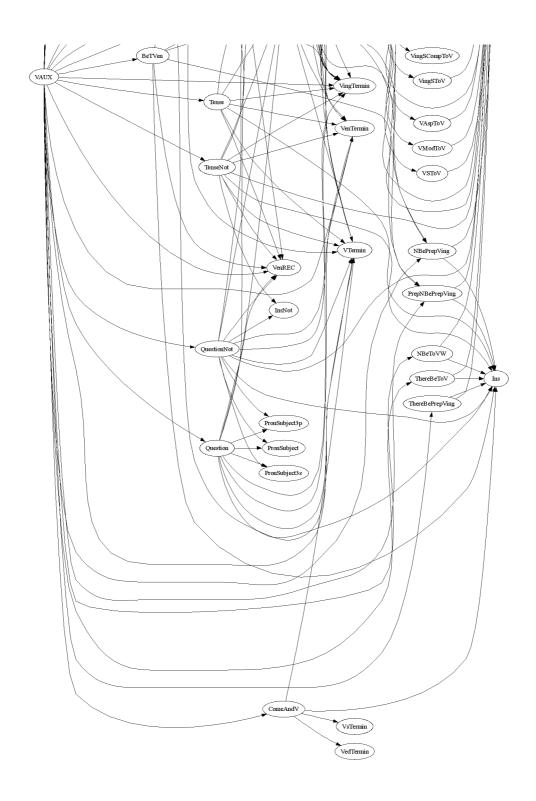


Abbildung 3.5: Übersicht der Lemmatisierungsgraphen aus Gross, 1998-1999 [44] (Teil 2)

#### 3.2 Jean Senellart

Jean Senellart, Directeur R&D<sup>11</sup>. Né en 1972. Diplômé de l'Ecole Polytechnique et titulaire d'un doctorat en Informatique linguistique de l'Université de Paris VII - LADL, Monsieur Senellart a commencé sa carrière comme chercheur et a enseigné à l'Ecole Polytechnique et à l'Université de Marne la Vallée.[88]

#### 3.2.1 Zur Person

Jean Senellart ist zur Zeit als Direktor der Forschungs- und Entwicklungsabteilung bei Systran S.A. in Paris beschäftigt. Systran ist eine der führenden Firmen im Bereich der maschinellen Übersetzung. Dem heute 36-Jährigen (Geburtsjahr 1972) wurde nach seinem Hochschuldiplom 1999 der Doktorgrad der Computerlinguistik (Informatique linguistique) am LADL (Laboratoire d'Automatique Documentaire et Linguistique) der Universität von Paris VII verliehen. Bevor er in die Wirtschaft ging, lehrte er u.a. an der der Universität von Marne la Vallée bei Paris. Des Weiteren wirkte er bei der Entstehung von GlossaNet<sup>12</sup> mit, einem online Konkordanzprogramm, dass wie eine Suchmaschine auf mehr als 100 elektronischen Zeitungsausgaben in 12 Sprachen arbeitet und auf dem System Unitex basiert, das an der Universität von Marne la Vallée entwickelt wurde.

#### 3.2.2 Erkennung von Eigennamen und Berufsbezeichnungen

#### 3.2.2.1 Motivation des Ansatzes

In seiner Arbeit "Locating noun phrases with finite state transducers" [81] beschreibt Jean Senellart einen Wörterbuch gestützten Ansatz zur Erkennung von Eigennamen mittels endlichen Transduktoren. Dafür hatte er sich zum Ziel gesetzt eine Lokale Grammatik zu entwickeln, die Nominalphrasen bestehend aus Eigennamen und/oder Berufsbezeichnungen beschreibt. Jedoch sollte sich die Erkennung von Eigennamen besonders von Personennamen - bzw. Berufsbezeichnern auf die Domäne der Zeitungsnachrichten beschränken. Dabei müssen aber auch semantische Relationen, wie Synonymie und Hyperonymie berücksichtigt werden, so dass Anfragen vom Typ "Find all newspaper articles in a general corpus mentioning the French prime minister." [81] oder "How is Mr. X referred to in the corpus; what have been his different occupations through out the period over which our corpus extends?" [81] verarbeitet werden konnten. Denn Antworten auf die erste Frage, werden wohl kaum Schlüsselworte aus der Query enthalten, sondern eher dazu passende Synonyme oder Eigennamen, die auf die Umschreibung "französischer Premierminister" zutreffen.

<sup>&</sup>lt;sup>11</sup>Recherche et Développement (R&D).

<sup>12</sup>http://glossa.fltr.ucl.ac.be/

<sup>&</sup>lt;sup>13</sup>Endliche Transduktoren werden im Englischen als Finite State Transducers (FST) bezeichnet.

#### 3.2.2.2 Im Vergleich mit anderen Information Retrieval Methoden

Der eben beschriebene Ansatz weicht stark von anderen gängigen Ansätzen der Informationsbeschaffung aus unstrukturierten Texten ab.

Weitere Konzepte zur automatischen Informationsgewinnung sind u.a.

- Algorithmen, die mit Schlüsselbegriffen arbeiten (Key-Word-Algorithms).
- Algorithmen, die nach Mustern exakt suchen (Exact-Pattern-Algorithms).
- Algorithmen, die die Statistik zu Hilfe nehmen (Statistical Algorithms).

Key-Word-Algorithms werden gerne von Suchmaschinen, wie z.B. Yahoo!, verwendet. Sie suchen nach Schlüsselbegriffen aus der Anfrage, die zusammen in einem Text auftreten. In der Regel werden noch leichte Abwandlungen in der Rechtschreibung, sowie verschiedene grammatikalische Endungen und Rechtschreibfehler akzeptiert und bei der Suche miteinbezogen.

Exact-Pattern-Algorithms bzw. Exact-String-Matching-Algorithms verwenden reguläre Ausdrücke aus Buchstaben, welche exakt auf dem Dokument suchen. Mit dieser Methode arbeitet u.a. das Oxford English Dictionary (OED). Bei der Eingabe des Suchstrings sind jedoch auch Wildcards wie das Fragezeichen? und der Asterisk \* erlaubt, wobei das Fragezeichen für einen beliebigen Buchstaben steht und der Asterisk eine beliebige Sequenz von Buchstaben repräsentiert. Des Weiteren beeinflusst die Groß- oder Kleinschreibung das Auffinden von Einträgen nicht, da case-insentive gesucht wird. Im Gegensatz zu den Key-Word-Algorithms muss jeder Term aus der Anfrage in der gegebenen Reihenfolge berücksichtigt werden.

Statistical Algorithms bieten dem Benutzer nur solche Dokumente als Ergebnis an, die sowohl Schlüsselwörter aus der Anfrage enthalten, aber auch statistisch gesehen semantisch nahe an den Anfragetermen liegen.

Am einfachsten zu implementieren sind wohl Algorithmen, die mit Schlüsselbegriffen aus der Anfrage arbeiten. Der Nachteil daran ist leider nur, dass die Ergebnisse sehr störanfällig sind, was nichts anderes heißt, als dass Homographen<sup>14</sup> der Anfrageterme im Text auftauchen können, oder dass Begriffe im Text gefunden werden, die sehr ähnlich zu den Anfragetermen sind.

Dagegen liefern Algorithmen, die mit Mustern bzw. regulären Ausdrücken arbeiten, ausgezeichnete Ergebnisse zurück. Jedoch sind die Muster hierbei so komplex, dass sich sogar Pattern spezifizieren lassen, mit denen man Synonyme der Anfrageterme finden kann. Des Weiteren lassen sich die verschiedenen grammatikalischen Endungen sehr präzise beschreiben. Nur wird es immer schwieriger die Muster zu konstruieren und zu verarbeiten, je komplexer die morphologischen Phänomene werden, welche es gilt zu beschreiben.

Ein Algorithmus, der auf statistischen Methoden basiert, kann lediglich für einfache Anfragen gute Resultate liefern, und man braucht große Dokumentenmengen, um statistisch repräsentative Ergebnisse zu bekommen. Doch dabei werden Terme mit niedriger Frequenz im Text meist ignoriert.

<sup>&</sup>lt;sup>14</sup>Ein **Homograph** ist ein Wort, das die gleiche Schreibweise wie ein oder mehrere andere Wörter hat, aber von unterschiedlicher Bedeutung ist und meist auch unterschiedlich ausgesprochen wird.

#### Welcher Ansatz wäre nun zur Erkennung von Eigennamen am idealsten?<sup>15</sup>

Der erste Ansatz würde funktionieren, wenn man entweder nur einen Vornamen oder einen Nachnamen in der Anfrage angeben würde, welcher auf keinen Fall ambig sein darf. D.h. dass beispielsweise Nachnamen wie "Major" nicht in der Query vorkommen dürfen, da sonst nicht nur ein Teil des Namens wie z.B. "John <u>Major</u>", sondern auch "Major Tom Stuart" erkannt würde. Auch könnte dieser Algorithmus im Falle von mehreren Suchbegriffen, wie z.B. John Major, alle Artikel finden, in denen jemand erwähnt wird, der "John" heißt und alle Texte, in denen das Wort "Major" (als Eigenname oder als militärischer Rang) auftritt. Natürlich sollten auch Artikel gefunden werden, in denen beide Begriffe auftauchen, doch müssten sie nicht direkt nebeneinander im Text stehen, aber sie könnten es theoretisch. Die Implementation des Algorithmuses schreibt nicht vor, dass im Text zuerst "John" gefolgt von "Major" auftreten muss, was die Ergebnismenge für diese Zwecke unnötig vergrößert und die Präzision der Treffer deutlich verschlechtert.

Jedoch könnte womöglich der dritte Ansatz, welcher die Statistik miteinbezieht, relativ gute Antworten liefern, wenn man noch zusätzlich die Begriffe prime und minister mit in die Anfrage aufnehmen und auf sehr langen Dokumenten arbeiten würde. Dabei könnte man beispielsweise Nominalphrasen von der Art wie "the prime minister, John Major" oder "the French prime minister" extrahieren. Das sind äußerst zufriedenstellende Ergebnisse, wenn man an das anfänglich gesteckte Ziel, die Erkennung von Eigennamen, denkt. Somit ist der statistische Ansatz, der auf keinerlei grammatikalischen Beschreibungsmethoden basiert, nicht zu verachten.

Deshalb hat Jean Senellart zusammen mit Maurice Gross versucht eine neue Methode zu entwickeln, um den statistischen Ansatz zu verbessern. In dem 1998 veröffentlichen Artikel "Nouvelles bases pour une approche statistique." [46] beschreiben sie die Möglichkeit einen Vorverarbeitungsschritt vor das statistische Matching zu schalten. Bei dieser Vorverarbeitung soll der Text zunächst nach Mehrwortlexemen, also mehreren Wörtern, die zusammen eine lexikalische Bedeutungseinheit bilden, durchsucht werden, so dass ungefähr 50% des Textes schon semantisch annotiert wurde. So kann es später beim statistischen Suchen auf dem Text nicht mehr möglich sein, dass z.B. die Wortgruppe "prime minister" oder "energy minister" bei der alleinigen Suche nach "minister" getrennt wird.

Obwohl diese erfolgreiche Zusammenarbeit von linguistischen mit statistischen Methoden einen sehr vielversprechenden Eindruck vermittelt, entschied sich Senellart bei seinem Vorhaben ganz auf die Dienste der Statistik zu verzichten und einen reinen Grammatik und Wörterbuch gestützten Ansatz zur Erkennung von Eigennamen und Berufsbezeichnungen in Nominalphrasen zu wählen.

Auf der Basis großer Lexika mit Eigennamen und Berufsbezeichnungen und unter Verwendung von Transduktoren sollten Grammatiken für die englische Sprache entstehen, welche Satzteile mit Personennamen oder Berufsbezeichnern formal und vollständig beschreiben.

<sup>&</sup>lt;sup>15</sup>vgl. [81] Seite 1213

#### 3.2.2.3 Funktionsweise des Algorithmuses

Der Algorithmus lässt sich in drei große Verarbeitungsschritte unterteilen.

- 1. Zunächst werden die Wörterbücher für die Eigennamen und die lokalen Grammatiken, welche die Berufsbezeichungen beschreiben, auf das Korpus angewandt. Dabei werden semantische Relationen wie Synonymie und Hyponymie und die Zeitlinie der Textsammlung formal definiert und als solches eingesetzt. Damit man die Ergebnisse dieses Schrittes in Echtzeit zurückgeliefert bekommt, wird auf einem zuvor konstruierten Index der Datensammlung gearbeitet.
- 2. In dieser Phase werden die erkannten Eigennamen im Transduktor durch Variablen ersetzt und die gefunden Eigennamen werden zum Auffinden neuer Eigennamen verwendet, die dann dem Benutzer als neue Wörterbucheinträge angeboten werden. Dadurch kann das Erstellen von neuen Transduktoren und die Erweiterung der Wörterbucheinträge weitgehend automatisiert werden.
- 3. Zum Schluss werden die erkannten Nominalphrasen automatisch in andere (natürliche) Sprachen übersetzt, indem entsprechende Transduktoren für die jeweilige Sprache generiert werden.

#### 3.2.2.4 Einblick in die formalen Beschreibungsmethoden

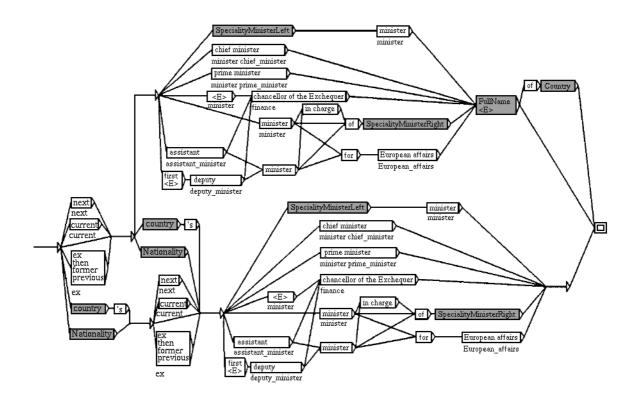


Abbildung 3.6: MinisterOccupation.grf aus Senellart, 1998 [81]

Abbildung 3.6 zeigt eine lokale Grammatik<sup>16</sup> in Form eines Finite-State-Graphen (FS Graph). Ein FS Graph ist im Grunde nur die graphische Repräsentation eines Finite-State-Transducer (FST). Jeder einzelne Knoten stellt die jeweilige Eingabesequenz dar, die der Automat an dieser Transition akzeptiert. Unter manchen Knoten befinden sich Markierungen, welche die Ausgabesequenzen für den entsprechenden Input im Knoten darüber illustrieren. Der Startzustand des Transduktors wird durch den Linkspfeil markiert, wohingegen der Endzustand als doppeltes Quadrat angedeutet wird. Hat ein Knoten einen leicht grauen Hintergrund, so heißt das, dass er einen Subtransduktor aufruft - eine Schreibweise, die es ermöglicht, die Übersichtlichkeit der Automaten zu gewährleisten. Natürlich ist es auch möglich, dass ein Subtransduktor einen Output hat, der dann in die Ausgabe des Haupttransduktors miteinbezogen wird. Ein Knoten, der ein <E> beinhaltet, symbolisiert die leere Transition.

Mit Hilfe dieser Darstellungsformalismen lassen sich linguistische Konstrukte recht einfach darstellen, da z.B. das System Unitex auch einen Grapheneditor bietet, mit dem sich solche Graphen leicht erstellen lassen. Außerdem sind diese FSTs besser als gewöhnliche FSTs, da die Subgraphen sich auf den Hauptgraphen - also auf den Kontext davor oder danach - beziehen können, so dass man mit ihnen auch kontextsensitive Wörter<sup>17</sup> des Typs  $a^nb^n$  erkennen kann.

An diesem konkreten Beispiel aus Abbildung 3.6 soll die Problematik, der formalen Beschreibung von Nominalphrasen, welche sich auf das Wort "minister" beziehen, behandelt werden. Dieser Graph erkennt beispielsweise die Sequenz minister for European affairs, aber er würde nicht French minister for agriculture matchen. Somit wäre dieser Graph sicher noch ausbaufähig.

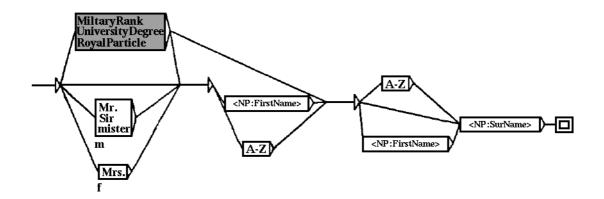


Abbildung 3.7: FullName.grf aus Senellart, 1998 [81]

Der Graph aus Abbildung 3.7 illustriert die Erkennung von Personennamen, wobei er sich "Wörterbuch-Look-Ups" zu nutzen macht.

 $<sup>^{16}</sup>$ vgl. [81] Seite 1214

<sup>&</sup>lt;sup>17</sup>siehe [49]

Die Knoten, welche  $\langle PN:FirstName \rangle$  oder  $\langle PN:SurName \rangle$  enthalten, stehen für Wörterbucheinträge, wo der Transduktor an dieser Stelle alle potentiellen Vornamen oder Nachnamen aus den Lexika mit dem Text abgleicht. Deshalb ist die Ausgabe dieses Automaten ein Nachname, eventuell ein Vorname, und wenn vorhanden m oder f für das ermittelte Geschlecht der Person, wobei das Geschlecht über die Anrede Mr, Sir, mister, Mrs ermittelt wird.

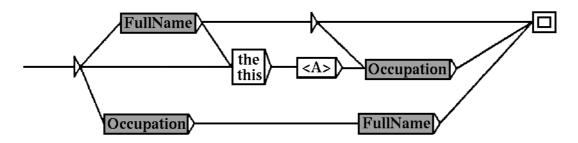


Abbildung 3.8: NounPhrases.grf aus Senellart, 1998 [81]

Der NounPhrases Graph in Abbildung 3.8 vereinigt die Subgraphen Occupation.grf und FullName.grf und stellt somit die syntaktische Beziehung zwischen diesen beiden semantischen Klassen der Eigennamen und Berufsbezeichnungen her. Dabei ist anzumerken, dass <A> stellvertretend für alle Adjektive steht, die dem Standardwörterbuch bekannt sind. Somit würde dieser Automat u.a. Phrasen wie software engineer Tom Mitchell oder Harry Smith the fantastic cook erkennen.

#### 3.2.2.5 Schwächen des Ansatzes

Natürlich beschreibt das komplette Graphenpaket nicht alle syntaktischen Möglichkeiten, wie Personennamen zusammen mit Berufsbezeichnungen auftreten können. Dennoch versucht es nahezu alle einfachen Konstruktionen abzudecken. Beispielsweise würde der NounPhrases Graph aus Abbildung 3.8 nicht auf dem Satz "Mr. Smith, who is since 1978, the chairman of ..." matchen, da der Nebensatz im Graphen nicht berücksichtigt wird. Auch andere Einschübe dieser Art sind kompliziert zu erfassen und werden in diesem Automaten außer Acht gelassen.

Eine andere Schwierigkeit besteht darin, dass eine Person mehrere Berufe ausüben kann, und somit besteht keine Möglichkeit, eine eindeutige Zuordnung zwischen Person und Beruf zu machen. Denn sie wird eventuell an einer Stelle im Text mit einer Berufsbezeichnung und an einer anderen Position im Dokument mit einem anderen Beruf referenziert. Dadurch kann das System nicht gewährleisten, dass bei folgender Zuordnung

SurName=Mitchell, FirstName=Tom, Gender=m, Occupation=cook SurName=Mitchell, FirstName=Tom, Gender=m, Occupation=hotel manager Tom Mitchell ein und dieselbe Person ist.

Auch sind Adverbiale, die an fast jeder Stelle im Satz auftreten können, schwer einer Berufsbezeichnung als deren Ergänzung zuzuordnen. Im Beispiel "In China, the first minister has ..." ist es selbst für den menschlichen Betrachter schwierig, die Ortsergänzung "In China" der Berufsbezeichung "first minister" zuzuordnen.

#### 3.2.3 Bootstrapping zur Erkennung von Nominalphrasen mit FSTs

Jean Senellart stellt in seinem Aufsatz "Tools for locating noun phrases with finite state transducers"[82] verschiedene praxisnahe Verfahren (Tools) vor, wie man relativ schnell eine große Datenbasis endlicher Automaten (FSTs) aufbauen kann, welche Eigennamen und Berufsbezeichner in Nominalphrasen lokalisieren.

#### 3.2.3.1 Die Vorgehensweise

Anfangs wird nur ein Wort ausgewählt, zu dem die Konkordanz über den Text erstellt wird. In Senellarts Fall war es das Wort "officer", was den Ausgangspunkt für die Graphenkonstruktion dargestellt hat. Mit Hilfe der Konkordanzen konnte er z.B. feststellen, welche militärischen Ränge im Korpus zusammen mit officer vorkamen und so Subgraphen erstellen, welche dies abdeckten. Des Weiteren traten auch Adjektive und Nomen im Kontext von officer auf, welche Staatszugehörigkeiten ausdrücken. Um diese Ergänzungsmöglichkeiten bzw. Spezifikationen nicht zu verlieren, entschied er sich sie in Form von Wörterbüchern zu kodieren. Auf diese Weise können die gewonnenen Erkenntnisse vielschichtig eingesetzt werden und sind nicht nur an diesen Kontext gebunden. Gleiches gilt für die gesammelten Kontexte von officer, denn auch diese konnten bei anderen Berufsbezeichnungen mit Erfolg angewandt werden.

Nachdem man Graphen zu einem bestimmten Schlüsselwort erstellt hat, ist es auch möglich diese Graphen zu verwenden,um neue Begriffe zu finden, die den gleichen Kontext, wie das Schlüsselwort aufweisen. Dafür muss lediglich der ursprüngliche Schlüsselbegriff durch eine Variable ersetzt werden. Der neue Graph liefert dann beim Matching alle Ergebnisse, die der alte Graph gefunden hat, und neue Treffer, wo nun an der Stelle von officer andere Berufsbezeichner gefunden wurden. Mit dieser Methode lassen sich leicht Synonyme, Hyponyme oder andere semantisch ähnliche Wörter zum Ausgangsbegriff finden.

Mit Hilfe dieser neuen Kandidaten lassen sich die beiden eben genannten Schritte wiederholen und so entsteht eine Dynamik in der Grammatikentwicklung, bei der aus alten Ergebnissen, immer neuere, bessere und ausbaufähigere Resultate entstehen. Genau diesen dynamischen Entstehungsprozess versteht man als Bootstrapping.

#### 3.2.3.2 Das Ergebnis

Insgesamt hat Jean Senellart mehr als 200 verschiedene Graphen konstruiert, um so viele Berufsbezeichnungen wie möglich abzudecken, und seine Lexika zu Nachnamen und Städten enthielten jedes für sich einige tausend Einträge. Seiner Meinung nach war das auch ein praktischer Beweis dafür, dass die Entwicklung lokaler Grammatiken sehr effizient sein kann, wenn diese nahe am Text verläuft. Doch diese Effizienz wurde auch durch die entsprechende Software gefördert, denn ein Graphen-Editor, ein Index basierter Parsing Algorithmus, sowie ein Konkordanzprogramm und diverse Debugging Möglichkeiten sind Dinge, die bei der Konstruktion von FSTs sehr hilfreich sein können.

#### 3.3 Nathalie Friburger

Nathalie Friburger - Maître de conférences, Enseignant et Chercheur. 18

#### 3.3.1 Zur Person

Nathalie Friburger hat derzeit am Institut für Informatik an der Universität François-Rabelais von Tours, Frankreich, einen Lehrauftrag inne. Bereits im Dezember 2002 wurde ihr an der Universität von Tours der Doktorgrad in Informatik verliehen. Ihr damaliger Doktorvater war Prof. Dr. Denis Maurel, mit dem sie zusammen bis zum Zeitpunkt ihrer Promotion einige Fachartikel zum Thema Methoden in der Eigennamenerkennung publiziert hat. Seit 2003 arbeitet sie in der Lehr- und Forschungseinheit Bases de données et Traitement des langues naturelles (BdTln), die Prof. Maurel seit 1996 leitet. Ihr Forschungsschwerpunkt ist die automatische Verarbeitung von Eigennamen, wobei sich dies bis 2005 auf das Projekt Prolex konzentrierte. Prolex war ein Projekt, dass die Entwicklung von multilingualen Lexika für Eigennamen zum Ziel hatte.

#### 3.3.2 Erkennung von Eigennamen in Zeitungstexten

Nathalie Friburger hat sich im Zuge ihrer Dissertation der Erkennung von Eigennamen in Nachrichtentexten gewidmet. Ähnlich wie Jean Senellart [81][82] wählt sie einen Wörterbuch gestützten Ansatz zur Erkennung von Eigennamen mittels Transduktoren. Um einen kurzen Überblick zu geben, wie sie an dieses Thema herangeht, werde ich nun die Arbeit von ihr und Denis Maurel Elaboration d'une cascade de transducteurs pour l'extraction des noms personnes dans les textes [26] und ihre Doktorarbeit Reconnaissance automatique des nomes propres - Application à la classification automatique de textes journalistiques [25] vorstellen.

#### 3.3.2.1 Kaskadierung<sup>21</sup> von Transduktoren zur Extraktion von Eigennamen [26]

Hierbei handelt es sich um ein Programm, welches Personennamen in französischen Zeitungsberichten erkennt.

Transduktoren sind im Grunde auch nur endliche Automaten, welche allerdings ein Eingabe- und ein Ausgabealphabet haben. In diesem Fall besteht das Eingabealphabet aus Mustern, die im Korpus gefunden wurden, und das Ausgabealphabet fügt den mit den Mustern erkannten Passagen die passende XML-Information hinzu. In der Regel

 $<sup>^{18}</sup> siehe \ \mathtt{http://www.li.univ-tours.fr/Perso\_Frames.asp?Num=44\&lg=frames.asp?Num=44\&lg=frames.asp?Num=44&lg=frames.asp.num=44&lg$ 

<sup>&</sup>lt;sup>19</sup>englische Übersetzung des Titels: Finite-state transducer cascades to extract named entities in texts.

<sup>&</sup>lt;sup>20</sup>englische Übersetzung des Titels: Automatic Recognition of Proper Names - An Application in Automatic Clustering of Journalistic Texts.

 $<sup>^{21} \</sup>mathrm{In}$  diesem Fall bedeutet Kaskadierung das Zusammenschalten von verschiedenen Transduktoren (Reihenschaltung von Transduktoren).

sind die gefundenen und XML annotierten Sequenzen Personennamen und ihre jeweiligen Kontexte, was folgendes Beispiel illustriert.

Doch bevor man die Transduktoren nacheinander auf das Korpus anwenden kann, sind einige Vorverarbeitungschritte notwendig, welche vom System INTEX [84] übernommen werden. <sup>23</sup> Zu den wichtigsten Phasen zählt u.a. die Satzenderkennung, welche auf dem ganzen Text durchgeführt wird. Im gleichen Schritt werden auch die Satzendmarkierungen in den Originaltext eingefügt. Später erfolgt die Anwendung sämtlicher Wörterbücher auf das Korpus, wobei jedes Wort mit allen Formen, die in einem der Lexika auftreten, markiert wird. An dieser Stelle wird noch keine Disambiguierung durchgeführt. Doch die annotierten Wörter besitzen nun sämtliche grammatikalischen und semantischen Informationen, die in den Lexika kodiert sind.

Dabei kamen die Wörterbücher

- *DELAS* ein Lexikon, welches die gesamte grammatikalische Information für "einfache Wörter"<sup>24</sup> festhält,
- *Prolintex*, ein Toponymlexikon<sup>25</sup>,
- Prenom-prolex, ein Wörterbuch für Vornamen,
- und ein Lexikon für Berufsbezeichnungen, das von Cédrick Fairon von der Universität Marne la Vallée erstellt wurde.

#### zum Einsatz.

Das Prinzip der Kaskadierung von Transduktoren ist im Grunde recht einfach zu erklären. Die Transduktoren müssen in einer aufeinander abgestimmten Reihenfolge nacheinander auf den Text angewendet werden. Denn oft ist der Output eines Transduktors, der Input - das zu Suchende - für den darauffolgenden Transduktor. Jede gefundene Sequenz wird markiert (siehe obiges Beispiel) und kann durch diese Markierung über den Index gefunden werden. Auch muss jedes erkannte Muster aus dem Text gelöscht werden, da sonst die Gefahr besteht, dass ein später geschalteter Transduktor es nochmal erkennt. So wird vermieden, dass Passagen mehrfach erkannt werden, und dass das System ineffizient arbeitet.

<sup>&</sup>lt;sup>22</sup>aus Friburger, N. & Maurel, D. (2001) [26]

<sup>&</sup>lt;sup>23</sup>Das System Unitex[75] bietet ähnliche Funktionen zur Korpusbearbeitung.

 $<sup>^{24},\!</sup>$ simple words" - im Gegensatz zu "compound words" (Mehrwortlexemen), welche das DELAC auflistet.

<sup>&</sup>lt;sup>25</sup>Ein **Toponym** bezeichnet einen Ortsnamen im allgemeinen Sinne. Hierunter versteht man also insbesondere die Bezeichnungen bestimmter Gebiete, Verwaltungseinheiten, Siedlungen, Verkehrswege, Gewässer und alle übrigen topographischen Objekte mit Eigennamen. [vgl. http://de.wikipedia.org/wiki/Toponym]

Eine Voraussetzung muss noch erfüllt werden, bevor die Transduktoren in Reihe geschaltet werden können. Es ist auch wichtig, sich einen Bestand an linken und rechten Kontexten der Personennamen, die in dem Zeitungskorpus vorkommen, aufzubauen. Denn das Matchen über die Kontexte von Personennamen stellte sich bei französischen Texten als äußerst hilfreich heraus, weil ungefähr 90% aller Personennamen in Nachrichtentexten über ihren linken Kontext erkannt werden können. Ein Grund dafür könnte sein, dass gewisse Stilkonventionen zur Behandlung von Personennamen in Printmedien bestehen, so dass eindeutige, fast standardisierte Muster erkennbar waren.

Mit Hilfe eines annotierten Korpus der französischen Zeitung  $Le\ Monde$ , der ungefähr 165000 Wörter umfasste ( $Ouest\ France$  enthielt 67000 Wörter) war es möglich, die häufigsten Kontexte von Personennamen im Text zu kategorisieren.

- In 25,9% (17,1% für *Ouest France*) der Fälle ging dem Personennamen ein Titel oder eine Berufsbezeichnung gefolgt von einem Vornamen und/oder einer Staatsangehörigkeit voran. (Fall 1)
- In 19,1% (16,3% für *Ouest France*) der Fälle ging dem Nachnamen ein Berufsbezeichner oder ein Titel zusammen mit einer Bezeichnung für eine Staatsangehörigkeit oder ein dem Wörterbuch unbekannter Vorname zusammen mit einer Nationalität voran. (Fall 2)
- Am häufigsten mit 43,4% (59,0% für *Ouest France*) hat der Kontext eher eine beschreibende Funktion zum jeweiligen Personennamen, d.h. er wird meist attributiv eingesetzt und der Personenname besteht in der Regel aus einem dem Lexikon bekannten Vor- und Nachnamen. (Fall 3)
- Setzt man den Kontext zur Erkennung der Eigennamen ein, so helfen Berufsbezeichnungen oder Verben der Äußerung wie z.B. sagen oder erklären dabei, 5,2% (2,2% für Ouest France) aller Personen im Text zu erkennen. Natürlich können Verben der Äußerung auch ohne ein menschliches Subjekt im Satz auftreten, d.h. diese Kontexte sind mit Vorsicht zu genießen. (Fall 4)
- Die übrigen 6,4% (5,4% für Ouest France) der Personennamen weisen keinerlei hilfreiche Kontexte auf, so dass diese nutzlos bei der Suche sind. Diese Personen sind in der Regel so berühmt, dass der jeweilige Autor es wahrscheinlich nicht für nötig gehalten hat, die Persönlichkeit vorzustellen, oder ein paar einleitende Worte zu ihr zu schreiben. Doch ca. die Hälfte dieser anscheinend nicht im Text zu findenden Leute, werden an anderen Stellen im Korpus nochmal namentlich erwähnt, so dass im Endeffekt nur noch 3,3% der ursprünglich 6,4% Personennamen unerkannt bleiben. Eventuell könnte ein Lexikon aller berühmten Personennamen diesen Prozentsatz auch noch verringern. (Fall 5)

Dass die Trefferquoten im ersten und im zweiten Fall für *Ouest France* kleiner als für *Le Monde* ausfallen, liegt wohl an strikteren Schreibkonventionen, die für die Journalisten von *Le Monde* bestehen. Somit können die vordefinierten Muster erfolgreicher auf *Le Monde* als auf *Ouest France* suchen.

Um die gute Trefferquote ihres Ansatzes der Kaskadierung von Transduktoren nachzuweisen, wandte Nathalie Friburger 14 Transduktoren in Reihe geschaltet nacheinander auf einen Teilkorpus von *Le Monde* an, der etwa 80000 Wörter umfasste.

Dabei ergaben sich je nach Fall (siehe vorige Seite) folgende Ergebnisse:

	Fall 1	Fall 2	Fall 3	Fall 4	Fall 5	Gesamt
Tatsächliche Anzahl der	253	187	424	50	64	977
Personennamen im Text						
Anzahl der gefundenen	245	187	413	32	32	909
Personennamen im Text						
Anzahl der korrekt gefundenen	242	186	410	30	31	899
Personennamen						
Recall	95,7%	99,5%	96,7%	60,0%	48,4%	91,9%
Precision	98,8%	99,5%	99,3%	93,8%	96,9%	98,7%

Tabelle 3.1: Statistische Auswertung der Extraktionsresultate aus Friburger (2001) [26]

Mit den Resultaten in den ersten der drei Fälle kann man sehr zufrieden sein. Doch leider weist Fall 4 einen schlechten Recall auf, was wohl mit der problematischen Erkennung von Eigennamen in ambigen Kontexten zu tun hat. Fall 5 behandelt nur Namen, die ohne einen spezifischen Kontext im Korpus auftreten, und die nur durch die Gesamtbedeutung des Satzes oder durch das Wissen eines menschlichen Lesers identifiziert werden können. Da einige dieser Personen schon in anderen Textpassagen vorkamen, ist es überhaupt möglich einen Recall von 48,4% zu erreichen.

Abschließend kann man sagen, dass das Prinzip der Kaskadierung von Transduktoren recht einfach und effektiv bei der Suche nach Personennamen sein kann. Dagegen gehen Nathalie Friburger und Denis Maurel davon aus, dass die Extraktion von anderen Eigennamen, wie Orts- und Organisationsnamen, wesentlich schwieriger ist, weil ihre jeweiligen Kontexte im Korpus nicht so schematisch sind wie die von Personennamen.

#### 3.3.2.2 Eigennamen bei der Klassifikation von Nachrichtentexten

Nathalie Friburger hat die Idee der Reihenschaltung von FSTs auch für ihre Dissertation eingesetzt. Vollständigkeitshalber möchte ich noch kurz die Thematik ihrer Doktorarbeit ansprechen, bei der das System casSys zum Einsatz kam, welches die Kaskadierung von Transduktoren implementiert. Das Ziel ihrer Arbeit war nicht nur die automatische Extraktion von Eigennamen, sondern auch die automatische Klassifikation von Zeitungstexten anhand der darin auftauchenden Namen. Das dafür eingesetzte Programm extractNP, welches casSys verwendet, ermöglicht es Ambiguitäten aufzulösen, sowie Eigennamen zu segmentieren und kategorisieren. Das System lieferte hervorragende Ergebnisse, so dass eine Präzision von 94% und ein Recall von 93% erzielt wurde. Des Weiteren entwickelte sie eine Anwendung, welche sich die verschiedenen Vorkommen von Personennamen zu nutze macht, um Zeitungsnachrichten nach Thematiken zu kategorisieren. Dabei stellte sich heraus, dass dieser Ansatz ein qualitativ gutes Clustering von Zeitungstexten möglich machte.

## 3.4 Friederike Mallchok

#### 3.4.1 Zur Person

Friederike Mallchok lebt heute in Flintham, Nottinghamshire, UK, und ist bei US amerikanischen und britischen Firmen unter Vertrag, die sich hauptsächlich mit der multikulturellen Namensverarbeitung oder der automatischen Wissensgewinnung beschäftigen. Die 1976 in München als Friederike Schmidt geborene heutige Frau Mallchok studierte am Centrum für Informations- und Sprachverarbeitung (CIS) der Ludwig-Maximilians-Universität (LMU) München Computerlinguistik und beendete 2001 ihr Studium mit dem akademischen Grad des Magister Artium (M.A.). In der Zeit ihrer Promotion setzte sie sich besonders mit der Thematik der automatischen Eigennamenerkennung in Zeitungstexten auseinander, wobei sie ihren Schwerpunkt in der Analyse von Organisationsnamen sah. Ihre Untersuchungen auf diesem Gebiet geben Aufschluss über das Verhalten von Organisationsnamen in Texten wie Wirtschaftsnachrichten und beschreiben ihre syntaktische und auch semantische Eingliederung in den Text. Im Juli 2004 wurde ihr der Doktorgrad der Philosophie an der Ludwig-Maximilians-Universität München verliehen.<sup>26</sup>

## 3.4.2 Erkennung von Organisationsnamen in Wirtschaftsnachrichten

#### 3.4.2.1 Motivation des Ansatzes

In ihrer Doktorarbeit "Automatic Recognition of Organization Names in English Business News" hatte sich Friederike Mallchok zum Ziel gesetzt, nachzuweisen, dass sich die Genauigkeit und Performanz der Eigennamenerkennung wesentlich verbessern lässt, wenn man einen sprachspezifischen Ansatz dafür wählt. Unter einem sprachspezifischen Ansatz versteht man einerseits die Beschränkung der Trainingskorpora auf eine bestimmte Domäne, wie z.B. den Bereich der Wirtschaftsnachrichten, und andererseits aber auch eine Einschränkung bei der Named-Entity-Recognition (NER). Wenn man eine Named Entity (benannte Entität), wie hier die Organisationsnamen, in den Vordergrund rückt, und dann ausgehend von dieser bestimmten Klasse der Eigennamen ihre Kontexte untersucht, finden sich weitere Eigennamen und noch weitere wertvolle Informationen in ihrem Umfeld. Um die Kontexte der Organisationsnamen syntaktisch und semantisch beschreiben zu können, wählte Friederike Mallchok die formale Repräsentation der lokalen Grammatiken. Dabei verzichtet sie vollständig auf statistische Methoden zur Extraktion von Eigennamen und verlässt sich ganz auf die Identifikation von Organisationsnamen durch ihre jeweiligen Kontexte und das in Wörterbüchern kodierte Zusatzwissen.

<sup>&</sup>lt;sup>26</sup>vgl. http://www.cis.uni-muenchen.de/~schmidt/FM\_Lebenslauf\_Jan\_2006.pdf

<sup>&</sup>lt;sup>27</sup>Automatische Erkennung von Organisationsnamen in Englischsprachigen Wirtschaftsnachrichten

#### 3.4.2.2 Einsatz von Ressourcen: Korpora und Lexika

Da Friederike Mallchok sich dazu entschlossen hatte, die für Wirtschaftsnachrichten typische "Subsprache"<sup>28</sup> zu untersuchen, fiel ihre erste Wahl auf das frei verfügbare Reuters Korpus<sup>29</sup>. Das Reuters Korpus enthält alle Nachrichtentexte (ca. 810000), welche die Nachrichtenagentur Reuters Ltd. vom 20. August 1996 bis einschließlich 19. August 1997 veröffentlicht hatte.

Nachdem dieses Korpus für Wirtschaftsnachrichten nicht mehr auf dem neuesten Stand war, ergänzte sie ihre Textsammlung durch Online-Ausgaben der "Financial Times", des "Wall Street Jounals", von "Newsday", der "New York Times" und durch aktuelle Artikel der "Reuters News". Denn gerade für die Erkennung von Firmennamen ist es wichtig, aktuelle Informationen über die Unternehmen vorliegen zu haben. Beispielsweise können einerseits dem Lexikon bekannte Firmen, welche nach 1997 gegründet wurden, nicht im Reuters Korpus gefunden werden, und andererseits werden nur Namen von Organisationen in diesem Text lokalisiert, welche in dieser Zeitspanne in den Nachrichten präsent waren. Unter Hinzunahme der eben angesprochenen elektronischen Nachrichtenausgaben konnten auch junge, aufstrebende oder immer noch bedeutende Unternehmen in den aktuellen Texten erkannt werden.

Wie eben kurz erwähnt, wurden mehrere semantische Lexika unterstützend zur Erkennung der Organisationsnamen in den Korpora eingesetzt. Mittels dieser Wörterbücher sollte das Auffinden von Firmennamen im Text wesentlich erleichtert werden.

Mit Hilfe diverser Internetressourcen konnte Friederike Mallchok ein beachtliches Begriffsinventar für ihr Organisationsnamenlexikon (ONL) und für ihr Organisationsbeschreibungslexikon (ODL) zusammenstellen. Die jeweiligen Namen der Lexika lassen natürlich schon auf ihren Inhalt schließen: Das ONL enthält ausschließlich Firmennamen, und das ODL führt eine Reihe an Beschreibungen für Unternehmen auf, welche oft in Wirtschaftstexten den Organisationsnamen einleiten.

Zudem ließen sich im Kontext von Organisationsnamen relativ oft Berufsbezeichner finden, welche in einem Berufsbezeichnerlexikon (HPL) archiviert wurden. Des Weiteren konnten Ortsbezeichnungen wie Länder, Städte und Staaten, sowie Zeitangaben in den entsprechenden Lexika gespeichert werden. Auch allgemeine Kontexte der Firmennamen extrahierte Friederike Mallchok aus den Korpora und bewahrte sie in Wörterbüchern auf. Dabei wurden nur die textuellen Umgebungen von Organisationsnamen ins Lexikon übernommen, welche besonders häufig in den Korpora vorkamen.

#### 3.4.2.3 Entwicklung lokaler Grammatiken

Basierend auf den eben genannten Lexika entwickelte sie lokale Grammatiken, welche einerseits die interne Struktur von Organisationsnamen repräsentierten und andererseits auf ihre Funktion im Satz eingehen bzw. ihr syntaktisches Verhalten in Wirtschaftsnachrichten wiederspiegeln. Die verschiedenen Grammatiken sollten so viele syntaktische

<sup>&</sup>lt;sup>28</sup>Einschränkung der Sprache auf eine bestimmte Bezugsdomäne wie z.B. Wirtschaftsnachrichten, sowie Dominanz von Fachvokabular.

<sup>&</sup>lt;sup>29</sup>http://about.reuters.com/researchandstandards/corpus/

Variationen wie möglich abdecken, in denen Unternehmen vorkommen können. Somit werden beispielsweise die Möglichkeiten, Organisationsbeschreibungen den Firmennamen voran- oder nachzugestellen, sowie Berufsbezeichner - eventuell in Verbindung mit einem Personennamen - im linken oder rechten Kontext der Organisationsnamen zu nennen, in den lokalen Grammatiken berücksichtigt. Den Schwerpunkt ihrer Untersuchungen legte sie auf die Wirtschaftereignisse "joint venture", "merger" und "partnership", für die sie zeigen wollte, dass eine erweiterte Indexierung durch lokale Grammatiken, welche diese Phänomene beschreiben, durchaus möglich ist und später für eine intelligente und effektive Suche eingesetzt werden kann.<sup>30</sup>

### 3.4.2.4 Bootstrapping und Akronymbildung

Wie bereits Maurice Gross (siehe Abschnitt 3.1.2) und Jean Senellart (siehe Abschnitt 3.2.3) sich des Bootstrappings bei der Entwicklung von Transduktoren (Finite-State-Transducern) bzw. lokaler Grammatiken mit dem System Unitex [75] bedient haben, verwendet auch Friederike Mallchok diese Methode zur Verbesserung ihrer Ergebnisse. Auf diese Weise stellte sich in mehreren Nachbearbeitungsschritten schnell heraus, welche Fehlerquellen noch in den Grammatiken vorlagen, und wie diese minimiert werden konnten. Zusätzlich generierte sie aus den Organisationnamen, die aus mehreren Wörtern zusammengesetzt waren, mögliche Akronymvarianten, welche später im Korpus verifiziert wurden. Außerdem wurden noch weitere Abkürzungsmöglichkeiten für die entsprechenden Firmennamen berücksichtigt und ihre Existenz auf dem Text überprüft. Bei der erfolgreichen Überprüfung wurden die Varianten der Organisationsnamen in das Lexikon aufgenommen und im Korpus annotiert.

#### 3.4.2.5 Fazit der Arbeit

Mit dem entwickelten System ist es Friederike Mallchok gelungen mit einer hohen Genauigkeit und guten Performanz, Organisationsnamen in englischsprachigen Wirtschaftsnachrichten zu erkennen. Dabei war es ihr möglich zu zeigen, dass das Ergebnis der Eigennamenerkennung signifikant verbessert werden kann, wenn jede Sprache, jede Domäne und jede Art von Entität getrennt behandelt wird. Des Weiteren widerlegte sie die Annahme von vergleichbaren NER<sup>31</sup>-Systemen, dass die Verwendung von Kontextinformationen nur zur Lokalisierung von Entitäten sinnvoll ist. Ihr Ansatz bewies, dass durch den Einsatz von lokalen Grammatiken weitere Informationen über die entsprechenden Entitäten aus den Korpora gewonnen werden können. So dienen beispielsweise semantisch kategorisierte Organisationsbeschreibungen dazu, die Entitäten, die sie beschreiben oder sogar die Texte oder Textabschnitte, in denen diese vorkommen, zu klassifizieren. Ihre Bemühungen auf dem Gebiet der automatischen Erkennung von Organisationsnamen in Wirtschaftsnachrichten brachten sie letztendlich zu dem Schluss, dass eine Weiterentwicklung dieser lokalen Grammatiken auf jeden Fall sinnvoll ist. Dadurch könnte später eine breitere Abdeckung auf der Domäne der Wirtschaftnachrichten erreicht werden.

<sup>31</sup>Named-Entity-Recognition

 $<sup>^{30}\</sup>mathrm{vgl.\ http://www.cis.uni-muenchen.de/``schmidt/lg/Deutsche\_Zusammenfassung.pdf}$ 

# 4 Beschränkungen im System

Wie bereits Friederike Mallchok so treffend in ihrer Doktorarbeit [55] bemerkt hatte, ist einer der größten Vorteile von lokalen Grammatiken die Modularität. Es ist keinenfalls ein Nachteil sich bei der Erstellung lokaler Grammatiken besonders auf eine bestimmte Entität zu konzentrieren und deren Kontext möglichst genau zu beschreiben.

Der hier vorgestellte Ansatz fokussiert zwar die Erkennung von Menschenbezeichnern und beschränkt sich auf biographische Relationen (siehe Abschnitt 1.2), doch wird sich bei der Entstehung des Systems zeigen, dass auch andere Entitäten in der Umgebung von Personen auftreten und somit berücksichtigt werden müssen.

Dafür wurden lokale Grammatiken entwickelt, die Personen näher spezifizieren, aber auch Organisationen und Toponymen eine gewisse Beachtung schenken. Überdies werden auch Verbrelationen in Form von lokalen Grammatiken beschrieben, welche die verschiedenen Entitäten miteinander logisch und syntaktisch verbinden.

So kann jede Grammatik separat erweitert und auch die darin verwendeten semantischen Hyperonymklassen können jederzeit durch weitere Wörterbucheinträge ergänzt werden. Des Weiteren lassen sich die entstandenen Grammatiken problemlos in andere Named-Entity-Recognition-Systeme integrieren. Auch zur Erkennung anderer Entitäten sollten die Informationen aus diesen lokalen Grammatiken herangezogen werden, so dass man auf diesem Wissen aufbauen und zugleich das System erweitern könnte.

Außerdem sollte die Entscheidung, sich bei der Erkennung von Personen innerhalb biographischer Kontexte auf die Domäne der Wirtschaftsnachrichten zu beschränken, kein Hindernis dafür sein, später die für diesen speziellen Bereich entwickelten Grammatiken für andere Themengebiete auszuweiten. Denn wie die Wahl meiner Korpora zeigen wird, gibt es biographische Relationen, welche äußerst selten in Wirtschaftsnachrichten auftreten, dagegen aber in einer richtigen Biographie kaum fehlen. Es ist nur natürlich, dass nicht jedes personenbezogene Prädikat in einem Wirtschaftstext eine biographische Relation verkörpert, und dass nicht jede biographische Relation in den Nachrichten veröffentlicht wird. Das ist für die Entwicklung lokaler Grammatiken nur insofern ein Problem, wenn Verbrelationen beschrieben werden, welche höchst selten auf dem Trainingskorpus vorkommen. Somit ist die Qualität einer lokalen Grammatik schwer zu messen und alternative Trainingskorpora werden benötigt. Man sollte sich dieser Tatsache immer bewusst sein, dass die Entwicklung lokaler Grammatiken stark von der Domäne des Korpuses abhängt und sein Einfluss auf die Grammatik nicht zu unterschätzen ist.

Auch wenn der hier präsentierte Ansatz sich hauptsächlich auf biographische Relationen konzentriert, die häufig in Wirtschaftsnachrichten vorkommen, soll das nicht heißen, dass diese lokalen Grammatiken nicht auf Texten anderer Bereiche gute Ergebnisse erzielen. Es werden lediglich die Relationen nicht abgedeckt, die kaum oder nie in Wirtschaftstexten genannt werden, was ein Ansporn wäre, das Konzept auszuweiten.

## 4.1 Sprachgebundenheit

Alle hier vorgestellten lokalen Grammatiken wurden für die englische Sprache entwickelt. Sicherlich ist die Entscheidung, für welche Sprache die Erkennung von Personen in biographischen Kontexten implementiert wird, nicht unbegründet getroffen worden. So wurde die Wahl der Sprache sicher durch die große Dominanz des Englischen als Sprache des Internets beeinflusst. Doch auch die Tatsache, dass für das Englische schon sehr viel im Bereich Named-Entity-Recognition (NER) erforscht und entwickelt worden ist, wovon man manches aufgreifen, verbessern oder mit seinem eigenen Ansatz vergleichen kann, spielte eine beachtliche Rolle bei dieser Entscheidung.

Soweit es das Gebiet der lokalen Grammatiken betrifft, wurde die meiste Vorarbeit bei der linguistischen Analyse der französischen und englischen Sprache geleistet.

Des Weiteren ist der Bereich der Wirtschaft ein von Anglizismen geprägtes Feld, was ebenfalls dafür sprechen würde, sich gleich auf die Originalsprache zu konzentrieren.

Überdies ist die Auswahl an Trainingskorpora wesentlich größer, wenn man sich für die Arbeit mit Englisch entscheidet, und bei der Erstellung von Lexika kann im Internet auf ein großes Spektrum an Ressourcen in Form von themenspezifischen Listen zurückgegriffen werden, so dass für das Englische in kürzerer Zeit als für eine andere Sprache eine enorme Wissensbasis zusammengestellt werden kann.

Trotz der Beschränkung auf das Englische bei der Entwicklung lokaler Grammatiken, können die entstandenen Grammatiken mit relativ wenig Aufwand auf andere Sprachen übertragen werden.

## 4.2 Schwerpunkt Wirtschaftsnachrichten

Für computerlinguistische Untersuchungen wurden immer schon gern Korpora herangezogen, welche aus Wirtschaftstexten zusammengestellt waren. Named-Entity-Recognition und Information Retrieval auf Wirtschaftsnachrichten sind in den letzten Jahren immer beliebter geworden, und wenn man an das frei verfügbare Reuters Korpus<sup>32</sup> denkt, das für Studien dieser Art sogar noch aufbereitet wurde, stellt man fest, dass der Bedarf an Informationsextraktion aus wirtschaftlich orientierten Texten bei weitem noch nicht gedeckt ist. Mit der immer stärker werdenden Verflechtung internationaler Wirtschaftsbeziehungen, dem ständig anwachsenden Trend der internationalen Fusionen und der Globalisierung der Wirtschaft, wächst die Nachfrage aus aktuellsten Wirtschaftsartikeln, kurz und prägnant interessante Information zu erhalten. Der Kreis der Suchenden beschränkt sich heute längst nicht mehr nur auf Betriebs- oder Volkswirte, sondern auf jeden, der in die Wirtschaft investieren möchte, und sich aufgrunddessen informiert. All diese Gründe machen Wirtschaftsnachrichten zu einer lukrativen und begehrten Domäne für die Informationsgewinnung und heben die Nachfrage nach qualitativ guten Systemen zur Wissensextraktion auf Nachrichtentexten.

 $<sup>^{32} {\</sup>rm http://about.reuters.com/researchandstandards/corpus/}$ 

## 4.3 Priorisierung von Entitäten

In Kapitel 3 wurden bereits unterschiedliche Ansätze zur Erkennung benannter Entitäten (Named Entities) mittels lokaler Grammatiken vorgestellt. All diese Ansätze haben nicht nur die Gemeinsamkeit, dass sie sprachbasierte statt statistische Methoden zur Lokalisierung von Eigennamen oder Verbgefügen anwenden, sondern auch dass keiner dieser Linguisten versucht hat, alle Kategorien von Entitäten in einem System zur Named-Entity-Recognition zusammenzufassen. Jeder von ihnen hat sich auf eine Entität konzentriert - einige auf Personen und andere auf Organisationen. Natürlich spielten immer wieder andere Entitäten, wie vorallem Toponyme, eine untergeordnete Rolle bei der Erkennung von Menschen oder Firmen. Doch waren sie dann nur Mittel zum Zweck, indem sie Teil des Kontextes der zu suchenden Entität waren.

Personen werden wohl immer eine der beliebtesten Entitäten für die NER sein, auch wenn die automatische Produktnamenerkennung inzwischen immer mehr in den Vordergrund rückt, wie es die Arbeit von Jeannette Roth [79] zeigt. Bis jetzt werden wohl die syntaktischen und semantischen Aspekte von Produktnamen noch unerforschter sein als die Eigenschaften von Organisationsnamen. Dennoch bewies auch die Arbeit von Friederike Mallchok [55], wie gut lokale Grammatiken das Problem der Organisationsnamenerkennung lösen können.

Gerade bei der Suche auf Wirtschaftsnachrichtentexten stößt man auf eine beträchtliche Anzahl von Organisationsnamen. Diese Kategorie der verschiedenen Entitäten wird jedoch in dem hier präsentierten Ansatz eine untergeordnete Rolle zu den Menschenbezeichnern haben. Da aber Personen in Wirtschaftsartikeln sehr häufig im Zusammenhang mit Firmen genannt werden und ihr Verhältnis zu diesen oft explizit beschrieben wird, sollte natürlich den Beziehungen zwischen diesen beiden Entitäten besonders viel Beachtung geschenkt werden. Obwohl diese beiden Gruppen - Personen und Organisationen - die frequentesten Entitäten in Wirtschaftstexten sein werden, gibt es dort noch viele weitere personenbezogene Relationen, in deren Kontext womöglich andere Entitäten wie Ortsbezeichnungen auftreten können.

Für alle Entitäten, die keine Menschenbezeichner sind, werden lokale Grammatiken erstellt, welche dazu dienen das Umfeld der Personen zu spezifizieren. Die Grammatiken entsprechen in ihrem Umfang und in ihrer Ausführlichkeit der Wichtigkeit der Relation, die zwischen der jeweiligen Entität und der Personenbezeichnung herrscht. Somit werden die Grammatiken für die Organisationsbezeichner umfassender als für die Toponyme sein, da sie eine größere Relevanz im Korpus haben.

Für diesen Ansatz gilt, dass die Menschenbezeichner die priorisierte Entität darstellen, doch es wäre für zukünftige Vorhaben keine Schwierigkeit die Gewichtung der Entitäten für die jeweiligen Zwecke abzuändern.

# 5 Ressourcen: Grundlagen des Systems

## 5.1 Korpora

Wie bereits mehrfach erwähnt wurde, sollten Menschenbezeichner innerhalb biographischer Relationen automatisch in Wirtschaftsnachrichten erkannt werden. Diese Vorgabe schränkt die Wahl der Texte, auf denen gearbeitet werden kann, zunächst auf die Wirtschaftsteile vieler englischsprachiger Zeitungen ein. Nur wer begnügt sich mit dem Wirtschaftsteil, wenn ganze Wirtschaftsblätter ihre Artikel online zur Verfügung stellen?

Ähnlich wie bei Friederike Mallchok (siehe Abschnitt 3.4, vgl. [55]) wäre das Reuters Korpus eine Option gewesen, da es eine Textsammlung aus Wirtschaftsartikeln ist. Doch die Tatsache, dass es für Wirtschaftsnachrichten relativ veraltet ist, machte es zu keinem Kandidaten für einen Testkorpus.

Dagegen war das Angebot vom Centrum für Informations- und Sprachverarbeitung der LMU München, mir eine Jahresausgabe der Financial Times zur Verfügung zu stellen, wesentlich interessanter. Vorallem handelte sich hierbei um die Jahresausgabe 2004 der FT, womit sicher gestellt ist, dass die darin enthaltenen Informationen relativ aktuell sind.

#### 5.1.1 Financial Times

Die Financial Times<sup>33</sup> ist eine Tageszeitung, welche fast täglich herausgegeben wird. In ihrem elektronischen Format ist jede Tagesausgabe eine XML-Datei und das Jahr 2004 umfasste 347 Tage, an denen die FT erschienen ist. Somit ergab sich eine Datenmenge von ungefähr 5,8 GB.

Um aus dieser Artikelsammlung einen Korpus zu erstellen, wurden zunächst alle Texte von ihrer XML-Information befreit, was die Größe der Daten auf 4,7 GB verminderte. Im Anschluss wurden die Tagesausgaben monatsweise zusammengefügt, so dass es für jeden Monat eine Datei der Financial Times gab.

Diese 12 Dateien wurden nun für die spätere Bearbeitung mit dem System Unitex [75] vorbereitet:

• Im ersten Schritt wurde die Satzenderkennung mit dem Tokenizer-Programm von Sebastian Nagel auf dem gesamten Text vorgenommen.

Ein Programmaufruf folgender Form

cat <korpus> | tokenizer -L en -SE {S} -P -o <korpus.eos>

 $<sup>^{33} {\</sup>rm http://news.ft.com/home/us}$ 

liefert einen Text mit Satzendmarkierungen, wie ihn das Programm Unitex fordert. Mir wurde dieses Programm in der Version 0.6 überlassen, so dass es für die Satzenderkennung im Englischen angepasst werden konnte, da die deutsche Satzenderkennung ausgereifter als die englische war. Herr Nagel übernahm diese Änderungen in seine Version 0.7, zu der eine Hilfestellung im Anhang A auf Seite 164 angeboten wird.

Das Programm Unitex bietet zwar auch eine Satzenderkennung für das Englische an, doch handelt es sich bis jetzt um die französische Satzenderkennung, die nur leicht für das Englische abgewandelt wurde und leider immer noch größtenteils die französischen Abkürzungen enthält. Somit war diese Satzenderkennung keine Alternative zum Tokenizer-Programm, was wirklich hervorragende Ergebnisse geliefert hat.

- Im nächsten Schritt wurde die Normalisierung und Tokenisierung des Textes mit den entsprechenden Programmen aus dem System Unitex vorgenommen.
- Im letzten Schritt wurde die gesamte grammatikalische und semantische Information aus den im nächsten Abschnitt angesprochenen Lexika im Korpus passend annotiert. Das heißt aber nicht, dass der Originaltext verändert wurde, sondern dass diese Zusatzinformationen in Wortlisten ergänzend zum Text gespeichert werden.

Somit ist der Financial Times (FT) Korpus für die Entwicklung lokaler Grammatiken mit den System Unitex bereit, welche anschließend darauf getestet werden können.

## 5.1.2 Biography.com

Dennoch ist das FT Korpus nicht die einzige Textsammlung, welche zur Validierung der hier vorgestellten lokalen Grammatiken verwendet wird.

Wie bereits in Kapitel 4 angedeutet wurde, ist nicht jede biographische Relation in Wirtschaftsnachrichten wie der Financial Times vertreten. Manche von ihnen sind typisch für Lebensläufe und könnten in einem Korpus basierend auf diversen Biographien sicher gefunden werden.

Aufgrunddessen wurde ein Crawler implementiert, der alle fast 25000 Biographien der Internetseite Biography.com<sup>34</sup> heruntergeladen hat. So standen nochmal knapp 100 MB Daten zum Testen der Grammatiken zur Verfügug.

Hierbei handelte es sich um HTML-Dateien, welche zunächst in reinen Text umgewandelt wurden und daraufhin durchliefen sie die gleichen Bearbeitungsschritte wie das Financial Times Korpus, so dass das Biography.com Korpus im System Unitex verwendet werden konnte.

34http://www.biography.com	
othttn://www.biography.com	
noop://www.biography.com	

## 5.2 Wörterbuchressourcen

Zur Verarbeitung dieser Korpora, wurden mehrere Lexika erstellt. Diese Wörterbücher könnten mit einigen Anpassungen in ihrem Format später auch ohne weiteres in das CISLEX-E [47] integriert werden. Das CISLEX-E ist ein am Centrum für Informations- und Sprachverarbeitung entwickeltes morphologisches, syntaktisches und semantisches Lexikon der englischen Sprache.

Natürlich wird die Erkennung von Menschenbezeichnern in biographischen Kontexten um einiges einfacher, wenn das System später auf eine relativ große Wissensbasis zurückgreifen kann. Durch qualitativ gute Lexikoneinträge lässt sich auch die Performanz und Genauigkeit des Systems deutlich verbessern, da Ambiguitäten leichter aufgelöst werden können, ohne dass der Kontext miteinbezogen werden muss. Das soll nichts anderes bedeuten, als dass ein im Text auftretender Personenname als dieser eindeutig identifiziert werden kann, wenn er bereits im Wörterbuch vorkommt. Wenn der Name also nicht in einem der Lexika kodiert wurde, muss man sich an bestimmten Kontextschemata orientieren, um herauszufinden, ob es sich hierbei um einen Personennamen handelt. Gibt es keinen eindeutigen Kontext, der darauf schließen lässt, wird das Auffinden dieses Namens fast unmöglich.

Muss man sich jedoch an den Kontexten, in denen ein Menschenbezeichner eingebettet sein kann, orientieren, so ist es durchaus hilfreich mehr Entitäten als nur Personennamen zu sammeln und in Lexika festzuhalten.

Ein Personenname kann beispielsweise aus einem Titel oder einer Anrede gefolgt von einem Nachnamen bestehen. So wäre es nur sinnvoll, Wörterbücher allein nur für Titel und Anredemöglichkeiten zu erstellen, sowie Vor- und Nachnamen gesondert aufzulisten, aber auch vollständige Personennamen zu archivieren.

Diese Hyperonymierelationen können auch für die allgemeinen Menschenbezeichner definiert werden, um festzulegen, welche englischen Nomina auf einen Menschen referenzieren, um später eine Übergenerierung der Graphen zu verhindern. Dadurch schränken wir die grammatikalische Klasse der Nomina auf diesen Teil ein und verhindern somit, dass später beispielsweise Tiere statt Menschen im Text gefunden werden.

Da in biographischen Texten häufig Beschäftigungsverhältnisse beschrieben werden, sollte man auch nicht auf die Kategorie der Berufsbezeichnungen verzichten. In diesem Zusammenhang sind Organisationsnamen bzw. Firmennamen und organisationsspezifische Attribute nicht außer Acht zu lassen. Auch eine Liste an Branchen, Fachbereichen und Industriesektoren kann von Vorteil sein, wenn nur die Arbeitsdomäne einer Person genannt wird.

Des Weiteren kommen in diesen Kontexten häufig Ortsbestimmungen und Beschäftigungszeiträume vor, was wiederum den Einsatz eines Wörterbuches mit geographischen Bezeichnungen nötig macht. Aber auch für biographisch typische Sätze wie "Harry Clifford was born in Dallas, Texas in 1956." sind Toponyme ein unverzichtbares Muss.

All diese Vorüberlegungen zeigen, welche weiteren Entitäten außer Personenbezeichnungen für diese Aufgabe notwendig sein werden. Deshalb habe ich verschiedene Lexika mit Hilfe von Listen aus dem Internet erstellt, welche dem System eine riesige Wissensbasis liefern werden. Insgesamt haben alle Wörterbücher zusammen mehr als 10,6 Millionen Einträge.

#### 5.2.1 Lexikon der Vornamen

Das Lexikon der Vornamen enthält weit über 38500 Einträge und hat den Namen FirstNames-.dic. Die darin enthaltenen Vornamen stammen aus dem CISLEX, sowie aus verschiedenen Listen des Internets (Wikipedia [97] ist nur eine von vielen Quellen.). Es wurden einfache Vornamen wie Mary, aber auch komplexe Vornamen wie Mary-Anne aufgelistet.

In Abbildung 5.1 wird deutlich, in welcher Form die Vornamen gespeichert werden. An dieser Stelle soll für uns das Format der Kodierung nicht von Interesse sein, aber die Bedeutung der Symbole nach dem jeweiligen Vornamen möchte ich schon einmal hier erläutern. Es ist möglich bei der Erstellung der Lexika anzugeben, welchen semantischen oder grammatikalischen Kategorien die jeweiligen Einträge angehören. Außerdem kann man eigene Kategorien selbst definieren und die Wörter nach diesen klassifizieren. Für FirstNames-.dic habe ich 4 Kategorien verwendet, von denen die Klassen N, PR und Hum dem System Unitex [75] bekannte Kategorien für die englische Sprache sind und FN wurde von mir eingeführt. Die Tabelle 5.1 erläutert die Bedeutung dieser Abkürzungen.

Annelie,.N+PR+Hum+FN
Anneliese,.N+PR+Hum+FN
Anne-Marie,.N+PR+Hum+FN
Barbara-Anne,.N+PR+Hum+FN
Bárbara,.N+PR+Hum+FN
Calvin,.N+PR+Hum+FN
Charly,.N+PR+Hum+FN
Delores,.N+PR+Hum+FN
Elizabeth,.N+PR+Hum+FN
Elizabeth,.N+PR+Hum+FN
Elizaveta,.N+PR+Hum+FN
Francesca,.N+PR+Hum+FN
Francesca,.N+PR+Hum+FN
Fredérique,.N+PR+Hum+FN
Grace,.N+PR+Hum+FN

Abbildung 5.1: Auszug aus dem Lexikon FirstNames-.dic

Abkürzung	Kategorietyp	wörtliche Bedeutung	Erläuterung
N	grammatikalisch	Noun	Nomen
PR	semantisch	Proper Noun	Eigenname
Hum	semantisch	Human	Menschenbezeichner
FN	semantisch	First Name	Vorname

Tabelle 5.1: Abkürzungen, die in FirstNames-.dic verwendet werden.

#### 5.2.2 Lexikon der Nachnamen

Das Lexikon der Nachnamen enthält ca. 1,25 Millionen Einträge und trägt den Namen LastNames-.dic. Die darin enthaltenen Nachnamen stammen wie schon bei dem Lexikon der Vornamen aus dem CISLEX, sowie aus verschiedenen Listen aus dem Internet (Wikipedia [97] ist nur eine von vielen Quellen.).

Auch hier wurden einfache Nachnamen wie Smith, aber auch komplexe Nachnamen wie  $Ames\"{o}der$ -Gerogan gesammelt.

Abbildung 5.2 zeigt einige Einträge aus dem Wörterbuch der Nachnamen, wofür ich auch 4 Kategorien verwendet habe, von denen wieder die Klassen N, PR und Hum dem System Unitex [75] bekannte Kategorien für die englische Sprache sind und SN wurde von mir eingeführt. Die Tabelle 5.2 erläutert die Bedeutung dieser Abkürzungen.

Aabel, .N+PR+Hum+SN Aabenhus,.N+PR+Hum+SN Aabenraa, .N+PR+Hum+SN Abagameh, .N+PR+Hum+SN Abaganova, .N+PR+Hum+SN Amesöder-Gerogan, .N+PR+Hum+SN Gabanelli, .N+PR+Hum+SN Gabaniova, .N+PR+Hum+SN Gabanji, .N+PR+Hum+SN MacMillan, .N+PR+Hum+SN Macmillen, .N+PR+Hum+SN Olsson, .N+PR+Hum+SN Oltay, .N+PR+Hum+SN Palandt-Schäfer,.N+PR+Hum+SN Smith, .N+PR+Hum+SN Töchert-Yildiz, .N+PR+Hum+SN Tolwinski, . N+PR+Hum+SN Tolxdorff,.N+PR+Hum+SN

Abbildung 5.2: Auszug aus dem Lexikon LastNames-.dic

Abkürzung	Kategorietyp	wörtliche Bedeutung	Erläuterung
N	grammatikalisch	Noun	Nomen
PR	semantisch	Proper Noun	Eigenname
Hum	semantisch	Human	Menschenbezeichner
SN	semantisch	Surname	Nachname

**Tabelle 5.2:** Abkürzungen, die in *LastNames-.dic* verwendet werden.

#### 5.2.3 Lexika der Personennamen

Die Lexika der Personennamen enthalten zusammen ungefähr 8,3 Millionen Einträge. Es sind insgesamt 9 Lexika, welche die Namen

- LongNamesBios-.dic,
- LongNamesFT-.dic,
- LongNamesSpecialistInfo-.dic,
- LongNamesZoominfo1-.dic,
- LongNamesZoominfo2-.dic,
- LongNamesZoominfo3-.dic,
- LongNamesZoominfo4-.dic,
- LongNamesZoominfo5-.dic und
- LongNamesAuthors-.dic

tragen. Einerseits sind die Namen nach ihrer Herkunft aus dem Internet gruppiert, und andererseits werden diese Lexika später vom System Unitex [75] binär für die Verarbeitung gespeichert, und diese Binärdateien sind auf eine Maximalgröße von 16 MB pro Datei beschränkt. Deshalb wurde nicht ein riesiges Lexikon für alle Personennamen angelegt, sondern 9 Stück, welche alle die gleichen grammatikalischen und semantischen Zusatzinformationen enthalten.

Wie in Tabelle 5.3 ersichtlich wird, wurden auch hier für die Kodierung der Wörterbucheinträge 4 Kategorien verwendet, von denen wieder die Klassen N, PR und Hum dem System Unitex [75] bekannte Kategorien für die englische Sprache sind und LN wurde von mir eingeführt. Leider konnte ich nicht die Abkürzung FN für Full Name wählen, da FN schon als Kategorie für First Name bzw. Vornamen im FirstNames-.dic belegt war, musste ich auf LN für Long Name ausweichen.

Abkürzung	Kategorietyp	wörtliche Bedeutung	Erläuterung
N	grammatikalisch	Noun	Nomen
PR	semantisch	Proper Noun	Eigenname
Hum	semantisch	Human	Menschenbezeichner
LN	semantisch	Long Name (Full Name)	Vollständiger Personen-
			name, bestehend aus
			Vor- und Nachname,
			evtl. mit Titel

Tabelle 5.3: Abkürzungen, die in LongNames\*-.dic verwendet werden.

#### 5.2.3.1 Personennamen von Biography.com

Das Internetportal Biography.com [4] stellt knapp 25000 englischsprachige Biographien zur Verfügung. Jede einzelne Biographie enthält in der Regel einen Personennamen. Somit bietet Biography.com indirekt fast 25000 Personennamen, bestehend aus Vorund Nachnamen im Netz an.

Variiert man nun die Stellung von Vor- und Nachname, so dass

einmal Aaron, Hank, .N+PR+Hum+LN und einmal Hank Aaron, .N+PR+Hum+LN

ins Wörterbuch aufgenommen werden, so verdoppeln sich die Einträge.

Auch mit der Abkürzung zweiter Vornamen, lassen sich weitere Lexikoneinträge gewinnen, wie z.B.

Aaron T\. Beck, .N+PR+Hum+LN, was Aaron T(emkin) Beck, .N+PR+Hum+LN

im Original lautete.

Des Weiteren können selten benutzte zweite Vornamen auch weggelassen oder geklammert werden, was zwei neuen Einträge entstehen lässt:

Aage Bohr, .N+PR+Hum+LN, sowie Aage (Niels) Bohr, .N+PR+Hum+LN und Aage Niels Bohr, .N+PR+Hum+LN

Mit Hilfe dieser Methoden lassen sich aus knapp 25000 Personennamen, mehr als 69000 Wörterbucheinträge generieren und im Lexikon *LongNamesBios-.dic* speichern.

Aaron Roy Weintraub, .N+PR+Hum+LN

Aaron Siskind,.N+PR+Hum+LN

Aaron Spelling,.N+PR+Hum+LN

Agnes Macphail,.N+PR+Hum+LN

Agnes Maud Royden, . N+PR+Hum+LN

Agnes Nestor, . N+PR+Hum+LN

Anton van Duinkerken, .N+PR+Hum+LN

Antón Villar Ponte, .N+PR+Hum+LN

Augusto César Sandino, .N+PR+Hum+LN

Augusto Pinochet, . N+PR+Hum+LN

. . .

Abbildung 5.3: Auszug aus dem Lexikon LongNamesBios-.dic

#### 5.2.3.2 Personennamen aus der Financial Times

Natürlich bietet das aus der Financial Times $^{35}$  zusammengestellte Korpus eine Fülle an Personennamen.

Um auf diese Namen nicht verzichten zu müssen, wurde zunächst nach zwei nebeneinander groß geschriebenen Wörtern, in deren Mitte evtl. noch ein einzelner Großbuchstabe auftreten darf, in dem gesamten Text gesucht. Diese Suche würde u.a.

George W Bush oder Henna Nordqvist

als Treffer in den Zeitungsberichten ergeben.

Da in englischen Texten in der Regel nur Eigennamen groß geschrieben werden, wenn der Satzanfang außer Acht gelassen wird, kann man mit großer Wahrscheinlichkeit davon ausgehen, dass alle Eigennamen im Korpus auf diese Weise gefunden werden.

Diese Eigennamen sind sicher nicht alle nur Personennamen, doch mehr als die Hälfte von ihnen fallen in diese Kategorie. Um Organisationsnamen wie *Peerless Industries* oder *Goodwill Industries International* aussortieren zu können, musste die Liste der potentiellen Personen von mir manuell durchgesehen werden. Aufgrunddessen konnten Firmennamen oder Filmtitel aus diesem Lexikon entfernt werden, welche zwar auf das Suchmuster gepasst haben, aber nicht als Treffer beabsichtigt waren.

Insgesamt blieben ca. 630000 richtige Personennamen von ursprünglich 701778 Kandidaten übrig, welche im Wörterbuch mit dem Namen LongNamesFT-.dic gespeichert wurden.

Dieses Ergebnis lässt nicht darauf schließen, dass knapp 7% von dieser Eigennamenliste Organisationsnamen gewesen sein müssen. Einerseits wären das viel zu wenig Firmennamen für eine Jahresausgabe der Financial Times, und andererseits wurden vor dem manuellen Durchsehen, alle Organisationsnamen, welche durch ihre Rechtsform (z.B. GmbH, AG, S.A., PLC, Inc., etc.) gekennzeichnet waren, aus dieser Kandidatenliste herausgenommen und gesondert aufbewahrt.

#### 5.2.3.3 Personennamen von SpecialistInfo.com

SpecialistInfo.com<sup>36</sup> bietet auf seiner Seite u.a. eine Liste von Personen, die als Berater tätig sind (Consultants), welche derzeit 31000 Personennamen mit ihrem jeweiligen Titel (akademischer Grad, militärischer Rang usw.) oder der jeweiligen Anredeform (Mr oder Mrs) enthält.

Hier hätte man auch die Möglichkeit gehabt, durch Entfernen der Titel und Anredeformen, welche im Grunde nur den Beruf bzw. das Geschlecht der Person kodieren, die Anzahl der späteren Lexikoneinträge zu verdoppeln. Doch da die meisten dieser Namen ohne Titel oder Anrede, in bereits genannten oder später genannten Lexika namentlich vorkommen, entschied ich mich dazu, auf diesen Schritt zu verzichten. Somit hat das Wörterbuch LongNamesSpecialistInfo-.dic nahezu 31000 Einträge in der Form, wie sie in Abbildung 5.5 gezeigt werden.

 $<sup>^{35} {\</sup>rm http://news.ft.com/home/us}$ 

 $<sup>^{36} \</sup>mathtt{http://www.specialistinfo.com/directory.php}$ 

George W Bush,.N+PR+Hum+LN
Helen A Donnelly,.N+PR+Hum+LN
Henna Nordqvist,.N+PR+Hum+LN
Iain Allan,.N+PR+Hum+LN
Iain A Macdonald,.N+PR+Hum+LN
Pablo Amorsolo,.N+PR+Hum+LN
Pablo Andrade,.N+PR+Hum+LN
Zhang Fan,.N+PR+Hum+LN
Zhang Feng,.N+PR+Hum+LN
Zoran Vucevic,.N+PR+Hum+LN

Abbildung 5.4: Auszug aus dem Lexikon LongNamesFT-.dic

Assoc Prof Neil R McLean,.N+PR+Hum+LN
Col Simon Mellor,.N+PR+Hum+LN
Cr Charmian M Kalic,.N+PR+Hum+LN
Dr A Alex Freeman,.N+PR+Hum+LN
Mr Adam Booth,.N+PR+Hum+LN
Mrs Amanda J Smith,.N+PR+Hum+LN
Wing Comm Andrew J Gibbons,.N+PR+Hum+LN
...

Abbildung 5.5: Auszug aus dem Lexikon LongNamesSpecialistInfo-.dic

### 5.2.3.4 Personennamen von ZoomInfo.com

ZoomInfo.com [100] ist eine Suchmaschine für Personen- und Firmennamen. Im Gegensatz zu den Firmennamen bieten sie für Personennamen eine Auflistung von A-Z an. Da sie über einen riesigen Datenbestand verfügt, lassen sich relativ leicht ca. 3,5 Millionen Personennamen herunterladen. Nach Anwendung der vorher genannten Methode des Vertauschens von Vor- und Nachnamen verdoppelt sich die Anzahl auf 7 Millionen potentielle Lexikoneinträge.

Wie bereits anfangs erwähnt, kann man aufgrund der Systembeschränkungen in Unitex [75] kein Wörterbuch mit über 7 Millionen Einträgen erstellen. Deshalb wurden die insgesamt 7,1 Millionen Einträge auf 5 Lexika verteilt, welche wie folgt benannt sind.

- LongNamesZoominfo1-.dic
- LongNamesZoominfo4-.dic
- LongNamesZoominfo2-.dic
- $\bullet \ \ LongNamesZoominfo 5 \text{-.} dic$
- LongNamesZoominfo3-.dic

Aacker\, David, N+PR+Hum+LN
A. Ackerman, N+PR+Hum+LN
Aaron Akwaboah, N+PR+Hum+LN
Douglas McCoy, N+PR+Hum+LN
Douglas McCracken, N+PR+Hum+LN
K. Bales-Blackburn, N+PR+Hum+LN
Keck\, Elizabeth, N+PR+Hum+LN
Keaton\, Penny, N+PR+Hum+LN
Pogue-Adm\, Debbie, N+PR+Hum+LN
W.O. Gray, N+PR+Hum+LN
Wolff\, Art, N+PR+Hum+LN
Wolff\, Arthur, N+PR+Hum+LN
Wolff\, Ashley, N+PR+Hum+LN

Abbildung 5.6: Auszug aus dem Lexikon LongNamesZoominfo[12345]-.dic

#### 5.2.3.5 Weitere Personennamen

Wie bereits Latanya Sweeney in ihrem Artikel "Finding Lists of People on the Web" [87] so treffend bemerkt hatte, kann man auch durch Zufall auf riesige Listen mit Personennamen im Internet stoßen. Als ich auf Google.fr nach einer Kurzbiographie von Nathalie Friburger suchte, bot mir Google u.a. zwei Listen<sup>37</sup> von Autoren<sup>38</sup> an. So erhielt ich für das Lexikon LongNamesAuthors-.dic insgesamt 449500 Einträge.

A.-A. A. Jucys,.N+PR+Hum+LN
Aabhas Paliwal,.N+PR+Hum+LN
Aabhas V. Paliwal,.N+PR+Hum+LN
A. Almansa-Martin,.N+PR+Hum+LN
Dae-Ghon Kho,.N+PR+Hum+LN
Dagmar Schönfeld,.N+PR+Hum+LN
Patrice Brémond-Grégoire,.N+PR+Hum+LN
Rabi N. Mahapatra,.N+PR+Hum+LN
Rachel Mason-Jones,.N+PR+Hum+L
Waleed A. Youssef,.N+PR+Hum+LN

Abbildung 5.7: Auszug aus dem Lexikon LongNamesAuthors-.dic

<sup>&</sup>lt;sup>37</sup>http://136.199.54.185/~ley/db/indices/AUTHORS

<sup>38</sup>http://sunsite.online.globule.org/dblp/db/indices/AUTHORS

#### 5.2.4 Lexika der Personentitel

Für den Fall dass ein Personenname im Text auftaucht, welcher nicht in einem der Wörterbücher registriert ist, muss man den Kontext zur Erkennung des Namens miteinbeziehen. Ein eindeutiger Indikator für einen Personennamen ist beispielsweise ein akademischer Grad, ein militärischer Rang, ein aristokratischer Titel oder einfach nur eine Form der Anrede. Aus diesem Grund habe ich mit der Hilfe der englischen Wikipedia [97] ein Wörterbuch von verschiedenen Titelbezeichnungen zusammengestellt. Dieses trägt den Namen *Titles-.dic* und enthält über 370 Einträge.

Abbildung 5.8 zeigt eine Auswahl der Wörterbucheinträge aus diesem Lexikon, wofür ich 4 Kategorien verwendet habe, von denen die Klassen N und XN dem System Unitex [75] bekannte Kategorien für die englische Sprache sind und Abbrev, sowie Title wurden von mir eingeführt. Die Tabelle 5.4 erläutert die Bedeutung dieser Abkürzungen.

Duke,.N+Title
King,.N+Title
Mrs\.,.Abbrev+Title
Mrs,.Abbrev+Title
MSci,.Abbrev+Title
MSci,.Abbrev+Title
Queen,.N+Title
Rabbi,.N+Title
Sir,.N+Title
Sister,.N+Title
Sultan,.N+Title
Sultan,.N+Title
ThD,.Abbrev+Title
The Right Honourable,.N+XN+Title
Tsar,.N+Title

Abbildung 5.8: Auszug aus dem Lexikon Titles-.dic

Abkürzung	Kategorietyp	wörtliche Bedeutung	Erläuterung
N	grammatikalisch	Noun	Nomen
XN	grammatikalisch	Extended Noun	Mehrwortlexem
Abbrev	semantisch	Abbreviation	Abkürzung
Title	semantisch	Title	Titel, wie z.B. akademische
			Grade, aristokratische Titel

**Tabelle 5.4:** Abkürzungen, die in *Titles-.dic* verwendet werden.

#### 5.2.4.1 Akademische Grade

DUniv,.Abbrev+Title
MSci,.Abbrev+Title
Prof\. h\.c\.,.Abbrev+Title
Prof hc,.Abbrev+Title
Prof\.,.Abbrev+Title
Psy\.D\.,.Abbrev+Title
PsyD,.Abbrev+Title

Titles-.dic enthält eine Fülle an Abkürzungen für akademische Grade. Meist besteht im englischen Sprachgebrauch die Möglichkeit diese Titel mit nachfolgendem Punkt oder ohne diesen zu schreiben. Um nun alle Möglichkeiten, wie eine Titelbezeichnung im Text auftreten kann, abzudecken, wurden natürlich auch diese Variationen ins Lexikon mitaufgenommen.

#### 5.2.4.2 Aristokratische Titel

Duke,.N+Title
Earl,.N+Title
Lord,.N+Title
Queen,.N+Title
Prince,.N+Title
Princess,.N+Title

Oft wird auch in den Nachrichten über adelige Persönlichkeiten berichtet. Deren Namen wird ein Adelstitel vorangestellt. Darum ist es ebenfalls wichtig, die Gruppe dieser Titel in das Wörterbuch aufzunehmen. So können später auch unbekannte Adelige, die irgendeiner Nebenlinie entstammen, und deshalb in keinem der Personenlexika aufgeführt sind, im Text gefunden werden. Auch werden bei Adelstiteln im Gegensatz zu akademischen Graden in der Regel keine Abkürzungen verwendet.

#### 5.2.4.3 Weitere Anredeformen

Captain,.N+Title
Mr,.Abbrev+Title
Patriarch,.N+Title
People's Commissar,.N+XN+Title
Pope,.N+Title
Saint,.N+Title

Es gibt weitere zahlreiche Anredemöglichkeiten oder Titel für die unterschiedlichsten Personengruppen, seien es jetzt Vertreter der Kirche, des Militärs oder einfach nur die Anredeform Mr, Mrs oder Ms. Da keine dieser Personengruppen dominanter als die andere in diesem Wörterbuch ist, werden sie in diesem Abschnitt zusammen erwähnt.

## 5.2.5 Lexika der allgemeinen Menschenbezeichner

Bis zu diesem Punkt wurden die Lexika der Personeneigennamen beschrieben, die sicherlich auch zu den Menschenbezeichnern zählen. Doch in diesem Abschnitt soll es nun ausschließlich um Bezeichnungen für Menschen gehen, die keine Eigennamen sind. Unter allgemeinen Menschenbezeichnern verstehe ich beispielsweise Begriffe, die Verwandtschaftsverhältnisse wie "mother" ausdrücken, Berufsbezeichnungen wie "salesperson" oder Personen durch ihren Charakter wie "idealist" bzw. ihre Neigungen wie "lesbian" näher bestimmen. Des Weiteren fallen auch Wörter, welche eine Staatszugehörigkeit oder die Zugehörigkeit zu einem Bezirk, einer Provinz, einer Stadt oder eines Stadtteils ausdrücken, in diese Kategorie der allgemeinen Menschenbezeichner. Aus diversen Internetverzeichnissen habe ich insgesamt über 54000 Begriffe zusammengetragen, die auf diese Beschreibungen zutreffen. Aus einer früheren Arbeit von Friederike Mallchok [55] konnte ich auf ein Lexikon mit mehr als 99000 Einträgen für Berufsbezeichnungen zurückgreifen, so dass mir ein Inventar von insgesamt ca. 153000 allgemeinen Menschenbezeichnern zur Verfügung stand.

#### 5.2.5.1 Allgemeine Menschenbezeichnungen aus WordNet

WordNet [98] ist eine echte Bereicherung, wenn man Begriffe mit Hilfe semantischer Relationen wie der Hyperonymie sucht.

WordNet [24] bezeichnet sich selbst als elektronische lexikalische Datenbasis. Im Grunde teilt WordNet das Lexikon in fünf Kategorien auf: Nomen, Verben, Adjektive, Adverbien und Funktionswörter [63]. Dabei verfolgt WordNet die Zielsetzung lexikalische Information mehr durch Wortbedeutungen zu organisieren als durch Wortformen. Die Lexikoneinträge für Nomen sind nicht alphabetisch sondern hierarchisch geordnet.

#### { Nomen-Oberbegriff (semantische Merkmale) }

Dabei werden die Nomen in Synsets klassifiziert, wobei ein Synset die Menge alle Synonyme eines Nomens ist. Es wurden alle wichtigen semantischen Relationen wie Synonymie, Meronymie, Holonymie, Hyperonymie, Hyponymie und Antonymie kodiert, und man könnte die semantische Relationen von WordNet als Relationen zwischen den einzelnen Synsets beschreiben. Außerdem ist die große Masse von Nomen hierarchisch in 25 eindeutigen initialen Synsets (25 unique beginners) organisiert (nach [28]).

Da Menschenbezeichner in der Regel als Nomen vorkommen, möchte ich nicht näher auf die anderen grammatikalischen Kategorien eingehen, welche in WordNet noch vertreten sind. Außerdem ist für uns momentan nur das Anfangssynset {person, human being} von Bedeutung, in dem alle Menschenbezeichner in WordNet gespeichert sind.

Die Gruppe der Menschenbezeichner aus WordNet enthält alle oben genannten Variationen an allgemeinen Menschenbezeichnungen und darüber hinaus noch ganz allgemeine Berufsbezeichnungen wie "engineer" oder "worker".

Abbildung 5.9 zeigt einen Ausschnitt aus dem Begriffsinventar von WordNet und die Tabelle 5.5 gibt einen Überblick zur grammatikalischen und semantischen Information des neu erstellten Wörterbuchs *MenbezWordnet-.dic*, was insgesamt ungefähr 6400 Einträge aufweisen kann.

godmother,.N+Hum
godparent,.N+Hum
godson,.N+Hum
heterosexual person,.N+XN+Hum
heterosexual,.N+Hum
husband,.N+Hum
idealist,.N+Hum
identical twin,.N+XN+Hum
islander,.N+Hum
lesbian,.N+Hum
madman,.N+Hum
mon,.N+Hum
woman,.N+Hum

Abbildung 5.9: Auszug aus dem Lexikon MenbezWordnet-.dic

Abkürzung	Kategorietyp	wörtliche Bedeutung	Erläuterung
N	grammatikalisch	Noun	Nomen
Hum	semantisch	Human	Menschenbezeichner

Tabelle 5.5: Abkürzungen, die in MenbezWordnet-.dic verwendet werden.

#### 5.2.5.2 Berufsbezeichnungen

Das Lexikon der Berufsbezeichnungen ist aus einer Vielzahl von Internetverzeichnissen<sup>39</sup>, den Listen von Job-Suchmaschinen<sup>40</sup> oder von Arbeitsämtern<sup>41</sup> aus dem englischsprachigen Raum entstanden.

Auch wurde ich noch bei der englischen Wikipedia [97] auf meiner Suche nach Berufsbezeichnern fündig und stellte so ein Wörterbuch von verschiedenen Berufsbezeichnungen zusammen. Dieses trägt den Namen *JD-.dic* und enthält nahezu 47000 Einträge.

In Abbildung 5.10 wird ein Ausschnitt aus diesem Lexikon gezeigt, wofür ich 4 Kategorien verwendet habe, von denen die Klassen N, XN und Hum dem System Unitex [75] bekannte Kategorien für die englische Sprache sind und JD wurde von mir eingeführt. Die Tabelle 5.6 erläutert die Bedeutung dieser Abkürzungen.

<sup>&</sup>lt;sup>39</sup>Guide to the World of Occupations [48], Occupational Outlook Handbook [70] [71], List of occupations [51], Standard Occupational Classification (SOC) [93], Dictionary Of Occupational Titles [19]

<sup>&</sup>lt;sup>40</sup>CareerBuilder.com [8], LabourMarket [50], Prospects.ac.uk [76]

<sup>&</sup>lt;sup>41</sup>Canadian Job Classification [30], Division of Professional Licensure [18], Ministry of Manpower [65], U.S. Department of Labor [92] [95]

accountant clerk,.N+XN+Hum+JD
accountant,.N+Hum+JD
adjuster and inspector,.N+XN+Hum+JD
adjuster & inspector,.N+XN+Hum+JD
aeronautical technician,.N+XN+Hum+JD
bakery products checker,.N+XN+Hum+JD
bank president,.N+XN+Hum+JD
bicycle salesperson,.N+XN+Hum+JD
blacksmith,.N+Hum+JD
cargo serviceman,.N+XN+Hum+JD
car salesman,.N+XN+Hum+JD
car saleswoman,.N+XN+Hum+JD

Abbildung 5.10: Auszug aus dem Lexikon JD-.dic

Abkürzung	Kategorietyp	wörtliche Bedeutung	Erläuterung
N	grammatikalisch	Noun	Nomen
XN	grammatikalisch	Extended Noun	Mehrwortlexem
Hum	semantisch	Human	Menschenbezeichner
JD	semantisch	Job Descriptor	Berufsbezeichner

**Tabelle 5.6:** Abkürzungen, die in *JD-.dic* verwendet werden.

#### 5.2.5.3 Einwohnerbezeichnungen

Eine relativ große Gruppe innerhalb der Menschenbezeichner, sind die Bezeichnungen nach Herkunftsland oder -ort. Bei der Zusammenstellung dieser Liste habe ich mich ganz darauf konzentriert möglichst alle Menschenbezeichner, die sich auf Staatsbürgerschaften beziehen, in das neue Lexikon aufzunehmen. Bei den Begriffen, welche ausdrücken, dass jemand Einwohner einer Stadt, eines Ortes, eines Bezirks, einer Provinz, eines Bundesstaates oder Bundeslandes ist, habe ich mich auf die wichtigsten und bekanntesten Regionen und Städte beschränkt, da Bezeichnungen wie "New Yorker" im Korpus häufiger als z.B. "Nottinghamian" auftreten.

Auch hier wurde ich bei der englischen Wikipedia [97] fündig und stellte so ein Wörterbuch zusammen, was den Namen Citizens-.dic hat und über 600 Einträge enthält.

Abbildung 5.11 zeigt einen Ausschnitt aus diesem Lexikon, wofür ich 8 Kategorien verwendet habe, von denen die Klassen N und Hum dem System Unitex [75] bekannte Kategorien sind und Citizen, AuProvinceCitizen, CaProvinceCitizen, USstateCitizen, NYCcitizen und Urbanite wurden von mir eingeführt. Die Tabelle 5.7 erläutert die Bedeutung dieser Abkürzungen.

Alaskans, Alaskan.N+Citizen+USstateCitizen+Hum
Albertans, Albertan.N+Citizen+CaProvinceCitizen+Hum
Americans, American.N+Citizen+Hum
Athenians, Athenian.N+Citizen+Urbanite+Hum
Atlantans, Atlantan.N+Citizen+Urbanite+Hum
Aucklanders, Aucklander.N+Citizen+Urbanite+Hum
Brooklyners, Brooklyner.N+Citizen+Urbanite+NYCcitizen+Hum
Bavarians, Bavarian.N+Citizen+Hum
Canberrans, Canberran.N+Citizen+AuProvinceCitizen+Hum
Christmas Islanders, Christmas Islander.N+Citizen+Hum
Ontarians, Ontarian.N+Citizen+CaProvinceCitizen+Hum
Queensites, Queensite.N+Citizen+Urbanite+NYCcitizen+Hum
Tasmanians, Tasmanian.N+Citizen+AuProvinceCitizen+Hum

Abbildung 5.11: Auszug aus dem Lexikon Citizens-.dic

Abkürzung	Kategorietyp	wörtliche Bedeutung	Erläuterung
N	grammatikalisch	Noun	Nomen
Hum	semantisch	Human	Menschenbezeichner
Citizen	semantisch	Citizen	(Staats)Bürger
AuProvinceCitizen	semantisch	Australian Province	Bewohner einer au-
		Citizen	stralischen Provinz
CaProvinceCitizen	semantisch	Canadian Province	Bewohner einer ka-
		Citizen	nadischen Provinz
USstateCitizen	semantisch	US State Citizen	Bewohner eines US
			Bundesstaates
NYCcitizen	semantisch	New York City Citizen	Bewohner von New
			York City
Urbanite	semantisch	Urbanite	Stadtbewohner

Tabelle 5.7: Abkürzungen, die in Citizens-.dic verwendet werden.

## 5.2.6 Lexikon der personenbezogenen Prädikate

In Abschnitt 1.2 wurde bereits der Begriff der biographischen Relation erläutert. Allen drei Kategorien von Relationen - der persönlichen, der öffentlichen und der zufälligen Relation - lassen sich Verben bzw. Verbalphrasen der englischen Sprache zuordnen. Obwohl anfangs die Priorisierung von öffentlichen Relationen in den Vordergrund gestellt wurde, möchte ich für die Wörterbucherstellung keine Differenzierung zwischen den drei Verbkategorien, welche sich durchaus überschneiden können, vornehmen.

Für die Erkennung von Menschenbezeichnern in biographischen Kontexten ist die Einschränkung der potentiellen Verben, welche in der Umgebung von Personen vorkommen, fast unumgänglich. Deshalb müssen Prädikate (insbesondere Verben) ausgewählt werden, welche zumindest an erster Argumentposition vornehmlich einen Menschenbezeichner fordern. Diese personenbezogenen Prädikate können natürlich auch als zweites Argument eine Person verlangen, doch reicht die erste Einschränkung in der Regel aus.

Um nun ein solches Lexikon erstellen zu können, hat man die Möglichkeit, alle relevanten Verben der englischen Sprache aufzulisten, oder die vorhandenen Ressourcen werden mit in die Erstellung eines personenbezogenen Verbalphrasenlexikons einbezogen.

Für die weitere Vorgehensweise wurden alle Personennamen des Financial Times Korpuses annotiert und die Menge der Sätze, die einen Namen enthalten, extrahiert. Danach wurden die lokalen Grammatiken zur Lemmatisierung komplexer englischer Verben (siehe Abschnitt 3.1.3, vgl. Maurice Gross, 1998-1999 [44]) auf den neuen Text angewandt und die entsprechenden Treffer markiert. Im Anschluss wurden alle erkannten Verbalphrasen, denen eine Nominalphrase mit einem Menschenbezeichner vorangestellt war, dem Korpus entnommen und noch einmal manuell überprüft, bevor man sie dem neuen Lexion HumVP-.dic hinzugefügt hat. Die noch fehlenden einfachen Verben wurden aus den zusammengesetzten Verbeinträgen rekonstruiert.

accused of committing, .HumVP aims to recruit, .HumVP agreed to marry, .HumVP announced, .HumVP apologize for misleading, .HumVP are resigning, .HumVP arrested for shoplifting, .HumVP attempted to establish, .HumVP became engaged to crown, .HumVP can join, .HumVP cannot predict, .HumVP cannot risk losing, .HumVP could be fired, .HumV

Abbildung 5.12: Auszug aus dem Lexikon Hum VP-.dic

. . .

#### 5.2.7 Lexika der Branchen

Auch Branchen treten recht häufig im Umfeld von Menschenbezeichnern bzw. Berufsbezeichnern auf. Meist dienen sie als Ergänzung oder Spezifikation eines Arbeitsbereichs. Diese Branchenbezeichnungen kann man aufgrund ihres syntaktischen Verhaltens in zwei Bereiche aufteilen: Fachbereiche und Lehrfächer, sowie Sektoren- und Branchenbezeichnungen.

#### 5.2.7.1 Fachbereiche und Lehrfächer

Die Fachbereiche und Lehrfächer zeichnet besonders aus, dass sie in der Regel im rechten Kontext von Berufsbezeichnern zu finden sind, welche akademische Berufe, Lehrtätigkeiten oder Spezialisierung von Berufen ausdrücken. Des Weiteren fordern die jeweiligen Fachbereiche selbst keine Ergänzungen für sich selbst, was nichts anderes heißt, als dass beispielsweise american literature in der Nominalphrase "professor of american literature" kein eigenes Argument mehr braucht - es dient selbst als Ergänzung zu professor.

Die gesammelten Wörterbucheinträge ließen sich aus Seiten der englischen Wikipedia [97] extrahieren und ergaben ein Begriffsinventar von ca. 580 Einträgen für das Lexikon Disciplines-. dic.

Abbildung 5.13 zeigt einen Ausschnitt aus diesem Lexikon, wofür ich 4 Kategorien verwendet habe, von denen die Klassen N und XN dem System Unitex [75] bekannte Kategorien sind und Sector, sowie Discipline wurden von mir eingeführt. Die Tabelle 5.8 erläutert die Bedeutung dieser Abkürzungen.

american literature,.N+XN+Sector+Discipline analytical chemistry,.N+XN+Sector+Discipline anatomy,.N+Sector+Discipline ancient history,.N+XN+Sector+Discipline bioinformatics,.N+Sector+Discipline cognitive science,.N+XN+Sector+Discipline dentistry,.N+Sector+Discipline descriptive linguistics,.N+XN+Sector+Discipline ...

Abbildung 5.13: Auszug aus dem Lexikon Disciplines-.dic

Abkürzung	Kategorietyp	wörtliche Bedeutung	Erläuterung
N	grammatikalisch	Noun	Nomen
XN	grammatikalisch	Extended Noun	Mehrwortlexem
Sector	semantisch	Sector	Sektor, Branche
Discipline	semantisch	Discipline	Fachbereich, Lehrbereich

Tabelle 5.8: Abkürzungen, die in *Disciplines-.dic* verwendet werden.

#### 5.2.7.2 Sektoren- und Branchenbezeichnungen

Ganz anders verhalten sich dagegen die Sektoren- und Branchenbezeichnungen. Sie fordern sehr wohl eine Ergänzung in Form von Begriffen wie "industry", "sector" oder "company" und bestimmen somit genauer, in welcher Branche eine Person beschäftigt ist. Somit sind im rechten Kontext dieser Branchenbegriffe, die unterschiedlichsten Schlüsselbegriffe zu finden.

administration sector
animal food industry
arts and leisure sector

Die eben genannten Sektoren könnten allerdings auch mit dem Schlüsselbegriff "services" versehen werden. Daraus wird ersichtlich, dass es eine Gruppe von Branchenbezeichnungen gibt, die flexibler als andere in der Wahl ihrer Ergänzungen sind. Es gibt eine Reihe von Begriffen, welche nur mit der Ergänzung "services" im Text auftreten. Es wäre beispielsweise unsinnig "animal physiotherapy services" ein anderes Argument als "services" mitzugeben. Solche Branchenbegriffe, die vornehmlich oder ausschließlich mit der Ergänzung "services" in Texten vorkommen, wurde die Ergänzung bei der Wörterbucherstellung mitgegeben.

In Abbildung 5.13 ist ein Ausschnitt aus dem Lexikon Sector-.dic zu sehen, welches mehr als 8400 Einträge vorzuweisen hat, und wofür ich 3 Kategorien verwendet habe, von denen die Klassen N und XN dem System Unitex [75] bekannte Kategorien sind und Sector wurde von mir eingeführt. Die Tabelle 5.9 erläutert die Bedeutung dieser Abkürzungen.

accident and health insurance,.N+XN+Sector administration,.N+Sector aircraft,.N+Sector aircraft part and auxiliary equipment,.N+XN+Sector animal food,.N+XN+Sector arts and leisure,.N+XN+Sector child services,child service.N+XN+Sector

Abbildung 5.14: Auszug aus dem Lexikon Sector-.dic

Abkürzung	Kategorietyp	wörtliche Bedeutung	Erläuterung
N	grammatikalisch	Noun	Nomen
XN	grammatikalisch	Extended Noun	Mehrwortlexem
Sector	semantisch	Sector	Sektor, Branche

Tabelle 5.9: Abkürzungen, die in Sector-. die verwendet werden.

## 5.2.8 Lexika der Organisationsnamen

Organisationsnamen sind Entitäten, welche ähnlich wie Branchenbezeichnungen im Kontext eines Arbeitsverhältnisses auftreten können. Wenn es darum geht, die Zugehörigkeit zu einem Betrieb, einer Firma oder einem Konzern auszudrücken, kann man auf Organisationsbezeichner nicht verzichten. Da häufig Personennamen im Wortlaut von Firmennamen vorkommen können, ist es von Vorteil ein Organisationsnamenlexikon aufzubauen. Somit kann später ausgeschlossen werden, dass z.B. bei "John Miller Systems" Systems als zweiter Nachname erkannt wird. Außer einem Eintrag im Organistionsnamenlexikon (ONL) gibt es noch weitere Faktoren, welche sicherstellen können, dass Organisationen auch als solche identifiziert werden. So kann beispielsweise die Nennung der Rechtsform nach dem Firmennamen oder eine kurze Beschreibung des Unternehmens beim Finden eines Organisationsnamen hilfreich sein.

### 5.2.8.1 Allgemeine Organisationsbezeichnungen

Wie bereits Friederike Mallchok im Zuge ihrer Dissertation [55] herausgefunden hatte, werden Organisationsnamen in Wirtschaftstexten meist durch eine kurze Beschreibung eingeführt. Sie fertigte für diesen Zweck ein Lexikon für Organisationsbezeichnungen an, welches insgesamt 23800 Einträge enthält.

Abbildung 5.15 gibt einen Einblick in dieses Organisationsbeschreibungslexikon, welches kurz ODL von ihr genannt wurde und ab jetzt im System unter dem Namen *orgbez-.dic* verwendet wird.

```
accessories retail company,.orgbez accident insurance company,.orgbez accountacy firm,.orgbez accountancy and consultancy firm,.orgbez agro-industry conglomerate,.orgbez airplane maker,.orgbez automobile seat manufacturer,.orgbez bankholding company,.orgbez bath-solution maker,.orgbez bicycle company,.orgbez
```

Abbildung 5.15: Auszug aus dem Lexikon orgbez-.dic

#### 5.2.8.2 Eigennamen von Organisationen

Des Weiteren hat Friederike Mallchok ein Lexikon für Organisationsnamen (ONL) mit fast 286500 Unternehmen angelegt, welches hier ebenfalls berücksichtigt werden soll und den Namen org-.dic trägt.

Intel Corp,.org
Lions Gate Films,.org
Kids Unlimited Ltd,.org
Lactalis Iberia S\.A,.org
Caterpillar Group,.org
Deutsche Telekom,.org
Exxon Mobile,.org
Walt Disney Pictures,.org

Abbildung 5.16: Auszug aus dem Lexikon org-.dic

Mit Hilfe diverser Listen aus dem Internet konnte das Organisationsnamenlexikon von Friederike Mallchok um einige Einträge erweitert werden. Auch die in Abschnitt 5.2.3.2 beschriebene Methode der Eigennamenerkennung, welche auf die Jahresausgabe der Financial Times 2004 angewandt wurde, lieferte weitere potentielle Firmennamen. Viele der damals aussortierten Begriffe waren Bezeichnungen für Unternehmen, welche an dieser Stelle von Nutzen sein können.

Firmen, die ihrem Namen die entsprechende Rechtsform anfügen, wurden allerdings doppelt in das neue Wörterbuch *Companies-.dic* aufgenommen - einmal mit dem Kürzel für die Rechtsform und einmal ohne dieses. Dadurch kann das jeweilige Unternehmen in beiden Varianten im Text auftreten und wird jedes Mal gefunden.

Um die neu gewonnenen Lexikoneinträge in das System aufnehmen zu können, wurde ein gesondertes Wörterbuch mit mehr als 229600 Einträge angelegt, bei dem die Einträge aus dem ONL von Friederike Mallchok nicht integriert wurden. Aus Organisationsgründen sollten die Inhalte beider Lexika nicht miteinander vermischt werden, was aber dem System keine Einbußen bringt, da ihm insgesamt 515500 Organisationsnamen zur Verfügung stehen.

Wie Abbildung 5.17 deutlich macht, wurden den Einträgen im Lexikon *Companies-.dic* grammatikalische sowohl als auch semantische Zusatzinformationen mitgegeben, welche in Tabelle 5.10 erklärt werden.

Abit AG,.N+XN+PR+Company
Abitex Resources Inc,.N+XN+PR+Company
Above Chase Nominees Ltd,.N+XN+PR+Company
Above Chase Nominees,.N+XN+PR+Company
A&C Black plc,.N+XN+PR+Company
Adler GmbH,.N+XN+PR+Company
Adolf Wurth & Co KG,.N+XN+PR+Company
Advanta Corp,.N+XN+PR+Company

Abbildung 5.17: Auszug aus dem Lexikon Companies-.dic

Abkürzung	Kategorietyp	wörtliche Bedeutung	Erläuterung
N	grammatikalisch	Noun	Nomen
XN	grammatikalisch	Extended Noun	Mehrwortlexem
PR	semantisch	Proper Noun	Eigenname
Company	semantisch	Company Name	Organisationsname bzw.
			Firmenname

Tabelle 5.10: Abkürzungen, die in Companies-.dic verwendet werden.

#### 5.2.8.3 Organisationsspezifische Adjektive

Es gibt gewisse Adjektive oder adjektivisch verwendete Partizipkonstruktionen, welche vornehmlich im linken Kontext von Organisationsbezeichnungen oder -namen auftreten. Die gebräuchlichsten organisationsspezifischen Adjektive extrahierte Friederike Mallchok schon 2004 aus den Kontexten der Organisationsnamen, die ihr NER-System [55] im Reuters Korpus erkannte, und fasste sie zu einem Lexikon mit 110 Einträgen zusammen. Dieses hat den Namen org\_adj-.dic und wird auszugsweise in Abbildung 5.18 dargestellt.

```
foreign,.org_adj
newly formed,.org_adj
privatized,.org_adj
profitable,.org_adj
worldwide,.org_adj
```

Abbildung 5.18: Auszug aus dem Lexikon org\_adj-.dic

#### 5.2.8.4 Organisationsspezifische Kontexte

Außerdem sammelte Friederike Mallchok noch die linken und rechten prädikativen Kontexte von Organisationsnamen. Da im Zuge meiner Untersuchungen der Kontexte von Menschenbezeichnern nur Organisationsnamen in Objektposition interessant sind, werde ich lediglich auf ihr Wörterbuch context\_before-.dic zurückgreifen, welche den linken Kontext von Firmennamen spezifiziert.

```
authorized the,.context_before
available from,.context_before
award from the,.context_before
awarded to,.context_before
```

Abbildung 5.19: Auszug aus dem Lexikon context\_before-.dic

## 5.2.9 Lexika der geographischen Begriffe

Toponyme sind aus Biographien eigentlich kaum wegzudenken. Mindestens einmal wird in jedem Lebenslauf eine Ortsbestimmung genannt - der Geburtsort. Doch manchmal sind Biographien wahre Bewegungsprofile, was nichts anderes heißt, als dass die betreffende Person vielleicht oft den Ausbildungsort, den Wohnort oder den Arbeitsort gewechselt hat. Um die Erkennung von geographischen Begriffen in biographischen Kontexten zu erleichtern, ist es von Vorteil ein Toponymlexikon zu erstellen.

Aus diversen Internetquellen [59] [97] wurden geographische Entitäten wie Länder, Kontinente, Bezirke, Provinzen, Regionen, Départements, Grafschaften, US Bundesstaaten, Bundesländer, Städte und Stadtteile extrahiert und in den beiden Wörterbüchern GeosMapplanet-.dic und GeosWikipedia-.dic gespeichert. Des Weiteren war Sebastian Nagel so freundlich, mir sein Toponymlexikon (GeosSebastian+.dic) zu überlassen, so dass dem System nun insgesamt 286900 geographische Entitäten zur Identifikation von Ortsbestimmungen im Kontext von Menschenbezeichnern zur Verfügung standen.

Die folgende Aufstellung 5.11 gibt einen detaillierten Überblick zu allen in den Lexika verwendeten Abkürzungen. Diese Übersicht beschränkt sich ganz auf die grammatikalischen und semantischen Zusatzinformationen für nominale Ortsbestimmungen. Toponyme, welche als Adjektive im Text fungieren, werden an späterer Stelle ausführlich behandelt.

Abkürzung	Kategorietyp	wörtliche Bedeutung	Erläuterung
N	grammatikalisch	Noun	Nomen
XN	grammatikalisch	Extended Noun	Mehrwortlexem
City	semantisch	City	Metropole, Stadt
NYCBourough	semantisch	New York City Bourough	Stadtteil von New
			York City
Borough	semantisch	Borough	Stadtteil, Stadtbezirk
CaProvince	semantisch	Canadian Province	Kanadische Provinz
Département	semantisch	Département	Département
County	semantisch	County	Grafschaft
Region	semantisch	Region	Region, Gebiet
USstate	semantisch	US State	US Bundesstaat
Nation	semantisch	Nation	Land
Continent	semantisch	Continent	Kontinent
GEO	semantisch	Geographical Term	Toponym, geographi-
			scher Begriff

Tabelle 5.11: Abkürzungen, die in Geos\*-.dic für Nomina verwendet werden.

### 5.2.9.1 Kontinente und Länderbezeichungen

Eine relativ kleine Teilmenge der Toponyme sind die Länderbezeichnungen, welche fast vollständig im Lexikon GeosWikipedia-.dic kodiert sind und ungefähr 2,4% der insgesamt 18000 Wörterbucheinträge ausmachen.

```
Africa, N+Continent+GEO
Algeria, N+Nation+GEO
America, N+Continent+GEO
Antigua and Barbuda, N+XN+Nation+GEO
Argentina, N+Nation+GEO
Austria, N+Nation+GEO
Bahrain, N+Nation+GEO
Belarus, N+Nation+GEO
Belgium, N+Nation+GEO
Bolivia, N+Nation+GEO
Bosnia and Herzegovina, N+XN+Nation+GEO
Europe, N+Continent+GEO
```

Abbildung 5.20: Auszug aus dem Lexikon Geos Wikipedia-. dic

#### 5.2.9.2 Städtenamen

Das Lexikon Geos Mapplanet-. dic [59] besteht ausschließlich aus Städtenamen und enthält ungefähr 25400 Einträge. Dazu kommen noch ungefähr 15100 Städtebezeichnungen aus Geos Wikipedia-. dic [97], welche die größte Subkategorie (nahezu 84%) in diesem Wörterbuch bilden.

```
Ahuachapán,.N+City+GEO
Aix-en-Provence,.N+City+GEO
Athens,.N+City+GEO
Bad Hersfeld,.N+City+GEO
Baden-Baden,.N+City+GEO
Bahçe,.N+City+GEO
Bahía Blanca,.N+City+GEO
Cap-Haïtien,.N+City+GEO
Cap-de-la-Madeleine,.N+City+GEO
Castellammare del Golfo,.N+City+GEO
```

Abbildung 5.21: Auszug aus dem Lexikon GeosMapplanet-.dic

#### 5.2.9.3 Grafschaften, Regionen, Bezirke und Départements

Eine weitere relativ kleine Gruppe bilden die regionalen Bezeichnungen im Lexikon Geos Wikipedia-.dic. Etwa 12% der Lexikoneinträge sind Bezeichnungen für englische Grafschaften, europäische Regionen, kanadische Provinzen, deutsche Bundesländer und französische Départements (Verwaltungseinheiten). Da Ortsnamen häufig durch den Verwaltungsbezirk, in dem sie liegen, spezifiziert werden, können diese gesammelten Einträge später durchaus von Nutzen sein.

Alsace,.N+Region+GEO
Andalusia,.N+Region+GEO
Basse-Navarre,.N+County+GEO
Brooklyn,.N+Borough+NYCBourough+GEO
Fairfield,.N+County+GEO
Manhattan,.N+Borough+NYCBourough+GEO
Newfoundland and Labrador,.N+CaProvince+GEO
Val-de-Marne,.N+Département+County+GEO
Warwickshire,.N+County+GEO

Abbildung 5.22: Auszug aus dem Lexikon Geos Wikipedia-.dic

### 5.2.9.4 US Bundesstaaten und ihre typischen Abkürzungen

Die US amerikanischen Bundesstaaten sind ihrem syntaktischen Verhalten ähnlich wie die in Abschnitt 5.2.9.3 genannten Gebiete. Dass sie gesondert aufgeführt werden, liegt an ihrer abgekürzten Schreibweise, die meist häufiger als die ausgeschriebene Form in Texten vorkommt. Somit ist es sinnvoll, beide Varianten - die offizielle Bezeichnung des Bundesstaates und die gebräuchliche Abkürzung - im Lexikon *USstates-.dic* aufzuführen.

Alaska,.N+USstate+GEO
Calif\.,.N+USstate+GEO
California,.A+AGEO+AState
Idaho,.N+USstate+GEO
Ida\.,.N+USstate+GEO
Louisiana,.N+USstate+GEO
Minnesota,.N+USstate+GEO
Minn\.,.N+USstate+GEO

Abbildung 5.23: Auszug aus dem Lexikon USstates-.dic

#### 5.2.9.5 Geographische Adjektive

Soll der Standort eines Unternehmens näher bestimmt werden, ist es in englischsprachigen Nachrichtentexten gebräuchlicher den geographischen Begriff als Adjektiv einer Firmenbezeichnung voranzustellen. In Regel gibt es zu jedem nominalen Toponym auch ein Adjektiv, welches in seiner grammatikalischen Form mit dem Nomen identisch sein oder von diesem abweichen kann. Auch existieren in der englischen Sprache grammatikalische Gesetze zur Bildung der Adjektivform einer Ortsbestimmung, die jedoch für die US Bundesstaaten relativ locker gesehen werden. Beispielsweise ist es grammatikalisch korrekt "the <u>Florida company Miller International"</u> in einem Text zu schreiben, doch hat sich auch die Variante "the <u>Floridian company Miller International"</u> eingebürgert. Eigentlich ist die zweite Möglichkeit grammatikalisch falsch, hat sich aber durch den ständigen Gebrauch bei Nicht-Muttersprachlern so weit verbreitet, dass man diese Form auf keinen Fall vernachlässigen sollte.

European,.A+AGEO
Florida,.A+AGEO+AState
Floridian,.A+AGEO+AState
Galway,.A+AGEO+ACity
Genevese,.A+AGEO+ACity
German,.A+AGEO+ANation
Icelandic,.A+AGEO+ANation

Abbildung 5.24: Auszug aus dem Lexikon GeosWikipedia-.dic

Abkürzung	Kategorietyp	wörtliche Bedeutung	Erläuterung
A	grammatikalisch	Adjective	Adjektiv
XA	grammatikalisch	Extended Adjective	lexikalische Einheit, welche
			die Funktion eines Adjek-
			tivs erfüllt
ACity	semantisch	City Adjective	Adjektiv für eine Metropo-
			le, Stadt
ABourough	semantisch	Bourough Adjective	Adjektiv für einen Stadtteil
AProvince	semantisch	Province Adjective	Adjektiv für eine Provinz
AState	semantisch	State Adjective	Adjektiv für einen Bundes-
			staat
ANation	semantisch	Nation Adjective	Adjektiv für ein Land
AGEO	semantisch	Geographical Term	Adjektiv für ein Toponym,
		Adjective	einen geographischen Be-
			griff

**Tabelle 5.12:** Abkürzungen, die in *Geos\*-.dic* für Adjektive verwendet werden.

## 5.2.10 Lexika der Temporalia

Was wäre ein Lebenslauf ohne Zeitbestimmungen? Fast jede biographisch relevante Information wird mit einer Zeitangabe oder einer Zeitspanne versehen. Natürlich wäre es unsinnig alle möglichen Datumsangaben in einem Lexikon zu sammeln. Die Aufgabe der Datumserkennung und des Auffindens von Zeitangaben aller Art wird später eine lokale Grammatik übernehmen, welche sich aber auf Begriffe stützt, die einen gewissen Zeitraum ausdrücken.

#### 5.2.10.1 Monatsnamen und -abkürzungen

Für die Lokalisierung von Datumsangaben sind Monatsnamen unverzichtbar. Innerhalb eines Datums kann ein Monat in ausgeschriebener oder abgekürzter Form vorliegen, wenn das Zahlenformat außen vor gelassen wird. Des Weiteren kann im Englischen eine Monatsabkürzung mit nachstehendem Punkt oder ohne diesen gebildet werden. Deshalb ist es hilfreich ein Lexikon für Monatsnamen (Month-dic) und deren Abkürzungen (Month-dic) anzulegen.

Die nachfolgenden Abbildungen 5.25 und 5.26 geben Einblick in die eben genannten Lexika. Dabei wird ersichtlich, dass jedem Eintrag das grammatikalische Merkmal N für Nomen und das semantische Merkmal Month für Monatsnamen bzw. MonthAbbr für Monatsabkürzungen mitgegeben wurde.

April,.N+Month
August,.N+Month
December,.N+Month
February,.N+Month
January,.N+Month
November,.N+Month
September,.N+Month

Abbildung 5.25: Auszug aus dem Lexikon Month-.dic

Apr,.MonthAbbr
Apr\.,.MonthAbbr
Aug,.MonthAbbr
Aug\.,.MonthAbbr
Dec,.MonthAbbr
Dec\.,.MonthAbbr
Feb,.MonthAbbr

Abbildung 5.26: Auszug aus dem Lexikon MonthAbbr-.dic

#### 5.2.10.2 Wochentage und ihre Abkürzungen

Analog zu den Monatsnamen wurde auch ein Lexikon für Wochentage und ihre jeweiligen Abkürzungsmöglichkeiten angelegt. Dabei wurde auf die Unterscheidung zwischen der Abkürzung und der Vollform verzichtet und lediglich DayOfWeek als semantisches Merkmal zur Ergänzung des Wörterbucheintrages gewählt.

Friday,.N+DayOfWeek
Monday,.N+DayOfWeek
Sat,.N+DayOfWeek
Saturday,.N+DayOfWeek
Sun,.N+DayOfWeek
Sun.,.N+DayOfWeek
Sunday,.N+DayOfWeek
Thursday,.N+DayOfWeek
Tue,.N+DayOfWeek
Tues,.N+DayOfWeek
Tues,.N+DayOfWeek

Abbildung 5.27: Auszug aus dem Lexikon DayOfWeek-.dic

### 5.2.10.3 Weitere zeitbezogene Nomina

Das Standardlexikon, auf welches das System Unitex [75] zurückgreift, kennt leider nur sehr wenige Nomina, die einen zeitlichen Aspekt ausdrücken. Um unbestimmte Zeitangaben, wie z.B. "in the afternoon", später in dem Korpus erkennen zu können, musste die semantische Klasse Ntime um einige Einträge erweitert werden.

afternoon,.N+Ntime
autumn,.N+Ntime
evening,.N+Ntime
fall,.N+Ntime
morning,.N+Ntime
night,.N+Ntime
today,.N+Ntime
tomorrow,.N+Ntime
winter,.N+Ntime
yesterday,.N+Ntime

**Abbildung 5.28:** Auszug aus dem Lexikon *Ntime-.dic* 

#### 5.2.11 Weitere Lexika

Sicherlich könnten noch mehr Wörterbücher erstellt werden, welche bei der Entwicklung von lokalen Grammatiken zur Erkennung von Menschenbezeichnern in biographischen Kontexten später recht hilfreich sind. Doch man sollte sich dabei auf das Wesentliche beschränken und nur Lexika für semantische Klassen (Kategorien) anlegen, welche unbedingt gebraucht werden.

Alle Wörterbücher, die bis zu diesem Zeitpunkt ausführlich beschrieben wurden, enthalten Entitäten oder Bezeichnungen, welche mit hoher Wahrscheinlichkeit im Umfeld von Menschenbezeichnern auftreten. Auf diese Weise vereinfachen sie die Lokalisierung von Menschenbezeichnern im Text und ermöglichen eine sehr genaue Beschreibung des biographischen Kontextes, in dem diese vorkommen.

Zudem sind diese Lexika aus diversen Vorüberlegungen, die mit der Analyse der Nachrichtentexte aus der Fincial Times (2004) einhergingen, entstanden. Dagegen wurde das folgende Lexikon während der Entwicklung verschiedener lokaler Grammatiken angelegt.

Im Verlauf der syntaktischen Studien des Korpuses zeigte sich, dass Zahlen in Wirtschaftsnachrichten häufig ausgeschrieben werden, wenn es darum geht, die Dauer eines Arbeitsverhältnisses anzugeben. Leider bietet das System Unitex grundsätzlich nur eine semantische Kategorie für arabische Zahlen, welche intern als NB bekannt ist.

Aus dieser Not heraus entstand das Wörterbuch *NBword-.dic*, welches alle englischen Zahlen von 1 bis 100 in ausgeschriebener Form, sowie Varianten von "hundred", "thousand", "million" und "billion" enthält. Analog zu NB für arabische Ziffern wird in diesem Wörterbuch die semantische Zusatzinformation NBword für Zahlen in Wortform eingeführt.

eight,.NBword
fifteen,.NBword
twenty-two,.NBword
thirty,.NBword
forty-seven,.NBword
fifty-six,.NBword
sixty-nine,.NBword
seventy-one,.NBword
million,.NBword
one billion,.NBword
a thousand,.NBword

Abbildung 5.29: Auszug aus dem Lexikon NBword-.dic

In Anhang B auf Seite 166 ist eine komplette Übersicht aller in den Wörterbüchern verwendeten grammatikalischen und semantischen Kategorien mit ihren jeweiligen Bedeutungen und entsprechenden Erläuterungen zu finden.

## 5.3 Verifikationsmöglichkeiten bei Google



Google stellt wohl die größte Menge englischsprachiger Texte zur Verfügung und ist somit eine Ressource, auf die gern zurückgegriffen wird.

#### Was bedeutet es nun, Google als Verifikationstool einzusetzen?

Wie die obige Abbildung zeigt, ist es möglich bei Google eine Anfrage mit Asterisk zu stellen, wie z.B. "\* was appointed as \*". Diese Schablone ermöglicht es die linken und rechten Kontexte des Prädikats "to be appointed as" aus Texten im Web zu ermitteln. Das führt schließlich zu Ergebnissen, wie diesen

In January 1990, Professor Chen San-ching was appointed as deputy director. ... and in 1997, he was appointed as Director.

Mr Thorn was appointed as Director, Office of Crime Prevention in August 2003.

Im Grunde bestätigen diese Ergebnisse nur die anfänglichen Vermutungen. Denn mit dieser Anfrage möchte man bestätigt bekommen, dass das Prädikat "to be appointed as" an seiner ersten Argumentposition eine Person in Form eines Eigennamens oder eines Personalpronomens hat, und dass auf die Präposition "as" eine Berufsbezeichnung folgen muss.

Ähnlich wie beim Bootstrapping mit lokalen Grammtiken könnte man auch hier neue Lexikoneinträge gewinnen. Nur dieses Mal geht es nicht darum, neue Informationen aus Texten zu extrahieren, sondern die Qualität einer Grammatik auf neuen Texten zu erproben. Wenn beispielsweise eine Grammatik für eine Verbrelation anhand eines bestimmten Korpuses entwickelt wird, werden einige syntaktische und semantische Schlussfolgerungen für die Umgebung des Verbs gemacht. Das führt dazu, dass die Grammatik mit der Zeit so gut wird, dass fast alle Vorkommen des Ausgangswortes in diesem Korpus gefunden werden. Dabei könnten wichtige Einschübe oder Variationen in der Satzstellung übersehen werden, die in anderen Artikeln durchaus gängig sind.

An dieser Stelle kommt nun Google ins Spiel, indem die "Aussage" der jeweiligen lokalen Grammatik konkretisiert und mit Wildcards an verschiedenen Positionen in der Phrase versehen wird. Unter der Konkretisierung einer lokalen Grammatik versteht man hier, dass ein Ausdruck der Form :date <LN> <HumVP> as <JD> bei Google als Anfrage nicht sinnvoll ist, dagegen aber "\* was appointed as \*" zugelassen wird und damit auch bestätigt werden, dass eine Instanz von <LN> oder <JD> im Umfeld einer Instanz von <HumVP> vorkommt.

Auf diese Weise lässt sich die Qualität der lokalen Grammatik an den entsprechenden Trefferquoten messen und die entsprechenden Passagen, welche mit Hilfe des Asterisks gefunden wurden, können entweder neu in der Grammatik berücksichtigt oder bestätigt werden.

# 6 Grammatik der Menschenbezeichner

Unter Verwendung der in Abschnitt 5.2 vorgestellten Wörterbuchressourcen lassen sich nun lokale Grammatiken zur Erkennung von Menschenbezeichnern innerhalb biographischer Relationen entwickeln.

Dabei stehen die Personen als zentrale Entität eindeutig im Vordergrund der linguistischen Untersuchungen. Aufgrunddessen werden sie am Anfang dieser strukturellen und semantischen Analyse stehen, gefolgt von weiteren Entitäten wie Organisationsnamen, Toponymen und Datumsangaben, welche im Umfeld von Personen vorkommen. Danach lassen sich die syntaktischen Zusammenhänge und biographisch relevanten Beziehungen zwischen den einzelnen Entitäten in Form von weiteren Grammatiken herstellen.

Wie bereits mehrfach angedeutet wurde, kann die semantische Kategorie der Menschenbezeichner in zwei Gruppen unterteilt werden. Diese Aufspaltung in Personennamen (Eigennamen) und in "allgemeine Menschenbezeichnungen" wird nicht ohne Grund gemacht. Natürlich wäre eine feinere Differenzierung innerhalb dieser beiden Gruppen noch möglich, doch zeichnen sich diese zwei Subkategorien durch klassentypische syntaktische Merkmale aus, so dass sie als zwei eigenständige Mengen betrachtet werden sollten.

## 6.1 Analyse von Personennamen

Eine dieser beiden Gruppen ist die der Personennamen. Leider handelt es sich hierbei um eine unendliche Menge, deren Elemente durch Kombination untereinander immer wieder neue Elemente in ihr hervorbringen. Auch können ihr jederzeit neue Namen von außen hinzugefügt werden, da im Laufe der Zeit immer mehr Bezeichnungen als neue Vornamen akzeptiert werden, welche zuvor nur bei fiktiven, literarischen Charakteren in Erscheinung getreten sind. So können beispielsweise "Nemo", "Smilla" oder "Aragorn" zu gebräuchlichen Vornamen werden<sup>42</sup> und somit zur Expansion dieser Menge beitragen. Deshalb ist es notwendig, Regeln anzugeben, welche die Struktur von menschlichen Eigennamen beschreiben, so dass auch lexikonfremde Personennamen in Texten gefunden werden können.

## 6.1.1 Syntaktische Variabilität bei Personennamen

Wie bereits in Abschnitt 1.3.1 auf Seite 14 angesprochen wurde, gibt es eine Vielzahl an syntaktischen Möglichkeiten, mit denen sich Personennamen darstellen lassen. Diese werden im Graphen person\_name.grf auf Seite 78 festgehalten.

Es gibt enorm viele Benennungen, welche sich auf eine einzige Person beziehen können.

<sup>&</sup>lt;sup>42</sup>Nachzulesen bei der Gesellschaft für deutsche Sprache e.V. http://www.gfds.de/namen1.html

Das folgende Beispiel illustriert im Fall von "Bill Gates" wie umfangreich diese Möglichkeiten sind, obwohl es längst nicht alle von ihnen erfasst.

Bill Gates
Bill Henry Gates
Bill Henry Gates III
Bill H. Gates III
Bill H. Gates
William Henry Gates III
William Gates
William Henry Gates
William H. Gates III
William H. Gates
W. H. Gates
B. Gates
Mr. Gates
Gates
Gates

#### 6.1.1.1 Abkürzung vs. Vollform

Im Beispiel von "Bill Gates" ist offensichtlich, dass eindeutig William Henry Gates III der Geburtsname ist, obwohl Bill Gates wohl die gebräuchlichere Form des Namens ist. Des Weiteren wird deutlich, dass die syntaktische Kombination aus Kurzformen (z.B. Bill), Abkürzungen wie W. H. und Vollformen (z.B. William Henry) der Vornamen mit dem entsprechenden Nachnamen sehr viele Variationen für einen Personennamen erzeugen. Auch der Einsatz von Anredeformen und Titelbezeichnungen in Verbindung mit dem Nachnamen oder dem Vor- und Zunamen, der wiederum in einer verkürzten Form oder in der Vollform vorliegen kann, erhöht die Anzahl der syntaktischen Möglichkeiten für den Namen "Bill Gates".

### 6.1.1.2 Synekdoche (Pars Pro Toto)

Wenn noch die Tatsache berücksichtigt wird, dass der Nachname bzw. der Vorname stellvertretend für den kompletten Personennamen im Text auftreten kann, wäre in einem formellen Dokument noch Gates bzw. Gates III und in einem persönlichen Schreiben Bill, William oder William Henry zu finden. Dieses Stilmittel ist in der Literatur unter zwei Bezeichnungen bekannt, wobei "Pars Pro Toto" der selbst erklärende Name ist, da ein Teil des Namens für den ganzen steht. Dennoch ist es auch unter dem Namen Synekdoche geläufig.

Wenn semantische Paraphrasierungen von Namen wie z.B. "the founder of Microsoft" zunächst außer Acht gelassen, und Formen mit der Anrede "Mr." nicht gesondert gezählt werden, so sind insgesamt über 100 Varianten des Namens "Bill Gates" möglich. (siehe dazu Anhang C auf Seite 168)

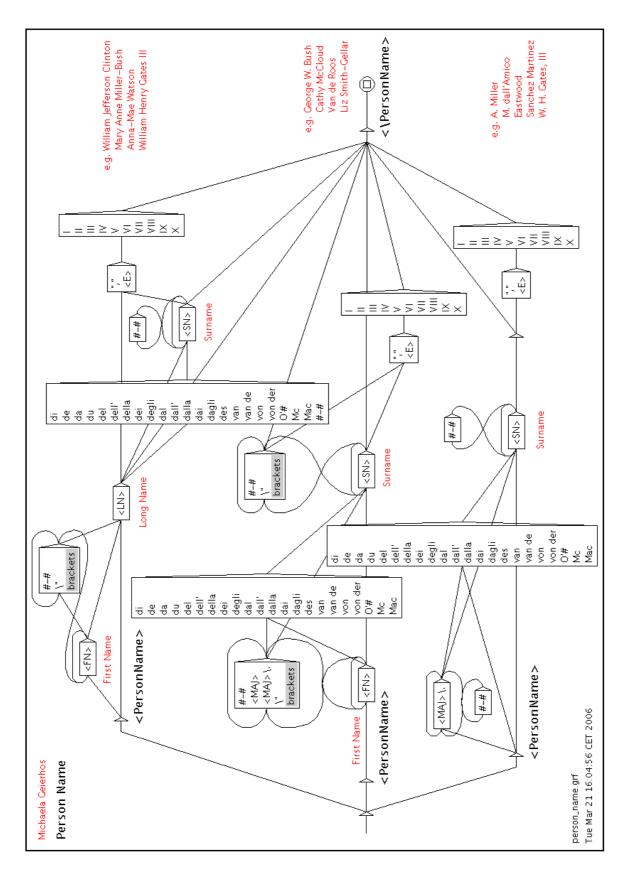


Abbildung 6.1: Graph zur Erkennung von Personennamen - person\_name.grf

## 6.1.2 Disambiguierung von "Scheinnamen"

Wird der Automat aus Abbildung 6.1 ohne weitere Kontextinformationen auf das Korpus angewendet, so besteht die Möglichkeit, dass Sequenzen als Personennamen identifiziert werden, welche sich in diesem Zusammenhang nicht direkt auf eine Person beziehen.

Ein klassisches Beispiel hierfür ist die Erkennung von Personennamen innerhalb von Firmenbezeichnungen.

```
at <PersonName>Gartner<\PersonName>, Inc.
```

Zwar wurde "Gartner" als Nachname im Text gefunden, doch fehlt die Information, dass es sich an dieser Stelle um eine Firma handelt, welche nach ihrem Gründer Gideon Gartner benannt ist.

Um diese Passage im Korpus der richtigen Entität - einer Organisation - zuzuweisen, muss der nachfolgende, sowie der vorangehende Kontext in die syntaktischen Untersuchungen miteinbezogen werden. Deshalb fungiert der Transduktor person\_name.grf als Subgraph im Graphen zur Lokalisierung von Organisationsnamen, wie in Abbildung 7.1 auf Seite 87 zu erkennen ist. Dieser endliche Automat mit dem Namen company.grf behebt damit den vermeintlichen Fehler der Personennamengrammatik und die gesamte Sequenz wird dementsprechend im Text annotiert.

```
at <ORG><PersonName>Gartner<\PersonName>, Inc.<\ORG>
```

Des Weiteren ist auch eine Verwechslung zwischen Personennamen und Toponymen möglich, welche nur über den Kontext aufgelöst werden kann, was manchmal keine leichte Aufgabe ist, wie das folgende Beispiel illustriert.

```
<PersonName>England<\PersonName>'s Queen
```

Hier wurde England fälschlicherweise als Nachname erkannt, obwohl es diese Funktion in anderen Fällen auch einnehmen kann:

```
<PersonName>England<\PersonName>'s mother, Terrie, told
    Secretary <PersonName>England<\PersonName> served as the
Mr. <PersonName>England<\PersonName> served as executive vice president
```

Wie diese Konkordanz zeigt, sollte eine nähere Spezifikation des Umfeldes von potentiellen Personennamen, die Fehlerquote deutlich eingrenzen, indem ein nachfolgender Menschenzeichner wie "mother", eine voranstehende Berufsbezeichnung wie "Secretary" oder eine Anredeform wie "Mr." sichere Indizien für die "Echtheit" eines Personennamens sind.

Doch wird sich das Problem mit "England's Queen" nicht ausschließlich über den direkten Kontext lösen lassen, denn "queen" ist ein Menschenbezeichner wie "mother" und nimmt in dieser Genetivkonstruktion die gleiche Rolle ein. Nun könnte man damit argumentieren, dass immer nur ein Land an der ersten Position dieser Sequenz stehen wird, doch wäre "Johnny's Queen"<sup>43</sup> schon ein Gegenbeispiel für diese Annahme.

An dieser Stelle muss entweder der Kontext sehr detailliert beschrieben, oder es sollten Prioritäten gesetzt werden, bei denen an der genannten Vermutung festgehalten wird.

<sup>43 &</sup>quot;Johnny's Queen" ist der Name eines finnischen Rennpferdes (http://hippos.ip-finland.com/ hippos/tulokset/120020829.php)

## 6.1.3 Vervollständigung des Personennamenlexikons

Der Graph auf Seite 78 verlässt sich beim Auffinden von Personennamen ausschließlich auf Vor- und Nachnamen, sowie auf vollständige Personennamen aus den entsprechenden Wörterbüchern. Dabei wurde gezielt darauf verzichtet, Variablen für mögliche Nachnamen oder Vornamen einzuführen, weil das Risiko zu hoch ist, andere Eigennamen zu erkennen, die nichts mit Personen zu tun haben. Nur unter Berücksichtigung des Kontextes kann dieses Verfahren dazu genutzt werden, neue Personennamen zu finden, welche noch nicht in den Lexika enthalten sind.

Der endliche Automat pot\_person\_name.grf arbeitet nach diesem Prinzip und versucht weitere Personennamen im Korpus zu lokalisieren.

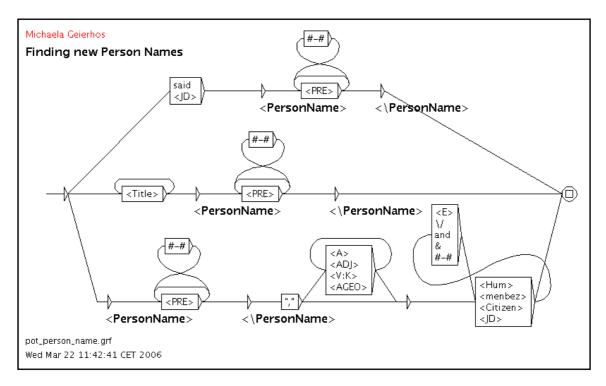


Abbildung 6.2: Graph zur Erkennung potentieller Personennamen - pot\_person\_name.grf

Bei der Suche nach möglichen Personennamen konzentriert sich dieser Graph auf die vier häufigsten Umgebungen, in denen Namen anzutreffen sind.

Es werden drei Varianten in Betracht gezogen, welche den linken Kontext eines Personennamens beschreiben.

Ähnlich wie schon Nathalie Friburger (siehe Abschnitt 3.3.2.1) Verben der Äußerung nutzte, um Personen im Korpus zu finden, wird hier stellvertretend eine Vergangenheitsform des Verbs "to say" verwendet.

said <PersonName>Allen Todd<\PersonName>, Director of Utilities
said <PersonName>Clive Davis<\PersonName>, Chairman and CEO
said <PersonName>James Moore<\PersonName>, analyst
said <PersonName>Steve Witt<\PersonName>, vice president and general manager

Außerdem sind Berufsbezeichnungen eindeutige Anzeichen für nachfolgende Personennamen, wenn die daran anschließenden Begriffe groß geschrieben werden.

Auch verschiedene Anredeformen oder Titel weisen mit hoher Wahrscheinlichkeit auf nachstehende Personennamen in Form von Nachnamen oder Vor- und Zunamen hin.

```
Mr. <PersonName>Jose Maria Aznar<\PersonName>
    Mrs. <PersonName>Diana Clark<\PersonName>
    Prof. <PersonName>Bahram Jalali<\PersonName>
Lord <PersonName>Andrew Lloyd Webber<\PersonName>
```

Überdies sind auch im rechten Kontext von Personennamen häufig Berufsbezeichner zu finden, welche nur durch ein Komma vom Eigennamen getrennt sind.

So lassen sich relativ gute Ergebnisse erzielen, welche mit Hilfe ihrer Markierung aus der Konkordanz automatisch extrahiert, mit den vorhandenen Wörterbüchern abgeglichen und als neue Einträge in die entsprechenden Lexika aufgenommen werden.

## 6.2 Allgemeine Menschenbezeichner

Der eben in Abbildung 6.2 vorgestellte Graph zur Erkennung potentieller Personennamen nutzte bereits Berufsbezeichnungen für die nähere Beschreibung des Kontextes von menschlichen Eigennamen. Daran sieht man, dass allgemeine Menschenbezeichner - in Form von Berufen, Nationalitäten, Verwandschaftsbezeichnungen, usw. - oft im direkten Umfeld von Personennamen zu suchen sind. Natürlich werden beide auch getrennt voneinander im Text auftreten. Trotz ihrer Eigenständigkeit besteht dennoch eine gewisse Verbundenheit, wodurch sich syntaktische Zusammenhänge zwischen den beiden Kategorien herstellen lassen. Unter Berücksichtigung dieser Tatsachen wurde ein gemeinsamer Automat zur Spezifizierung von Personenbezeichnungen entwickelt.

Der Transduktor in Abbildung 6.3 auf Seite 82 vereint die Grammatik der Personennamen person\_name.grf mit der Kontextinformation aus dem Graphen pot\_person\_name.grf und verflechtet sie mit weiteren Möglichkeiten, wie allgemeine Menschenbezeichner im Korpus vorkommen können.

Dabei werden diese über die Symbole <hum>, <menbez>, <Citizen> und <JD> in den entspechenden Wörterbüchern nachgeschlagen bzw. über den Subgraphen jd.grf zur Erkennung komplexer Berufsbezeichnungen (siehe Abbildung 11.5 auf Seite 137) genauer spezifiziert.

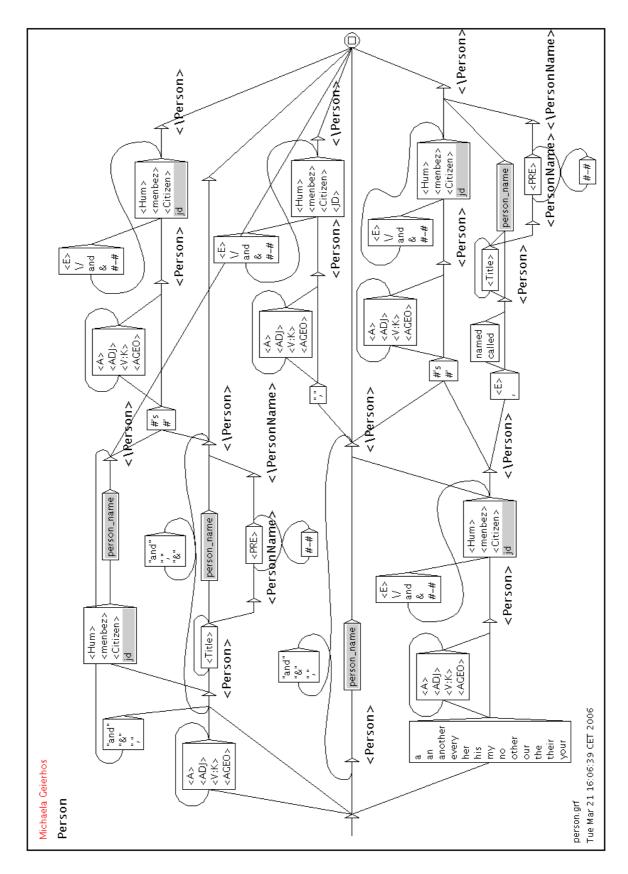


Abbildung 6.3: Graph zur Erkennung von Menschenbezeichnern - person.grf

```
said Panthers <Person><JD>defensive coordinator<\JD><PersonName>Mike Trgovac<\PersonName><\Person>
                         said <Person><JD>Irish manager<\JD><PersonName>Brian Kerr<\PersonName><\Person>
           said <Person><JD>lawyer<\JD><PersonName>Romeo Alcantara<\PersonName><\Person>, Pasig election officer
                said <Person><JD>president<\JD> and <JD>CEO<\JD><PersonName>Barry Shaked<\PersonName><\Person>
<Person><JD>Spanish Prime MinisterAria) AznarAznarAznarAznarAznarAznarAznarAznarAznarAznarAznarAznarAznarAznarAznarAznarAznarAznarAznarAznarAznarAznarAznarAznarAznarAznarAznarAznarAznarAznarAznarAznarAznarAznarAznarAznarAznarAznarAznarAznarAznarAznarAznarAznarAznarAznarAznarAznarAznarAznarAznarAznarAznarAznarAznarAznarAznarAznarAznarAznarAznarAznarAznarAznarAznarAznarAznarAznarAznarAznarAznarAznarAznarAznarAznarAznarAznarAznarAznarAznarAznarAznarAznarAznarAznarAznarAznarAznarAznarAznarAznarAznarAznarAznarAznarAznarAznarAznarAznarAznarAznarAznarAznarAznarAznarAznarAznarAznarAznarAznarAznarAznarAznarAznarAznarAznarAznarAznarAznarAznarAznarAznarAznarAznar
          <Person><ORG>Senate Intelligence Committee<\ORG><JD>Vice Chairman<\JD><PersonName>John D. Rockefeller IV
                                                            (<\PersonName><\Person>
               warns <PersonName>Avery Shenfeld<\PersonName><\Person>, a CIBC World Markets economist
                              (She) danced with <Person><PersonName>Baryshnikov<\PersonName><\Person>
     said <Person><PersonName>Bill Wise</PersonName>on>, <Person>presidentPerson>, North American Operations
 stated <Person><PersonName>Cary Losson<\PersonName><\Person>, <Person>Founder and President<\Person> of 1031 Exchange
            said <Person><PersonName>Dan Rothfeld<\PersonName><\Person>, <Person>senior vice president<\Person>
           says <Person><PersonName>Daniel Cohn-Bendit<\PersonName><\Person>, who led students to the barricades
    stated <Person><PersonName>Eric Kuhn<\PersonName><\Person>, <Person>Chairman and Chief Executive Officer<\Person>
said <Person><PersonName>R. Richard Fontaine</personName></person>, <Person>Chairman & Chief Executive Officer</Person>
     said <Person><PersonName>Steve Witt<\PersonName><\Person> , <Person>vice president and general manager<\Person>
           four years ago [Russian President] <Person><PersonName>Vladimir Putin<\PersonName><\Person> announced
                <Person>«JD>president<\JD> and <JD>CEO<\JD><PersonName>Doron Inbar<\PersonName><\Person> said
                                said <Person>lawyer<PersonName>Romeo Alcantara<\PersonName><\Person>
                         state immunity in the same way as a <Person><JD>foreign minister<\JD><\Person>
                               her face her father by pretending to be her <Person>husband<\Person>
        would cost the same for us, whether the <Person>traveller<\Person> was booking it from Australia or France
the normally unexcitable <Person><ORG>Nokia<\ORG><JD>chief executive<\JD><PersonName>Jorma Ollila<\PersonName><\Person>
```

Abbildung 6.4: Konkordanz zum Graphen person.grf

## 6.3 Anaphern auflösen

Bis jetzt wurde davon ausgegangen, dass gezielt nach Menschenbezeichnern in Form von Eigennamen oder allgemeinen Benennungen innerhalb diverser biographischer Kontexte in Wirtschaftsnachrichten gesucht wird. Doch sollte man nicht außer Acht lassen, dass bei einer Aneinanderreihung biographischer Informationen, der Personenname wahrscheinlich nur zu Anfang und im späteren Verlauf der Ausführungen nur noch selten vorkommen wird. Stattdessen wird ein Personalpronomen die Position des Eigennamens einnehmen. Obwohl sich der Schwerpunkt dieser Arbeit auf die Erkennung expliziter Personenbezeichnungen konzentiert, muss man an dieser Stelle eingestehen, dass biographische Relationen, welche ein Personalpronomen als Subjekt haben, nicht von den hier präsentierten Grammatiken erfasst werden. Dieses Defizit ließe sich durch den Transduktor aus Abbildung 6.5 beheben, da dieser - ähnlich wie in der Diskursanalyse - versucht, die Anaphern aufzulösen, und somit die Rückbezüge der Personalpronomen auf vorangegangene Personennamen wiederherzustellen. Dabei könnte das Korpus durch den Automaten so manipuliert werden, dass das jeweilge Personalpronomen durch den entsprechenden Eigennamen ersetzt wird. Bei der Konkordanz in Abbildung 6.6 auf Seite 85 wurde allerdings noch ein Zwischenschritt vorgenommen, welcher die Zuordnung zwischen einem Personalpronomen und dem jeweiligen Bezugsnamen illustriert.

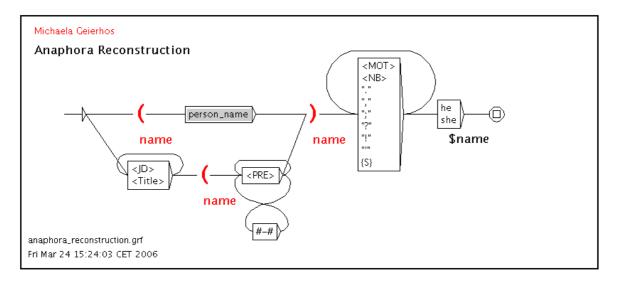


Abbildung 6.5: Graph zum Auflösen von Anaphern - anaphora.grf

Dennoch ist das Auflösen von Anaphern ein sehr risikobehaftetes Unterfangen, falls in den Lexika nicht das Geschlecht zu den jeweiligen Namen vermerkt wurde. Ohne diese Information lässt sich bei einem Namen nur feststellen, ob er zu einem Mann oder einer Frau gehört, wenn zuvor die Anrede "Mr.", "Mrs." oder "Miss" gefallen ist. So besteht die Gefahr, dass ein falscher Bezug zwischen dem Pronomen und dem Namen hergestellt wird. Dabei könnte beispielsweise ein weiblicher Personenname mit dem Personalpronomen "he" in Verbindung gebracht wird, da die Kongruenz der Geschlechter aufgrund fehlender Informationen nicht gewährleistet werden kann, und immer der letztgenannte Eigenname aufgegriffen wird.

```
said <PersonName>Christine Aguilera<\PersonName> as (Christine Aguilera) she accepted the best female pop vocal performance cooperation of <PersonName>Alfred Allington<\PersonName> who invited us in like (Alfred Allington) he did with everyone else I see no reason why <PersonName>Charles Allen<\PersonName> should not succeed, but if (Charles Allen) he does not, then sadly that is a fact of life
```

I have been talking to <a href="PersonName">PersonName</a>, the Meath secretary, and (Barney Allen) he confirmed to me that Gary was getting an invitation

I think <PersonName>Andy Blignaut<\PersonName> has been as quick of some of the others, (Andy Blignaut) he is up in the 140s I am sure his coach <PersonName>Malcolm Arnold<\PersonName> would have worked that out and I am sure (Malcolm Arnold) he has done a hard regime of quantity and quality

His son <PersonName>Aidan Barclay<\PersonName> . . . gave helpful testimony but admitted at many points that (Aidan Barclay) he did not know what his father had actually done or why.

If <PersonName>Tony Blair<\PersonName> thought (Tony Blair) he had scars on his back from trying to reform the public services two or three years ago

Every year -- every single year <PersonName>George Bush<\PersonName> has promised to create jobs and every year (George Bush) he is ended up losing them

How can we trust President <PersonName>Bush<\PersonName> to create 2.6 million jobs when (Bush) he has the worst record since Herbert Hoover

If <PersonName>Francisco Carrasquero<\PersonName> has a military complex, then (Francisco Carrasquero) he cannot be chairman of the CNE

Abbildung 6.6: Konkordanz zum Graphen anaphora\_reconstruction.grf

# 7 Grammatik der Organisationsnamen

In biographischen Texten nehmen Beschäftigungsverhältnisse oft einen relativ hohen Stellenwert ein, so dass die Spezifikation eines Arbeitsverhältnisses unumgänglich ist, wenn biographisch relevante Informationen aus Nachrichten extrahiert werden sollen. Zunächst kann die Beziehung zwischen einer Person und einer Firma auf diese beiden Entitäten reduziert werden. Wie diese Relation nun genau aussieht, soll erst zu einem späteren Zeitpunkt geklärt werden. Deshalb muss an dieser Stelle nur die syntaktische Struktur von Firmennamen analysiert und die darin verborgene semantische Klassifizierung von Organisationen vorgenommen werden.

Natürlich werden auch die Wörterbucheinträge aus Abschnitt 5.2.8 (siehe Seite 65) in die Entwicklung einer solchen Grammatik für Organisationsnamen einbezogen.

## 7.1 Syntaktische Variabilität bei Organisationsnamen

Ähnlich wie schon bei den Personennamen tritt eine gewisse syntaktische Variabilität in der Struktur von Organisationsnamen auf. Allerdings kann davon ausgegangen werden, dass die Vielfalt an Variationen bei Firmennamen deutlich geringer ist, als es bei den Personennamen der Fall war (siehe Abschnitt 6.1.1).

Der Graph aus Abbildung 7.1 auf Seite 87 behandelt diverse Möglichkeiten, in welcher Form eine Organisationsbezeichnung im Text auftreten kann.

So können Unternehmen einerseits mit einem Zusatz im Namen genannt werden, der Aufschluss über ihre jeweilige Rechtsform oder die Art des Gewerbes gibt.

```
<ORG>Novartis AG<\ORG> declined to comment on newspaper reports
Kim Manley, Chief Marketing Officer, <ORG>Allied Domecq PLC.<\ORG>.{S}
Ashkin, Executive Vice President of <ORG>America Online, Inc.<\ORG>.{S}
vice president, partner services, for <ORG>Choice Hotels.<\ORG>
company news service from the <ORG>London Stock Exchange Freeport PLC<\ORG>
Mike Gausling, President and CEO of <ORG>OraSure Technologies<\ORG>.{S}
```

Andererseits kommt ein Firmenname meist ohne diese Zusatzinformation im Text vor und wird dabei nur über das entsprechende Wörterbuch identifiziert, weil der Kontext nicht zur Disambiguierung herangezogen werden kann. Der Grund hierfür liegt am linguistischen Charakter des Organisationsnamens, der häufig im gleichen Zusammenhang wie ein Personenname verwendet wird. Außerdem gilt die Großschreibung sowohl für die Namen von Unternehmen als auch für die der Menschen. Somit bestehen kaum Chancen, eine eindeutige Abgrenzung zwischen Personen und Organisationen vorzunehmen, wenn nicht auf gewisse Lexikoninformationen oder firmentypische Indikatoren zurückgegriffen wird.

Des Weiteren erkennt der folgende Graph auch Ausdrücke, welche sich auf Institutionen, Abteilungen, Zweigstellen oder Lehreinrichtungen beziehen.

Jim Saunders of the <ORG>Department of Public Safety<\ORG> said said Brian Moyne, Project Manager, <ORG>Department of Transport<\ORG> According to a 1999 report by the <ORG>Institute of Medicine<\ORG> says Sir Howard Davies, director of the <ORG>London School of Economics<\ORG> said Claes Fornell, director of the <ORG>University of Michigan<\ORG>'s National Quality Research Center

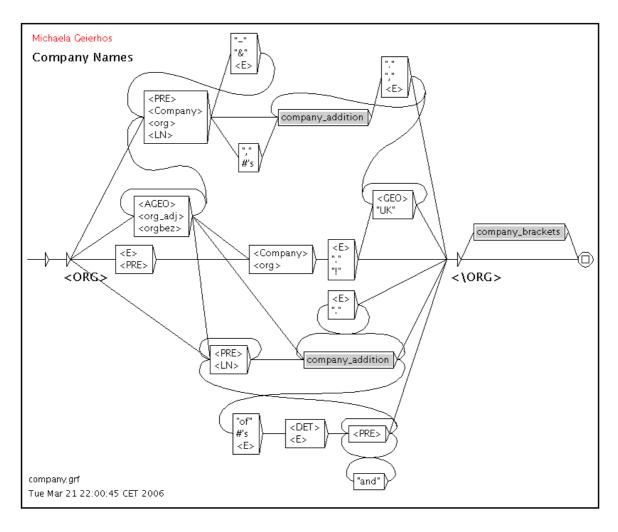
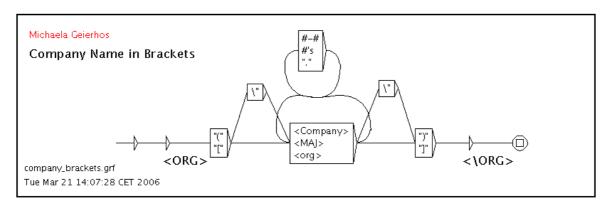


Abbildung 7.1: Graph zur Erkennung von Firmennamen - company.grf

Wenn das Akronym eines Firmennamens mindestens genauso bedeutend wie der eigentliche Organisationsname selbst ist, kann es durchaus vorkommen, dass es der Unternehmensbezeichnung in Klammern nachgestellt wird. In diesem Fall kommt der Subgraph company\_brackets.grf zum Einsatz.



**Abbildung 7.2:** Graph zur Erkennung von Firmennamen in geklammerten Ausdrücken - company\_brackets.grf

Dabei können beispielsweise folgende Textpassagen im Korpus gefunden und als Organisationsnamen gekennzeichnet werden:

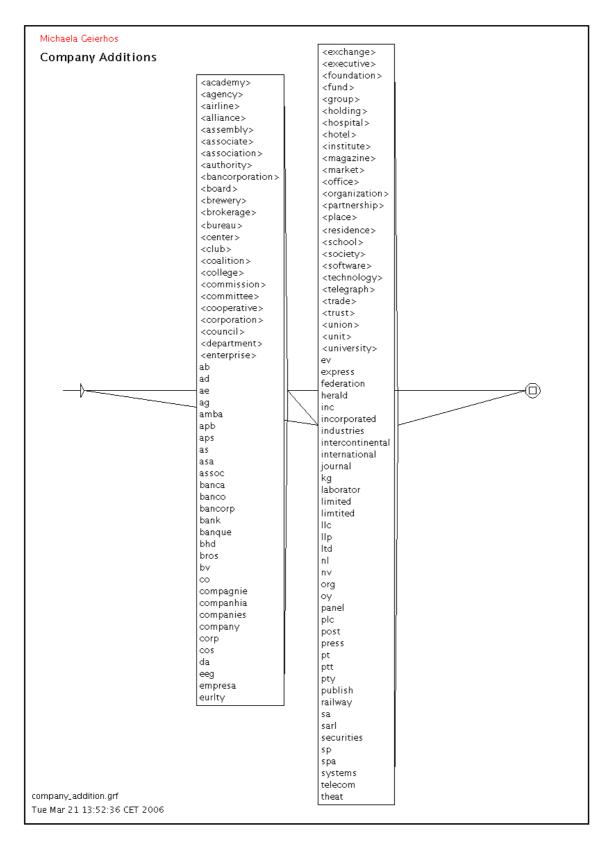
more than 21 million <ORG>American Online (AOL)<ORG> and CompuServe members with production at <ORG>Maschinenfabrik-Augsburg-Nuernberg (MAN)<ORG>

## 7.2 Abgrenzung von unechten Organisationsnamen

Da sich der endliche Automat *company.grf* aus Abbildung 7.1 auf sehr eingeschränkte Kontextinformationen stützt, besteht die Möglichkeit, dass Sequenzen irrtümlich als Firmennamen erkannt werden.

Demnach wären für diesen Graphen folgende Phrasen potentielle Fehlerquellen:

Diese vermeintlichen Organisationen entsprechen dem Schema eines Personennamens gefolgt von einem Firmenzusatz. Doch war in diesem Fall keiner der nachgestellten Firmennamenindikatoren für eine eindeutige Identifikation ausreichend. Da der Subgraph company\_additions.grf aus Abbildung 7.3 auf Seite 89 bei der Erkennung von Rechtsformen und anderen möglichen Zusätzen keine Rücksicht auf Groß- und Kleinschreibung nimmt und nicht im Text "nach vorn blickt", werden auch nachfolgende Mehrwortlexeme zerrissen bzw. nicht erkannt. Jedoch kann dies nur passieren, wenn der Transduktor ohne zusätzliche Informationen zum textuellen Umfeld des Unternehmens eingesetzt wird. Für die hier vorgestellten linguistischen Untersuchungen von biographischen Relationen, werden Firmennamen lediglich innerhalb bestimmter Satzstrukturen auftreten. Deshalb wird die Grammatik zur Erkennung von Organisationsnamen immer nur als Subgraph eingebettet in einem umfassenden Kontext - aufgerufen werden.



**Abbildung 7.3:** Graph zur Erkennung von Rechtsformen und weiteren Firmenzusätzen -  $company\_additions.grf$ 

## 7.3 Vervollständigung des Organisationsnamenlexikons

Ähnlich wie schon bei den Personennamen lassen sie Grammatiken auch für die Suche nach neuen Organisationsnamen einsetzen und erweitern auf diese Weise die Lexika.

Der folgende Graph bedient sich einerseits der linkstypischen Kontexte für Firmennamen und andererseits berücksichtigt er nachstehende Firmenzusätze aus dem Subgraphen company\_additions.grf. Dabei werden die von Friederike Mallchok gesammelten Kontextinformationen in Bezug auf Organisationen verwendet, welche über die Symbole <contextbefore>, <org\_adj> und <orgbez> abgerufen werden können.

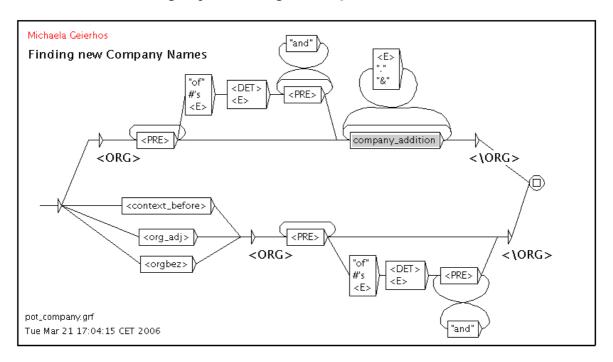


Abbildung 7.4: Graph zur Erkennung potentieller Firmennamen - pot\_company.grf

Im Grunde ist dieser Automat schon im Graphen *company.grf* in Abbildung 7.1 auf Seite 87 enthalten. Hier ist seine Aufgabe aber, potentielle Organisationsnamen zu lokalisieren, die nicht unbedingt in einem der Wörterbücher vorkommen.

So findet er beispielsweise folgende Firmennamen im Korpus:

```
<ORG>Above All Software<\ORG> was selected from the hundreds of companies
   a partner in <ORG>Accenture's Insurance Solution Group<\ORG>.{S}
   said Phil Flynn, an analyst with <ORG>Alaron Trading Corp.<\ORG>
   stepping in as a white knight to rescue <ORG>Aventis SA<\ORG> from
        President and CEO of <ORG>OraSure Technologies<\ORG>.{S}
   president of <ORG>Mobile Competency Inc.<\ORG>, a consulting firm.{S}
```

Die Treffer lassen sich nun automatisch aus der Konkordanz extrahieren, sollten aber noch sorgfältig durchgesehen werden, bevor sie mit den bestehenden Lexikoneinträgen abgeglichen werden. Als neue Wörterbucheinträge können sie später dabei helfen, die Qualität der Grammatik zu verbessern.

# 8 Grammatik der Ortsangaben

Geographische Begriffe treten in verschiedenen Zusammenhängen in Texten auf. Dabei ist die häufigste Aufgabe eines Toponyms die Angabe eines Ortes, an dem ein Ereignis stattgefunden hat oder stattfinden wird. In Anbetracht dessen, dass in dieser Arbeit Menschenbezeichner innerhalb biographischer Kontexte untersucht werden sollen, lassen sich die hier verwendeten Lokativa eher auf Geschehnisse in der Vergangenheit beziehen.

## 8.1 Biographische Relationen mit Ortsangaben

In der Menge der biographischen Relationen gibt es einige Prädikate, die meist zusammen mit einer Ortsangabe genannt werden.

Darunter fallen beispielsweise die Verben "to be born" und "to be raised":

```
Artemis Wines president Eric Smith, <u>born and raised</u> in <GEO>Chile<\GEO>, heads
Artemis Wines International

Born in <GEO>Long Island, New York<\GEO> on March 27, 1970, Carey moved to New York
City at the age of 17
```

Die hier erkannten Toponyme wurden mit Hilfe des Automaten aus Abbildung 8.1 auf Seite 93 im FT Korpus gefunden. Jedoch fungiert diese Grammatik nie als eigenständige Komponente bei der Lokalisierung von Ortsangaben, da sie grundsätzlich im Kontext anderer Transduktoren aufgerufen wird. Im obigen Beispiel wurde sie als Subgraph in den Graph der Verbalphrase "to be born (and raised)" eingebunden (siehe Abbildung 10.2 auf Seite 108).

Für das Auffinden geographischer Begriffe sind in der Grammatik einige "Wörterbuch-Lookups" notwendig. Die entsprechenden Lexika, welche dazu herangezogen werden, wurden bereits in Abschnitt 5.2.9 auf Seite 68 eingeführt.

So wird einerseits mit folgenden semantischen Wörterbuchkategorien gearbeitet:

- <Bourough> sucht nach einem Stadtteil.
- <CaProvince> sucht nach einer kanadischen Provinz.
- <City> sucht nach einer Stadt.
- <Continent> sucht nach einem Kontinent.
- <County> sucht nach einer Grafschaft.
- < Département > sucht nach einem französischen Département.
- <GEO> sucht nach Toponymen, die keine besondere Kennzeichnung haben.

- <NYCBourough> sucht nach Stadtteilen von New York City.
- <Nation> sucht nach Ländern.
- < Region > sucht nach Regionen.
- <USstate> sucht nach U.S. Bundesstaaten.

Andererseits werden noch folgende Schlüsselwörter zur Disambiguierung des Kontextes verwendet:

beach, bourough, city, county, country, district, province, region, state, town

In der Umgebung dieser Wörter sind mit hoher Wahrscheinlichkeit weitere geographische Begriffe zu finden, die eventuell noch nicht in einem der Lexika enthalten sind. So können im selben Schritt neue und bekannte Toponyme im Korpus gefunden werden.

Ein weiteres Indiz für eine nachfolgende Ortsangabe ist eine Himmelsrichtung wie North, West, East, South bzw. Northern, Western, Eastern, Southern oder andere ortsbegrenzende Adjektive wie Middle und Central.

Auf diese Weise lassen sich die unmittelbaren Kontexte von Toponymen auf die wesentlichen Bestandteile eingrenzen und tragen dazu bei, die Quote der falsch erkannten geographischen Begriffe zu minimieren.

# 8.2 Ortsangaben in ihrer Funktion als Attribute

Des Weiteren beziehen sich Lokativa nicht nur auf Ereignisse, sondern sie bestimmen auch Nominalphrasen näher, indem beispielsweise Aufgabengebiete von Menschen spezifiziert bzw. eingegrenzt werden, oder der Standort einer Firma als Zusatz in deren Namen wiedergegeben wird.

## 8.2.1 Toponyme als Attribut einer Berufsbezeichnung

Für die Erkennung von Berufsbezeichnungen ist zwar der Graph jd.grf auf Seite 137 zuständig, doch ruft dieser in seinem linken und rechten Kontext die Toponymgrammatik geo.grf auf. Somit werden beispielsweise folgende Berufsbezeichnungen im Umfeld von Ortsangaben lokalisiert:

```
said Peter Gregory, <GEO>England<\GEO>'s chief medical officer.{S}
said EMC product manager for <GEO>South Asia<\GEO> Ajaz Munsiff
```

# 8.2.2 Toponyme als Attribut eines Organisationsnamens

Analog dazu sucht der endliche Automat company.grf auf Seite 87 im rechten Kontext von Firmennamen nach geographischen Begriffen, welche den Sitz des Unternehmens angeben.

```
or of travel and fleet services for \underline{J\&J} <GEO>Europe<\GEO>.{S} Dean Tang, president and CEO of \underline{ABBYY} <GEO>USA<\GEO>.{S}
```

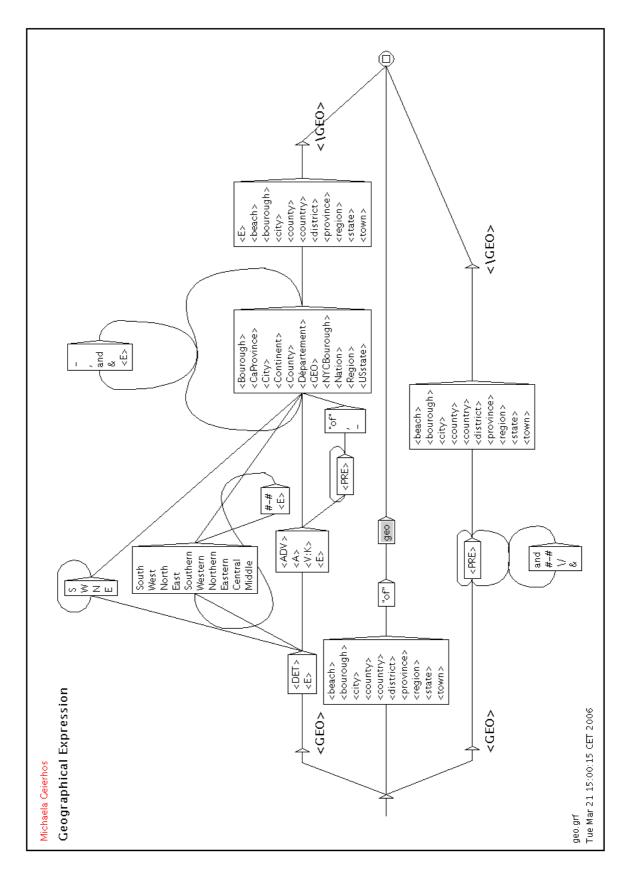


Abbildung 8.1: Graph zur Erkennung von Toponymen - geo.grf

```
AMC-10 will be delivering some of <GEO>America<\GEO>'s leading cable programs
               Caricom observers be in <GEO>Antigua<\GEO> and <GEO>Barbuda<\GEO> at this time
                       said EMC product manager for <GEO>South Asia<\GEO> Ajaz Munsiff
                 the traveller was booking it from <GEO>Australia<\GEO> or <GEO>France<\GEO>
                         for example to add a few branches in <GEO>California<\GEO>
                     said Peter Gregory, <GEO>England<\GEO>'s chief medical officer.{S}
                 the Tripartite Aggression against <GEO>Egypt<\GEO> was launched in 1956.{S}
             officer of Prescribed Solutions, a New York<\GEO>-based cosmeceuticals company.{S}
  the blowing up of homes in <GEO>Moscow<\GEO> and <GEO>Volgodonsk<\GEO> (which has never been proved).{S}
                          Dean Tang, president and CEO of ABBYY <GEO>USA<\GEO>. {S}
    Naomi Schwartz, SBCAG Chairwoman and <GEO>Santa Barbara<\GEO> County Supervisor (First District).{S}
                 Economic Group, a research firm in <GEO>Lansing<\GEO>, <GEO>Mich.<\GEO>.{S}
         companies have said they plan to list in <GEO>London<\GEO> but they are a long way from it
                       an analyst with Alaron Trading Corp. in <GEO>Chicago<\GEO>.{S}
              Dr Tien Wu, President of ASE Americas, <GEO>Europe<\GEO> and <GEO>Japan<\GEO>.{S}
     A mission of Delta Dental Plan of <GEO>Tennessee<\GEO> is to support programs of demonstrated value
manager Brian Little, aiming to emulate <GEO>Chesterfield<\GEO> and <GEO>Wycombe<\GEO>, Second Division clubs
           A commodities boom threatens to give <GEO>Canada<\GEO> a bad case of the Dutch disease
   House Democratic leader Nancy Pelosi (<GEO>Calif.<\GEO>) and other Democratic congressional leaders.{S}
                   a market research consultancy based in <GEO>Austin<\GEO>, TX<\GEO>.{S}
                        Sales Vice President Diana Clark, <GEO>Los Angeles<\GEO>.{S}
   that Italy's exclusion from the <GEO>Berlin<\GEO> meeting reflects a loss of influence over EU affairs
    number of companies in both <GEO>Germany<\GEO> and <GEO>Switzerland<\GEO> have expressed an interest
                        A car in <GEO>San Francisco<\GEO> would cost the same for us
```

Abbildung 8.2: Konkordanz zum Graphen geo.grf

# 9 Grammatik der Datumsangaben

Bereits von Anfang an war klar, dass ein System zur automatischen Erkennung biographischer Relationen nicht ohne eine lokale Grammatik für Datumsangaben auskommen wird. In Anlehnung an die Grammatiken, welche Maurice Gross zur Beschreibung genauer und ungefährer Datumsangaben erstellt hatte (siehe Gross, 1993 [40]), wurde ein Finite State Graph entwickelt, welcher Datumsangaben innerhalb biographischer Relationen als solche identifiziert.

Da in biographischen Kontexten meist präzise Datumsangaben gemacht werden, wird nur der entsprechende Referenzgraph von Maurice Gross in Abbildung 9.1 gezeigt. Wie nun der Graph für die Erkennung von Daten innerhalb biographischer Relationen in Wirtschaftsnachrichten aussieht, illustriert Abbildung 9.2 auf der folgenden Seite.

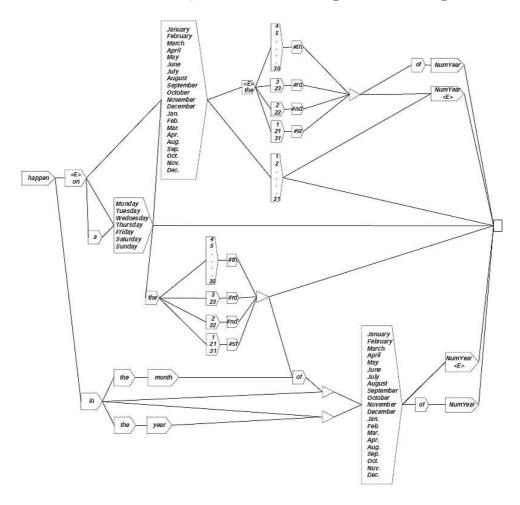


Abbildung 9.1: Graph zur Erkennung genauer Datumsangaben aus Gross, 1993 [40]

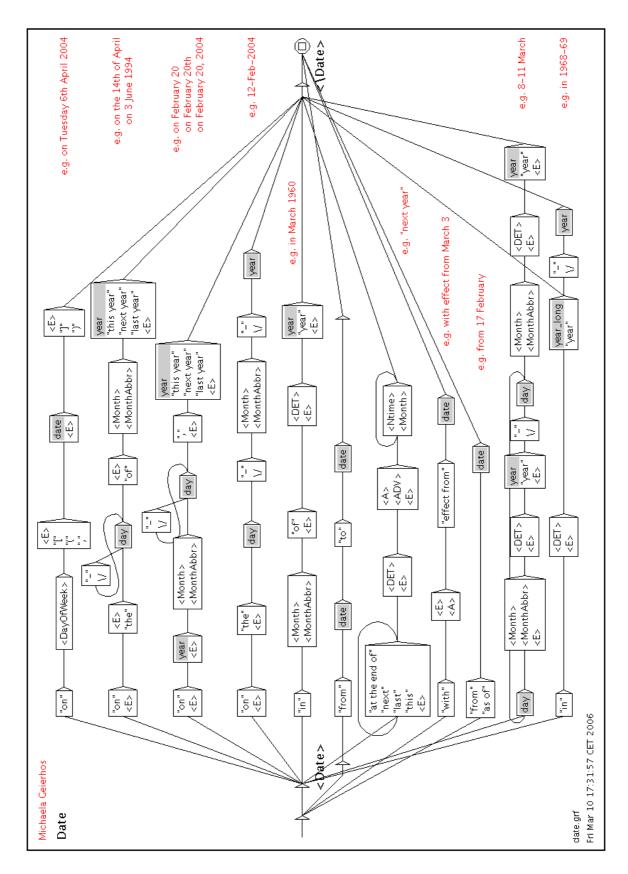


Abbildung 9.2: Graph zur Erkennung von Datumsangaben - date.grf

Da der Graph date.grf aus Abbildung 9.2 ein sehr komplexes Gebilde ist, lässt sich seine Vorgehensweise bei der Lokalisierung von Datumsangaben zunächst schwer nachvollziehen. Aufgrunddessen ist es durchaus sinnvoll die Arbeitsweise dieses Transduktors Schritt für Schritt nachzuvollziehen.

Dabei werden die wichtigsten Pfade des Graphen virtuell durchlaufen, um die Funktionsweise des Automaten zu erläutern.

Betrachtet man den ersten (obersten) Pfad des Graphen date.grf isoliert, so lassen sich deutlich zwei Varianten von Zeitangaben erkennen, welche dieser Automat beschreibt. Einerseits werden damit Textpassagen entdeckt, die Aufschluss darüber geben, an welchem Wochentag etwas geschehen ist, und andererseits wird diese Information noch mit einem genauen Datum versehen.

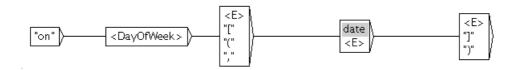


Abbildung 9.3: Erster Pfad aus dem Graphen date.grf

So würde dieser Ast des Graphen beispielsweise Ausdrücke der Form

on Saturday

on Monday, September 9, 2004

on Tuesday [Feb 4, 2005]

on Friday (January 7th, 1996)

on Wednesday May 3rd, 1998

matchen, wobei in diesem Pfad ein rekursiver Aufruf des gesamten Graphen date.grf erfolgt. Das heißt nun, dass dieser Ast andere Wege im Graphen miteinbezieht, welche wiederum die Erkennung des konkreten Datums vornehmen.

Dabei ist der reguläre Ausdruck on <DayOfWeek> für die Spezifikation des Wochentages zuständig, so dass beim Knoten <DayOfWeek> im entsprechenden Lexikon nachgeschlagen wird, welche Wochentage an dieser Stelle möglich sind (siehe Abschnitt 5.2.10.2). Die nachfolgenden Knoten können entweder leer durchlaufen werden (<E>) oder es wird im grau hinterlegenden Knoten date der Graph wieder selbst aufgerufen, um ein genaues Datum zu erkennen, welches möglicherweise geklammert auf den Wochentag folgt.

```
rname Beier, was due to appear in court lanthropist who came to their community sitional government, Bahnam Ziya Bulus, ry: President Mohammad Khatami's office results on record by the VLL board here charya.{S} Talking to presspersons here rls and four boys to the transit center charya.
```

Abbildung 9.4: Konkordanz zum ersten Pfad aus dem Graphen date. grf

Der zweite Pfad (von oben gezählt) im Graphen date.grf beschreibt nur eine Möglichkeit, wie ein konkretes Datum im Englischen aufgebaut sein kann.

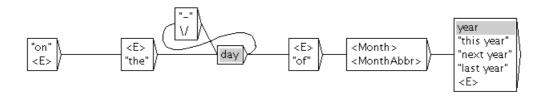


Abbildung 9.5: Zweiter Pfad aus dem Graphen date.grf

So würde er beispielsweise Textpassagen wie

31 May 1983 4th Sept 2004 on the 2nd of March 2005 on 5 Jan. 1998 on 19-20 June 2001 on 22/23 April last year

im Korpus erkennen.

Dabei werden zwei Subgraphen miteinbezogen, welche den Tag und die Jahreszahl spezifizieren. Auf diese Weise wird sichergestellt, dass keine Zahl größer 31 irrtümlicherweise als Tagesangabe identifiziert wird, und dass nur richtige Jahreszahlen als solche erkannt werden. Auf Seite 103 zeigen die Abbildung 9.16 und 9.17 die entsprechenden Automaten zur Erkennung von numerischen Tagesangaben und Jahreszahlen, welche hier nur über day und year in den jeweiligen Knoten referenziert werden.

Des Weiteren werden über die Ausdrücke <Month> und <MonthAbbr> alle möglichen Monatsnamen und Monatsabkürzungen in den beiden Wörterbüchern *Month-.dic* und *MonthAbbr-.dic* nachgeschlagen (siehe Abbildung 5.25 und 5.26 auf Seite 72).

Dabei sind auch wieder  $\epsilon$ -Transitionen möglich (<E>), welche beispielsweise garantieren, dass sich ein Datum eingeleitet von "on", als auch ein Datum ohne vorangehende Präposition finden lässt.

Die Konkordanz zu diesem Pfad des Transduktors illustriert, welche Form die Daten haben müssen, um von diesem Ast des Graphen entdeckt zu werden.

```
e new rates will remain in effect until

vice - United Kingdom Government News 14

(Date>9 February<\Date>
The minister for trade and industry, D

News National Security Summits & Talks 8

(Date>9 February<\Date>
Peshawar: Shah Raza, a Pakistani broad

measures against the chemical beginning
meyni landed at Tehran Mehrabad airport
ews Service Sports General News 90 site
an newspaper The Sunday Vision web site

(Date>0 1 February<\Date>
A mysterious disease has killed thr

The International Criminal Court (I
```

Abbildung 9.6: Konkordanz zum zweiten Pfad aus dem Graphen date. grf

Der nächste Weg durch den Automaten ist sehr ähnlich zum Muster des zweiten Pfades, nur dass hier die einzelnen Komponenten des Datums untereinander vertauscht wurden. Wie schon der vorherige Zweig des Graphen die Angabe mehrerer Tage innerhalb eines Datums berücksichtigt hat, lässt auch dieser die Verbindung von zwei Tagen durch '/' oder '-' zu.

Der vierte Pfad des Graphen date.grf entspricht in seiner internen syntaktischen Struktur der des zweiten Astes. Jedoch wird hier auf eine andere Verknüpfungsmethode zwischen Tag, Monat und Jahr Wert gelegt. In Nachrichtentexten erscheint oft am Anfang oder am Ende des Artikels die Datumsangabe als Tag-Monat-Jahr, so dass auch diese Möglichkeit für ein gültiges Datum in Betracht gezogen werden muss.

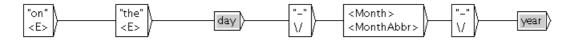


Abbildung 9.7: Vierter Pfad aus dem Graphen date.grf

So beschränkt sich dieses Suchmuster auf Daten folgender Form:

9-July-1956 on 24-Sept-2001 12/Mar/1995 on 23-October-2003 02/Jan./1972

Im Financial Times Korpus wiesen die Datumsangaben zu Beginn oder Ende des Berichtes grundsätzlich dieselbe Struktur auf, wie auch der folgende Ausschnitt aus der Konkordanz für diesen Pfad des Automaten zeigt.

```
CDate>31-Jan-2004<\Date>
CDate>31-Jan-2004<\Date>
Talk of the Town# Why remember the CPP purges?{S} Business Services Political
The Nuclear Noose Around Pakistan's Neck Executive Legislative & General Gene
CDate>31-Jan-2004<\Date>
CDate>31-Jan-2004<\Date>
The Price Of Winning At Any Cost washingtonpost.com Elections Government News
They Strut, They Fret, They Build! washingtonpost.com 1821 BAGHDAD As Ayad Al
Travel# Lanuza through the rain Philippine Daily Inquirer Government News 120
CDate>31-Jan-2004<\Date>
Turn of the tide - Prepackaged Software Chemical Preparations NEC Plastics Ma
CDate>31-Jan-2004<\Date>
Why We Did not Get the Picture Executive Legislative & General Legislative Bo
```

Abbildung 9.8: Konkordanz zum vierten Pfad aus dem Graphen date.grf

Der folgende Zweig des Graphen date. grf beschränkt sich auf Daten, in denen nur Monat und Jahr genannt werden, und ganz auf den Tag verzichtet wird.

Dabei werden diese Datumsangaben mit der Präposition "in" eingeleitet, worauf der Monatsname und eventuell im Anschluss eine Jahreszahl folgt.

Außerdem lässt dieses Suchmuster auch relative Jahresangaben zu, so dass nicht nur explizite Jahreszahlen wie z.B. 2006 sondern auch "this year", "last year" usw. erlaubt sind. Diese Eigenschaft des Automaten wird einerseits durch den Knoten mit der

Markierung <DET> und andererseits durch den folgenden Knoten, der das Wort "year" enthält, gewährleistet. Hierbei gleicht der Transduktor jede Instanz von <DET> - eines Determinativs - aus dem entsprechenden Wörterbuch mit dem Text ab.

```
me responsibility for the peace process
id.(S) The current nuclear crisis began ic" about holding six-party Korea talks
the International Atomic Energy Agency
Thica exploded its first nuclear device ed his formal announcement in the state sannually.(S) WASA officials said that

| Ababasia November<| Date> in October of 2002<| Date> after Us officials said the November<| Date> in Ebruary<| Date> in Eventure Legislative & General International Atomic Energy Agency
| Cate>in May 1974<| Date> | S} Executive Legislative & General Gen
```

Abbildung 9.9: Konkordanz zum fünften Pfad aus dem Graphen date.grf

Der sechste Pfad im Graphen date.grf ist nicht nur für die Erkennung von Datumsangaben sondern auch für Zeitspannen zuständig. Da dieser Weg durch den Automaten zwei Aufrufe von date.grf selbst enthält, werden konkrete Daten durch andere Zweige des Transitionsnetzes gefunden, wenn ihnen die Schlüsselworte "from" und "to" vorangehen.



Abbildung 9.10: Sechster Pfad aus dem Graphen date.grf

So lokalisiert dieser Ast des Graphen Textpassagen der Form:

```
agestan, takes up the story: "Overnight from <Date> 30 November<\Date> to <Date> 1 December<\Date> , t

In the latest figures, tourist arrivals from <Date> January 2003<\Date> to <Date> September 2003<\Date>
ipt of applications for IT scholarships from <Date> January 31<\Date> to <Date> March 31<\Date> . Acco

S} The Johor Arts Festival will be held a leading car and truck rental company, from <Date> November 1997<\Date> to <Date> April 4<\Date> . It will

to <Date> February 2000<\Date>
```

Abbildung 9.11: Konkordanz zum sechsten Pfad aus dem Graphen date.grf

Der nächste Pfad im Automaten beschreibt verschiedene Möglichkeiten, wie Monatsangaben im Text zur Spezifikation von Zeitpunkten eingesetzt werden können.

Dabei würde dieser reguläre Ausdruck zu folgenden Zeitangaben passen:

at the end of May at the end of this beautiful spring yesterday tomorrow morning last summer

Durch seine  $\epsilon$ -Transitionen (siehe Abbildung 9.12 auf Seite 101) werden verschiedene Variationen ungefährer Zeitangaben zugelassen.

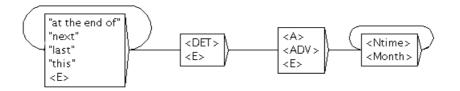


Abbildung 9.12: Siebter Pfad aus dem Graphen date.grf

So können Monatsangaben, sowie relative Tagesangaben wie "tomorrow" oder Zeitangaben wie "morning" und Saisonangaben wie "spring" für Jahreszeiten berücksichtigt werden. All diese Nomina werden über die beiden Kategorien <month> und <mtime> in den entsprechenden Lexika nachgeschlagen und mit dem Korpus verglichen. Welche Einträge das Lexikon Ntime-.dic enthält, wurde bereits in Abschnitt 5.2.10.3 (siehe Seite 73) erläutert.

Dagegen beschäftigt sich der achte Pfad im Graphen date. grf mit der Erkennung von exakten Datumsangaben, welche durch Varianten der Floskel "with effect from" eingeleitet werden.

Unter Variationen dieses feststehenden Ausdrucks versteht man:

with effect from with immediate effect from with backdated effect from

Ähnlich wie schon im vorherigen Ast des Graphen, ruft sich im neunten Pfad der Graph date.grf rekursiv selbst auf. Dieser Weg durch den Automaten unterscheidet sich vom achten Pfad nur darin, dass eine andere Floskel bzw. Präposition dem Datum vorangeht.

Der zehnte Pfad (siehe Abbildung 9.13) beschäftigt sich wieder mit der Erkennung von Zeiträumen. Ähnlich wie beim "from … to …" Zweig werden auch hier zwei Daten erkannt, die im Zusammenhang miteinander stehen. Dabei muss es sich nicht um zwei vollständige Datumsangaben handeln, denn die  $\epsilon$ -Übergänge (<E>) lassen zu, dass beispielsweise vom ersten Datum nur der Tag genannt wird und dafür das zweite komplett im Text erscheint. So ist es möglich, Daten der Form 3/4 February 2005 , sowie 25 Oct – 10 Nov 1993 oder 2nd – 3rd May 2001 im Korpus zu finden, was die Konkordanz aus Abbildung 9.14 bestätigt.

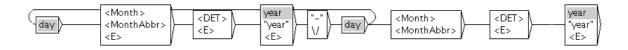


Abbildung 9.13: Zehnter Pfad aus dem Graphen date.grf

Die nun folgenden syntaktischen Variationen zu einer bestimmten Zeitspanne, sollen anschaulich deutlich machen, wie eine  $\epsilon$ -Transition  $\leq$ E> dazu beitragen kann, möglichst viele Möglichkeiten eines "Von-Bis-Zeitraumes" zu erfassen.

1	February		2005	-	5	February		2005
1st	February		2005	-	5th	February		2005
1	Feb		2005	-	5	Feb		2005
1	Feb		05	-	5	Feb		05
1				-	5	Feb		2005
1st				-	5th	Feb.		2005
1st	February	last	year	-	5th	February	next	year
1				/	2	Feb		2005
1				/	2	Feb	this	year

Der Zyklus im Graphen, welcher durch einen Übergang vom Knoten mit der zweiten Referenz auf den Subgraphen day auf den ersten Knoten mit derselben Markierung entsteht, ermöglicht es, auch Daten mit folgender Struktur im Text zu lokalisieren:

```
1/2/3 Feb 2005
1/2/3/4 Feb 2005
```

Auf diese Weise können auch Datumsangaben gefunden werden, welche eigentlich eine Zeitspanne ausdrücken, in denen die einzelnen Tage separat aufgeführt sind.

```
Ex-Date 30 January 2004 Period Covered <a href="Mailto:Annuary 2004"><a href="Mailto:Date"><a href="Mailto:Date">Mailto:Date</a><a href="Mailto:Date"><a href="Mailto:Date">Mailto:Date</a><a href="Mailto:Date">Mailto:Date</a><a href="Mailto:Date">Mailto:Date</a><a href="Mailto:Date">Mailto:Date</a><a href="Mailto:Date">Mailto:Date</a><a href="Mailto:Date">Mailto:Date</a><a href="Mailto:Date">Mailto:Date</a><a href="Mailto:Date</a><a href="Mailto:Date</a><a href="Mailto:Date</a><a href="Mailto:Date</a><a href="Mailto:Date</a><a href="Mailto:Date</a><a href="Mailto:Date</a><a href="Mailto:Date</a><a href="Mailto
  ing Char Chinu District on the night of <a href="Nata">Abate>17/18 January</a> [S] It is worth remembering that th
  ws agency weekly schedule of events for <Date>2 - 8 Feb 04<\Date>
                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                              BBC Monitoring Service - United King
ons this year will last for three days, <a href="Academics.color: blue">Academics.color: Academics.color: blue days, <a href="Academics.color: blue">Academics.color: blue days, <a href="
  ayhem.(S) The drama will be staged from  <Date>22 - 26 April<\Date> .(S) For chamber music enthusiasts,
 pay an official visit to Turkey between  <-Date>-22-24 February 2004<\Date> _.{S} Schroeder is scheduled to
  -2004 German chancellor to visit Turkey <a href="Abel-22-24 February">(Date></a> Executive Legislative & General Ge
n Conference and Exhibition, Singapore, <a href="Abde-22nd-24th March<\Date">Ade-22nd-24th March<\Date</a> .(S) Norwegian Ship Finance Confer
  4 BBC Monitoring Iranian media roundup (<a href="Logo Space">Logo Space 
 Feb-2004 Russian TV highlights for week <a href="Maintain-2004"><a href="Maint
                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                               Television Broadcastin
  tion is on a working visit to Laos from <<u>Oate>27-31 January<\Date></u> at the invitation of the LPRP's Dep
  bruary].{S} Chen, who visited Indonesia <a href="mailto:ADate-29 January - 1 February">ADate-</a> at the invitation of the
         a NATO high delegation to Macedonia on <a href="Action-10-16-bell-was-the-reason"><a href="Action-10-bell-was-the-reason"><a href="Action-10-bell-was-the-was-the-was-the-was-the-was-the-was-the-was-the-was-the-was-the-was-the-was-the-was-the-was-the-was-the-was-the-was-the-was-the-was-the-was-the-was-the-was-the-was-the-was-the-was-the-was-the-was-the-was-the-was-the-was-the-was-the-was-the-was-the-was-the-was-the-was-the-was-the-was-the-was-the-was-the-was-the-was-the-was-the-was-the-was-the-was-the-was-the-was-the-was-the-was-the-was-the-was-the-was-the-was-the-was-the-was-the-was-the-was-the-was-the-was-the-was-the-was-the-was-the-was-the-was-the-was-the-was-the-was-the-was-the-was-the-was-the-was-the-was-the-was-the-was-the-was-the-was-the-was-the-was-the-was-the-was-the-was-the-was-the-was-the-was-the-was-the-was-the-was-the-was-the-was-the-was-the-was-the-was-the-was-the-was-the-was-the-was-the-was-the-was-the-was-the-was-the-was-the-was-the-was-the-was-the-was-the-was-the-was-the-was-the-was-the-was-the-was-the-was-the-was-the-was-the-was-the-was-the-was-the-was-the-was-the-was-the-was-the-was-the-was-the-was-the-was-the-was-the-was-the-was-the-was-the-was-the-was-the-was-the-was-the-was-the-was-the-was-the-was-the-was-the-was-the-was-the-was-the-was-the-was-the-was-the-was-the-was-the-was-the-was-the-was-the-was-the-was-the-was-the-was-the-was-the-was-the-was-the-was-the-was-the-was-the-was-the-was-the-was-the-was-the-was-the-was-the-was-the-was-the-was-the-was-the-was-the-was-the-was-the-was-the-was-the-was-the-was-the-was-the-was-the-was-the-was-the-was-the-was-the-was-the-was-the-was-the-was-the-was-the-was-the-was-the-was-the-was-the-was-the-was-the-was-the-was-the-was-the-was-the-was-the-was-the-was-the-was-the-was-the-was-the-was-the-was-the-was-the-was-the-was-the-was-the-was-the-was-the-was-the-was-the-was-the-was-the-was-the-was-the-was-the-was-the-was-the-was-the-was-the-was-the-was-the-was-the-was-the-was-the-was-the-was-the-was-the-was-the-was-the-was-the-was-the-was-the-was-the-was-the-was-
```

Abbildung 9.14: Konkordanz zum zehnten Pfad aus dem Graphen date. grf

Der letzte und unterste Pfad des Automaten ist für die Erkennung von Jahresangaben zuständig. Dabei werden einzelne Jahre, sowie Jahrzehnte, und auch Jahresspannen von diesem regulären Ausdruck berücksichtigt.

```
hastened.{S} "I was supposed to retire \frac{\langle \text{Date} \rangle \text{in 2005} \langle \text{Date} \rangle}{\langle \text{Date} \rangle \text{in a year} \langle \text{Date} \rangle}," Tancangco told the Inquirer in a phone ,000 ticket-buyers could be entertained \frac{\langle \text{Date} \rangle \text{in a year} \langle \text{Date} \rangle}{\langle \text{Date} \rangle \text{in the 1920s} \langle \text{Date} \rangle}. (S) Allawi's crime was that he was a
```

Abbildung 9.15: Konkordanz zum elften Pfad aus dem Graphen date.grf

Des Weiteren ruft dieser Pfad des Transduktors den Subgraphen year\_long.grf auf, welcher im Gegensatz zu year.grf nur vierstellige Jahreszahlen erkennt. Analog dazu gibt es auch den Subgraphen year\_short.grf, welcher hier nicht explizit zum Einsatz kommt, weil year.grf schon seine Funktion, zweistellige Jahresangaben zu finden, übernimmt. Die entsprechenden Abbildungen zu diesen Graphen befinden sich auf Seite 103 und 104.

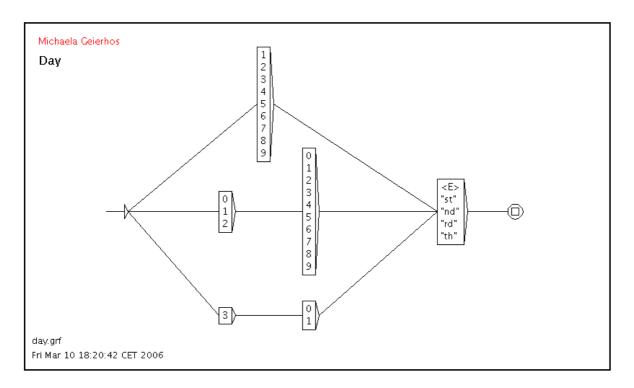


Abbildung 9.16: Graph zur Erkennung von numerischen Tagesangaben - day.grf

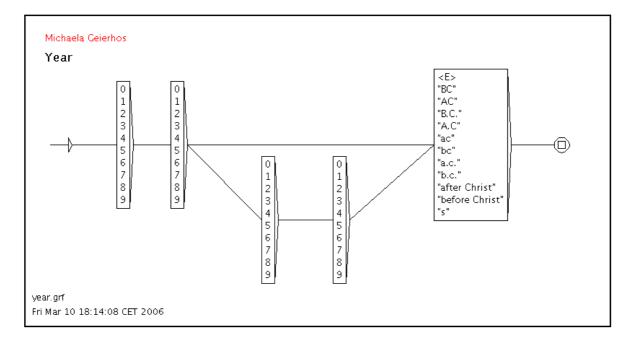


Abbildung 9.17: Graph zur Erkennung von Jahreszahlen - year.grf

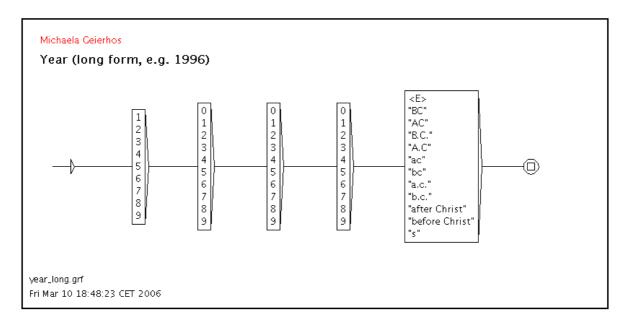


Abbildung 9.18: Graph zur Erkennung von vierstelligen Jahreszahlen - year\_long.grf

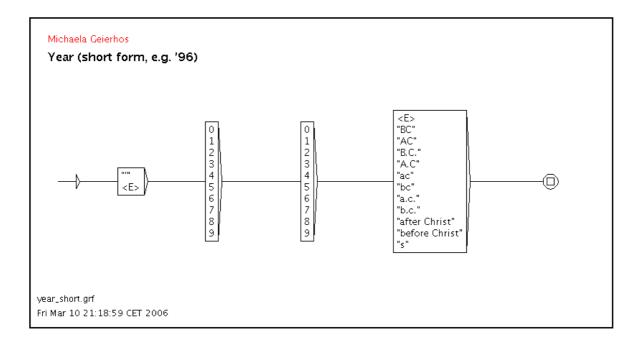


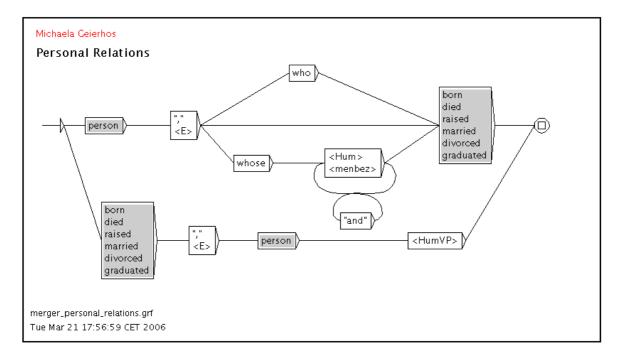
Abbildung 9.19: Graph zur Erkennung von zweistelligen Jahreszahlen - year\_short.grf

# 10 Grammatik persönlicher Relationen

Nachdem in den vorangegangenen Kapiteln die Grundlagen zur Erkennung diverser Entitäten gelegt wurden, können diese nun in Beziehung zueinander gesetzt werden. Wie bereits in Abschnitt 1.2 angesprochen wurde, vermitteln die persönlichen Relationen als eine Untermenge der biographischen Relationen eine Fülle an Informationen über die verschiedensten Leute.

Da es den Rahmen dieser Arbeit übersteigen würde, auf alle englischen Prädikate einzugehen, welche in persönlichen Kontexten auftreten können, werden die wichtigsten Verben stellvertretend ausgewählt. An ihnen soll aufgezeigt werden, wie sich die unterschiedlichen Menschenbezeichner zusammen mit anderen Entitäten syntaktisch und semantisch in die Struktur dieser Relationen einfügen.

Für diesen Zweck wurde eine auf den ersten Blick sehr kompakt erscheinende Grammatik entwickelt, deren Aufgabe es ist, die einzelnen Prädikatrelationen in nicht allzu komplexen Satzgefügen zu lokalisieren.



**Abbildung 10.1:** Graph zur Erkennung von ausgewählten persönlichen Relationen -  $mer-ger\_personal\_relations.grf$ 

Hierbei beschränkt sich der Graph aus Abbildung 10.1 auf zwei Typen von Sequenzen. Einerseits können sie mit einer Nominalphrase beginnen, in der ein Menschenbezeichner vorkommt, auf die anschließend ein Relativsatz folgt, welcher wiederum von "who" oder

"whose" eingeleitet wird. Andererseits ist es auch möglich, dass die Verbalphrase, welche die persönliche Relation beschreibt, am Satzanfang steht und erst danach die betreffende Person zusammen mit einem anderen Prädikat genannt wird.

Um diese Grammatik so übersichtlich wie möglich zu halten, wurde sie sehr modular aufgebaut, so dass jeder einzelnen Relation ein eigener Graph zugewiesen wurde. Dabei decken die Subgraphen nicht nur die Verben ab, welche schon im Namen der Automaten vorkommen, sondern auch noch deren Synonyme.

Auf diese Weise behandelt der Graph merger\_personal\_relations.grf 50 Verbkonstruktionen in verschiedenen Zeitformen und Variationen:

- to be born
- to be born and raised (up)
- to be born and brought up
- to see the light of day
- to grow up
- to be brought up
- to be raised (up)
- to spend one's childhood
- to become a graduate
- to graduate
- to take one's degree
- to receive one's degree
- to get one's degree
- to complete one's studies
- to become man and wife
- to join in marriage
- to marry so.
- to be married to so.
- to get married to so.
- to plight one's troth to so.
- to pledge one's troth to so.

- to lead so. to the altar
- to take so. to wife/husband
- to wed so.
- to be wedded to so.
- to divorce so.
- to be divored from so.
- to file for a divorce from so.
- to sue for a divorce from so.
- to get a divorce from so.
- to part from so.
- to part company with so.
- to separate from so.
- to split from so.
- to split up with so.
- to break up with so.
- to end one's marriage to so.
- to dissolve one's marriage to so.
- to annul one's marriage to so.
- to breath one's last
- to decease
- to depart one's life

- to die (off)
- to expire
- to lay down one's life
- to lose one's life

- to meet one's death
- to meet one's end
- to pass away
- to perish

Damit erfasst dieser Automat die wichtigsten persönlichen Ereignisse im Leben eines Menschen, wozu die Geburt, die Kindheit, der Schulabschluss, die Heirat und eventuelle Scheidung, sowie der Tod gehören.

Des Weiteren besteht auch die Möglichkeit jeden der Subgraphen einzeln im Kontext des Graphen merger\_personal\_relations.grf (siehe Seite 105) aufzurufen, um sich ein Bild von der jeweiligen Relation im Satz zu verschaffen.

Die folgenden Graphen, welche hier nicht abgedruckt wurden, übernehmen diese Aufgabe:

- merger\_born.grf
- merger\_died.grf
- $\bullet \ \ merger\_divorced.grf$

- merger\_graduated.grf
- merger\_married.qrf
- merger\_raised.grf

In den nächsten Abschnitten werden nun die verschiedenen Grammatiken vorgestellt, welche die Verbalphrasen mit den bereits genannten Prädikaten, beschreiben.

## 10.1 Die Geburt: "to be born"

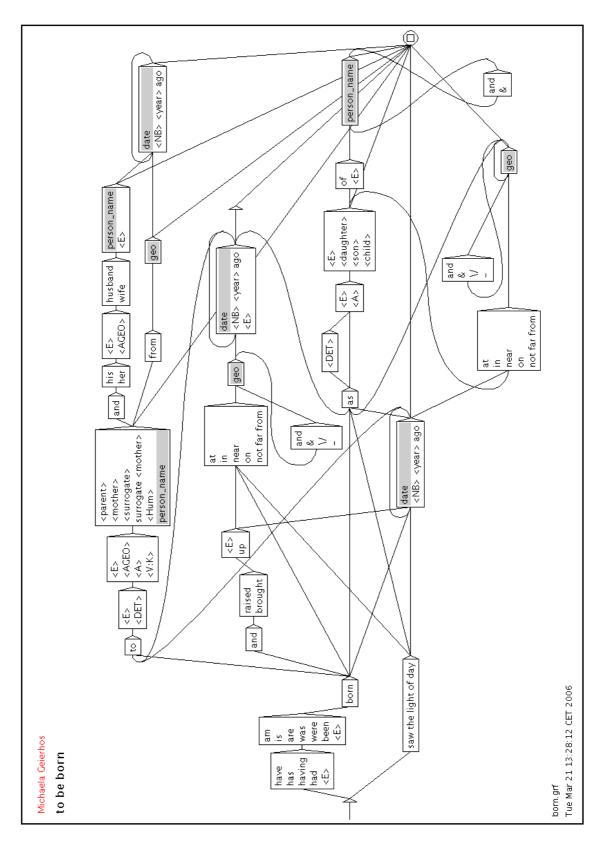
Das Ereignis der Geburt ist wohl das am häufigsten beschriebene Geschehen in Lebensläufen. Meist werden dabei der Geburtstag, der Geburtsort und eventuell noch die Eltern der betreffenden Person erwähnt. Manchmal kann es allerdings auch vorkommen, dass jemand an dem Ort aufgewachsen ist, an dem er geboren wurde. In diesem Fall fällt der Geburtsort mit dem Ort zusammen, wo derjenige seine Kindheit verbracht hat. Somit werden auch die beiden Prädikate "to be born" und "to be raised (up)" miteinander verknüpft.

Die folgende Grammatik in Abbildung 10.2 auf Seite 108 beschreibt die syntaktischen Strukturen der Verben

• to be born

- to be born and brought up
- to be born and raised (up)
- to see the light of day

und berücksichtigt die diversen Kombinationsmöglichkeiten von Ort und Zeit im rechten Kontext der Prädikate.



**Abbildung 10.2:** Graph zur Erkennung von Verbalphrasen mit dem Verb "to be born" und seiner Paraphrasierung "to see the light of day" - born.grf

Außerdem werden mit der Anwendung dieser Grammatik (siehe Seite 108) auf das Korpus einige Fragen beantwortet, welche im Zusammenhang mit der Geburt eines Menschen aufkommen.

1. Wann wurde die betreffende Person geboren?

Tunku Abdul Rahman was born <Date>in 1903<\Date>
Cecil Beaton, who was born <Date>in 1904<\Date>
their first child, Aidan Gering Pollack, born <Date>last week<\Date>

2. Wo wurde die betreffende Person geboren?

Sigman, born in  $\langle GEO \rangle$  Brooklyn $\langle GEO \rangle$  in 1909 the former mayor of Tel Aviv who was born in  $\langle GEO \rangle$ GEO $\rangle$  McQueen, who will be 35 next month, was born in  $\langle GEO \rangle$ the East End of London $\langle GEO \rangle$ 

3. Wer sind die Eltern der betreffenden Person (des Neugeborenen)?

a child was born to <PersonName>Basir<\PersonName> and his Japanese wife
Prince Michael II, was born to an unknown surrogate
twin baby boys born to an American surrogate mother
Jones, who was born to Welsh parents in Papua New Guinea
She was born as daughter of Matheus Klaas and Ida Lysse in Zürich

4. Wie lautet der Geburtsname der betreffenden Person?

1934 BORN IN Poland as <PersonName>Manya Sklodowska<\PersonName> Captain Jack was born as <PersonName>Francisco Gutierrez<\PersonName> in Havana/Cuba
Martika was born as <PersonName>Marta Marrero<\PersonName> in Whittier, CA, on
May 18, 1969

Jedoch die eigentliche Frage nach dem Namen der geborenen Person wird nicht von dieser Grammatik, sondern vom Graphen  $merger\_born.grf$  übernommen, welcher somit indirekt auch die eben genannten Fragestellungen beantwortet.

So würde nun dieser Graph Phrasen folgenden Typs beachten und darin unter anderem markieren, wer geboren wurde.

Born and raised in Flushing, Queens, <PersonName>Woodbridge<\PersonName> had attended the School of Visual Arts <PersonName>Lycia Danielle Trouton<\PersonName>, who was born in Belfast <PersonName>James Saunders<\PersonName> was born in London in 1925

Eine Konkordanz mit sämtlichen Annotationen des Transduktors born.grf, welcher nur Verbalphrasen mit Formen des Verbs "to be born" berücksichtigt, ist auf der nächsten Seite zu finden.

was born <Date> on Aug. 20, 1939</Date> as an illegitimate son of <PersonName>Allan</PersonName> was born in <GEO> Brazil</GEO> and has his family roots in Lebanon. {S} The cost consultants, D .{S} 28-Feb-2004 28-Feb-2004 Be was born to Basir and his Japanese wife <Date> last December<\Date> , and he intends to file f born to an American surrogate mother, a ministry official said Friday, raising the prospect tha born to smokers with the claim that 'Some women would prefer having smaller babies'. (S) However mpany News Sales 582 GUILIN OF MALAYSIA HAVING been bern and raised in <GEO> Ipoh</GEO> , I am very proud of my hometown and being asso married with one son, Abdul Karim, who was born <Date> in June 2002<\Date> . (S) He is more likely to be at home, he says, "reading bo to Fernando Pou, a Filipino, and Bessie Kelly, an Ame was born in <GEO> Birkenhead<\GEO> and even turned out for Tranmere as a substitute in the LDV was born in <GEO> Dartford, Kent. < \GEO> {S} In my view there are no other reasons for dismissin and took a degree in Chemistry at South was born in <GEO> the East End of London<\GEO> , the youngest of six children, and attended Rok was born to an unknown surrogate). [S] If you were going from A to B, and A was Harvey Nicks and born in <GEO> the Ukraine<\GEO> , is joined by Greg Anderson, the personal trainer of the Ameri , the son of settlers who had nearly died crossing the Sierra, born and raised in <GEO> Chile<\GEO> , heads Artemis Wines International and will oversee all U <pr , Carey moved to Ne was born in <GEO> Belfast<\GEO> and grew up in Canada, has recently come to Australia to take was born in <GEO> Derry</GEO> and spent her early childhood there.{S} It was given to Glasgow has been born in <GEO> Sydney<\GEO> .(S) BBC Monitoring Service - United Kingdom International was born <Date> 7 October 1952<\Date> .{S} "This deal allows Elcoteq to be even more vigorous ates the centenary of Cecil Beaton, who was born <Date> in 1904<\Date> , started taking photographs aged 11, and died in 1980.{S} The BORN IN <GEO> Poland</GEO> as <PersonName> Manya Sklodowska</PersonName> , the scientist moved born in <GEO> Germany</GEO> to parents from <GEO> Russia</GEO> .(S} "This is really a wait and great operas: Born and raised in <GEO> Flushing, Queens, Woodbridge<\GEO> had attended the School of Visual ssell Ellis - Hunter was a stage name - was born <Date> in 1925<\Date> in <GEO> Glasgow<\GEO> but shortly afterwards was sent to liv born in <GEO> St Austell<\GEO> 37 years ago, is in select company and even at this late stage born and brought up in <GEO> Belfast<>GEO> , film-maker Maeve Murphy remembers being puzzled , is best known as a composer of was born to <PersonName> Basir<\PersonName> and his Japanese wife <Date> last Born in <GEO> Long Island, New York<\GEO> <Date> on March 27, 1970<\Date> <Date> on July 17, 1947<\Date> was born in <GEO> London<\GEO> <Date> in 1925<\Date> was born 150 years ago <Date> this year<\Date> was born <Date> on Aug. 20, 1939<\Date> near <GEO> Sacramento</GEO> was born in <GEO> London<\GEO> Born in <GEO> Iowa</GEO> born ight suggest otherwise. {S} Janacek, who all the candidates only Vladimir Putin o register the births of twin baby boys dismissed claims that a cloned baby boy ro noted the "undisputed fact" that Poe atural-born Filipino citizen because he Lycia Danielle Trouton, who es well at this level. (S) Andy Robinson the company's chairman in waiting - who s christened Graham Leaver, and that he rge Bernard Shaw's wife, Charlotte, who hat he wanted to say. (S) James Saunders der McQueen, who will be 35 next month, reporters after the ruling that a child star's third child, Prince Michael II, reporters after the ruling that a child about the lower birth-weight of babies ' Irish Independent 861 \* Camilla Shand d blood vessels. (S) As a young teenager inner who discovered radium 1867 - 1934 to elite athletes. (S) The 71-year-old, 19th-century philosopher Josiah Royce, (S) Artemis Wines president Eric Smith, died after being bitten by a dog. {S} ation in Berlin," said Lana Bruschtina, Geology and Mining Industry in 1984. (S) nerian roles in the 1920s and 1930s. {S} t was a few years ago. {S} Nigel Martyn, is week, {S}

Abbildung 10.3: Konkordanz zum Graphen born.grf

# 10.2 Die Kindheit: "to be raised (up)"

Wie bereits im vorherigen Abschnitt angesprochen wurde, ist das Aufwachsen manchmal im unmittelbaren Kontext des Geburtsortes zu suchen. Doch kann es natürlich auch als eigenständiges Geschehen auftreten und muss somit durch eine gesonderte Grammatik behandelt werden.

Der Automat, welcher nun dafür zuständig ist, herauszufinden, wo jemand seine Kindheit verbracht hat, wie derjenige aufgezogen wurde und in welchem Umfeld er aufgewachsen ist, beschäftigt sich mit den folgenden Verbkonstruktionen:

• to grow up

• to be raised (up)

• to be brought up

• to spend one's childhood

Des Weiteren berücksichtigt dieser Graph (siehe Abbildung 10.4 auf Seite 112) die gebräuchlichsten englischen Synonyme für den Begriff der "Kindheit". An dieser Stelle macht er sich folgendes Vokabular zunutze:

adolescence, babyhood, boyhood, childhood, early days, early life, early years, girlhood, immaturity, infancy, infanthood, pre-teens, prepubescence, teenage years, teens, young adulthood, youth

Wenn es darum geht, die Art oder den Ort der Erziehung näher zu beschreiben, stützt er sich auf folgende Ausdrücke, welche zuvor mittels Bootstrapping auf den Korpora ermittelt wurden.

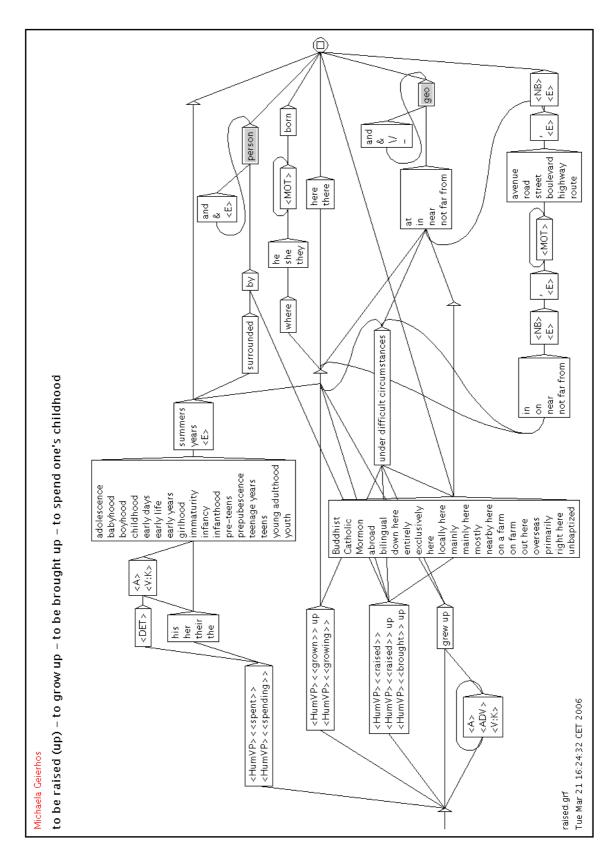
Buddhist, Catholic, Mormon, abroad, bilingual, down here, entirely, exclusively, here, locally here, mainly, mainly here, mostly, nearby here, on a farm, on farm, out here, overseas, primarily, right here, unbaptized

Dagegen werden die Umstände, unter denen jemand aufgewachsen ist, meist mit der Floskel "under difficult circumstances" abgetan, und die Familie oder die Personen, welche in der Kindheit eine maßgebliche Rolle gespielt haben, werden häufig mit den Worten "surrounded by" eingeleitet.

Der Ort, an dem jemand aufgewachsen ist, kann allerdings auch sehr detailliert wiedergegeben werden, indem sogar der Straßenname genannt wird. Manchmal lässt sich auch wieder ein Bezug zum Geburtsort der betreffenden Person herstellen. Dafür wird meist die Phrase "where so. was born" verwendet. Zudem lässt die Grammatik noch folgende englische Synonyme für den Ausdruck "Straße" zu:

avenue, road, street, boulevard, highway, route

Welche Verbalphrasen von diesem Graphen beispielsweise im FT Korpus gefunden werden, illustriert die Konkordanz in Abbildung 10.5 auf Seite 113. Dagegen findet der Graph  $merger\_raised.grf$  (ohne Abbildung) wieder heraus, wessen Kindheit und Aufwachsen hier beschrieben wird.



**Abbildung 10.4:** Graph zur Erkennung von Verbalphrasen mit dem Verb "to be raised (up)" und seiner direkten Synonyme - raised.grf

```
113
```

```
two young Kenyans born and brought up in <GEO>Kenya<\GEO>'s coffee growing areas.{S}
 the magazine publisher James Brown, who grew up in <GEO>Leeds<\GEO> and has supported the team ever since.{S}
  The 47-year old actress, who grew up near <GEO>Londonderry<\GEO> and plays forensic pathologist Dr Sam Ryan
                  has been raised by <Person>Mrs <PersonName>Pat O'Dwyer<\PersonName><\Person>
 have also been raised by <Person><JD>D.C. Council member<\JD><PersonName>Phil Mendelson<\PersonName><\Person>
      A woman brought up in <GEO>Australia<\GEO> who became a heroine of the French resistance has finally
        As a young teenager born and brought up in <GEO>Belfast<\GEO>, film-maker Maeve Murphy remembers
Rosemary Canavan - born in Scotland, brought up in <GEO>Northern Ireland<\GEO>, with an English, Welsh, Irish and
                                               Huguenot background
            says Matt Hollander, who grew up in <GEO> Belfast<\GEO>, but lived in London for a while
 Patrick Devlin grew up in <GEO>Ireland<\GEO> and believes that this is the most fun band he is ever heard.{S}
Situell took over Weston Hall, where he had spent his childhood, and was a much-loved member of the village. (S)
          Having spent her childhood in 48 Cherryfield Avenue, Ranelagh, just a few hundred yards from
                       Kevin Lewis spent his childhood being beaten by his mother at home
  Nigel, who was born in 1958, spent his early years at the Mamre Brook house at the Saltram winery in Barossa
        Born and raised in <GEO>Flushing, Queens<\GEO>, Woodbridge had attended the School of Visual Arts
                A Jewish musician raised in <GEO>Israel<\GEO> and a Palestinian intellectual who
                    Charlotte, who was born in Derry and spent her early childhood there. {S}
Shirley Morgan (now Shirley Larsen) was being brought up by her white <Person>parents<\Person>, <Person> Jim and
                                                Jean<\Person>.{S}
             I was brought up by <Person><PersonName>Bill Nicholson<\PersonName><\Person> at Spurs
                                I was brought up in <GEO> north London<\GEO>.{S}
```

Abbildung 10.5: Konkordanz zum Graphen raised.grf

the house where Katherine Mansfield spent her early years is closed only one day of the year. {S}

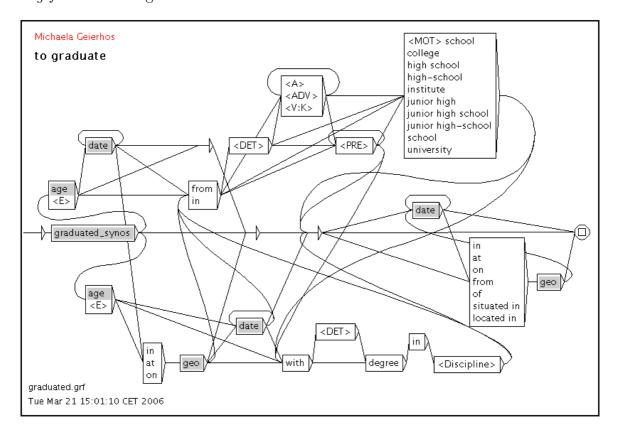
### 10.3 Der Schulabschluss: "to graduate"

Ein weiterer "Meilenstein" im Lebenslauf kann der Schulabschluss sein. Für das Ereignis, erfolgreich von einer Schule abzugehen, werden folgende Verbkonstruktionen näher untersucht:

- to become a graduate
- to graduate
- to take one's degree

- to receive one's degree
- to get one's degree
- to complete one's studies

Diese werden zusammen mit ihren jeweiligen rechten Kontexten im Graphen graduated.grf in Abbildung 10.6 detailliert beschrieben.



**Abbildung 10.6:** Graph zur Erkennung von Verbalphrasen mit dem Verb "to graduate" und seiner direkten Synonyme - graduated.grf

Dabei geht diese Grammatik einerseits auf die Schule bzw. den Namen der Lehreinrichtung ein, an welcher der jeweilige Abschluss erzielt wurde und versucht mit Hilfe der Schlüsselwörter

college, high school, high-school, institute, junior high, junior high school, junior high-school, school, university

vorangehende groß geschriebene Wörter als den Schulnamen zu identifizieren.

Des Weiteren ist oft interessant, mit welchem Alter jemand von der Schule abgegangen ist. Um diesen Faktor in die Grammatik aus Abbildung 10.6 einzubinden, wurde ein eigener Graph zur Lokalisierung von Altersangaben entwickelt.

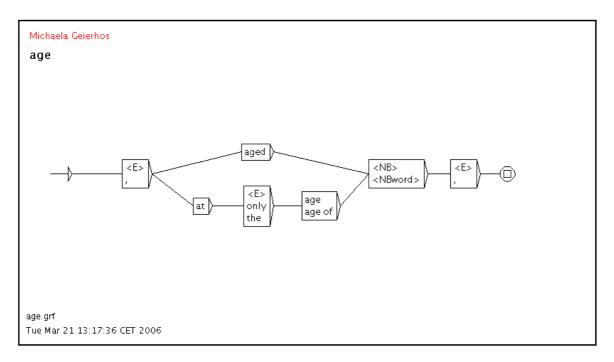


Abbildung 10.7: Graph zur Erkennung von Altersangaben - age.grf

So würde dieser Automat in Abbildung 10.7 z.B. folgende Altersangaben erkennen:

```
at the age of 18
aged 19
at age of 20
```

Dagegen werden die möglichen Verbkonstruktionen für den Ausdruck "einen Abschluss erlangen" vom Graphen graduated\_synos.grf beschrieben.

Doch bevor der Aufbau dieses Subgraphen genauer betrachtet wird, sollte zunächst einmal eine Beispielkonkordanz zum Hauptgraphen graduated.grf einen Einblick in die Art der Treffer geben.

```
but will return home to complete his studies.{S}

tudents offered places in MRSM did not complete their studies.{S}

Andrew Gilligan had graduated from Cambridge with a degree in history.{S}

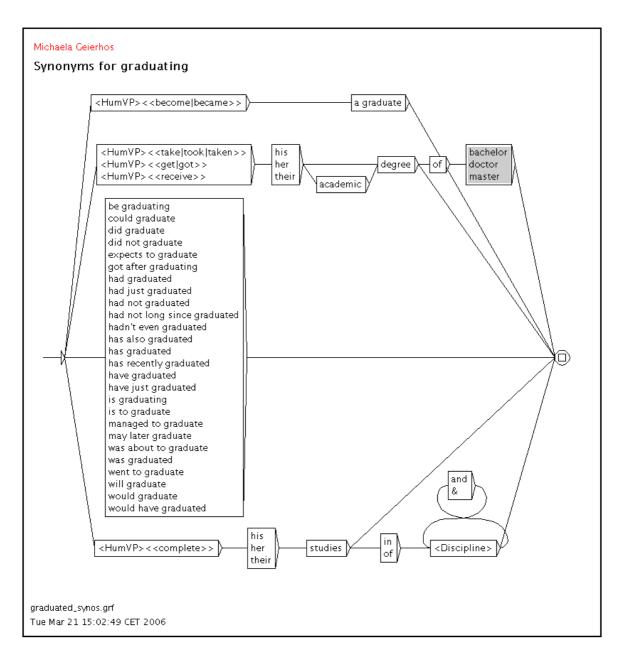
Conal Keaney who has also graduated from the county's All-Ireland winning under-21s but those who have graduated are putting their talent and skill to good use

More than 51,000 teachers have graduated from the Enrique Jose Varona Higher Teaching

Institute, which is celebrating its 40th anniversary this year.{S} He will receive his degree <DATE>in June<\DATE> in University College Cork along with Kerry footballer Mick O'Connell
```

An dieser Konkordanz wird deutlich, dass Informationen über die Art des Schulabschlusses, sowie die Fachrichtungen und der Name der Lehreinrichtung im direkten Umfeld von Verben mit der Bedeutung von "to graduate" zu finden sind.

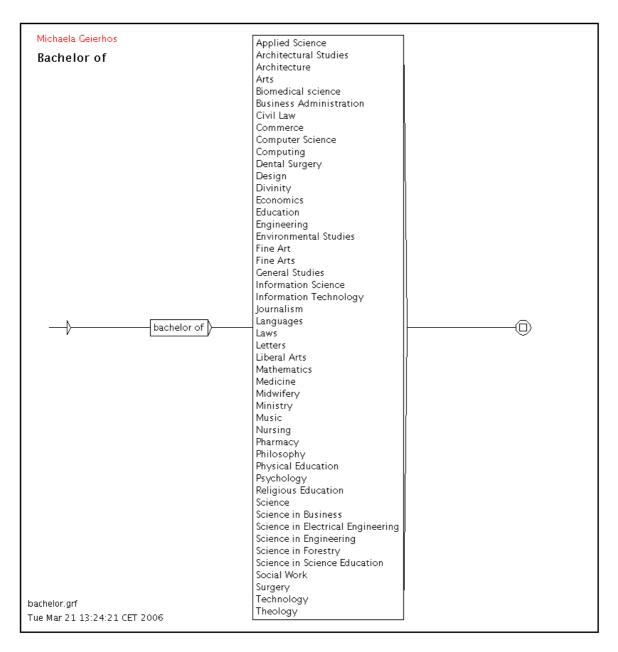
Da die Erkennung von Schulen bereits in der Grammatik graduated.grf (siehe Abbildung 10.6 auf Seite 114) behandelt wird, werden die unmittelbar an das Verb anschließenden Phrasen, welche einen Abschluss charakterisieren, im Subgraph graduated\_synos.grf aus Abbildung 10.8 beschrieben.



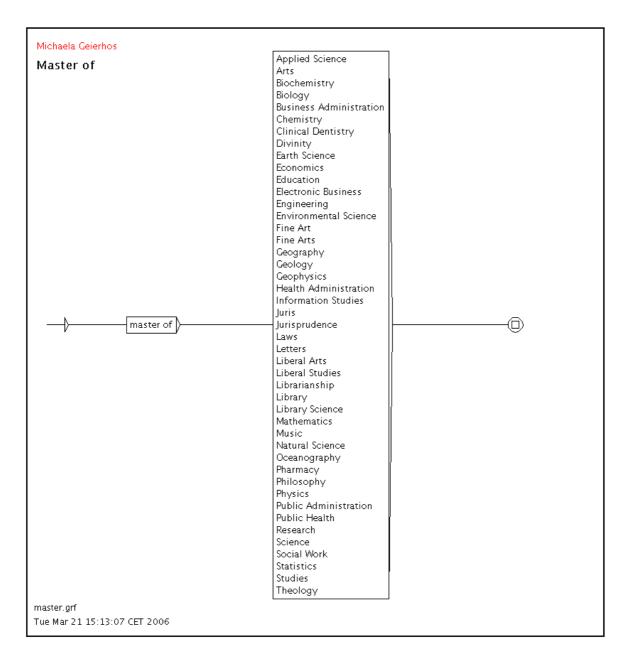
**Abbildung 10.8:** Graph zur Erkennung von Verben mit der Bedeutung von "to graduate" - graduated\_synos.grf

Um die Erlangung eines akademischen Grades zu spezifizieren, wurden drei Graphen entwickelt, welche die gängigsten Bezeichnungen für den Bachelor, den Master und den Doktor abdecken. Diese Graphen werden auf den folgenden Seiten 117-119 dargestellt.

Im Gegensatz dazu wird ein Abschluss in einem bestimmten Fach über das Symbol <Discipline> im entsprechenden Lexikon nachgeschlagen (vgl. Abschnitt 5.2.7.1 auf Seite 63).

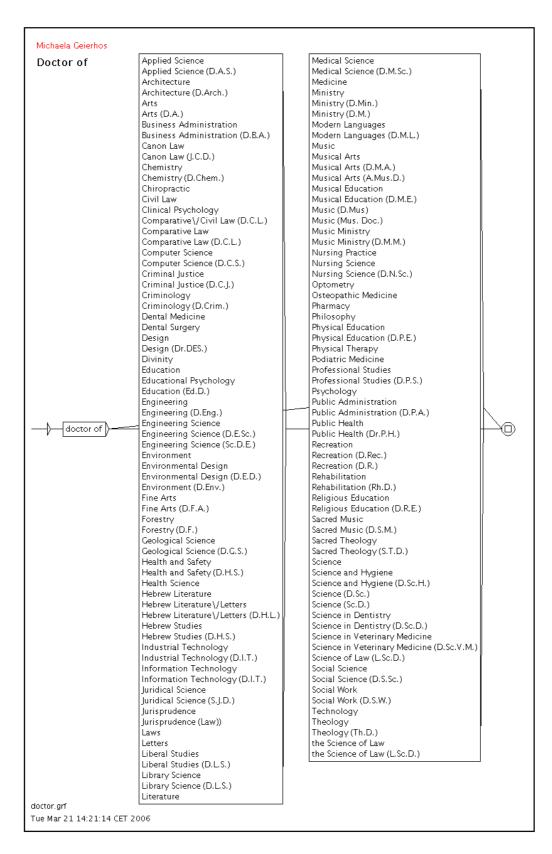


**Abbildung 10.9:** Graph zur Erkennung von Abschlussbezeichnungen für den akademischen Grad des Bachelors - bachelor.grf



**Abbildung 10.10:** Graph zur Erkennung von Abschlussbezeichnungen für den akademischen Grad des Masters - master.grf

Natürlich ist das Vokabular des Automaten, welcher Bachelorabschlüsse in bestimmten Fächern erkennt, sehr ähnlich zu dem des Masterabschlussgraphen. Jedoch wurde darauf Wert gelegt, dass die darin genannten Abschlussmöglichkeiten wirklich existieren, und die in den Transduktoren enthaltenen Begriffe aus fundierten Quellen stammen (wie z.B. der Wikipedia [97]). Zudem ist es für eine mögliche Weiterentwicklung der Graphen ratsam, den Bachelor- und den Master-Graphen getrennt zu halten, falls neue Abschlussfächer hinzukommen sollten, in welchen man beispielsweise den akademischen Grad des Bachelors, aber nicht den des Masters erreichen kann.

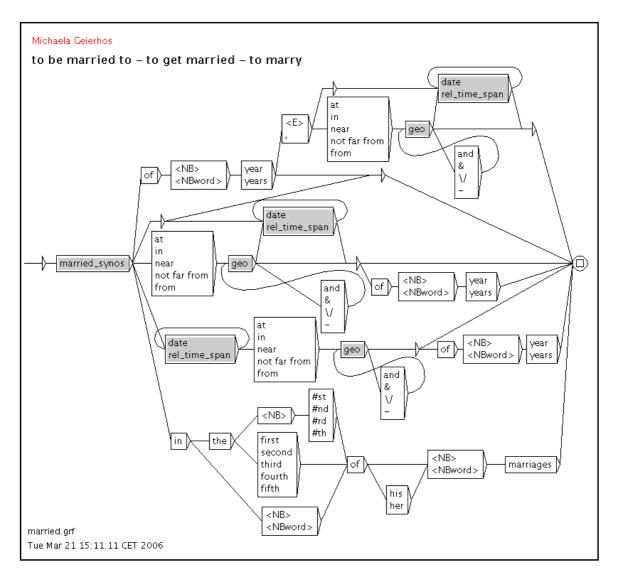


**Abbildung 10.11:** Graph zur Erkennung von Abschlussbezeichnungen für den akademischen Grad des Doktors - doctor.grf

### 10.4 Die Heirat: "to be married"

Ein weiteres Ereignis im Leben eines Menschen kann die Eheschließung sein. Aus linguistischer Sicht, ist sie für die Analyse von Personenbezeichnungen in biographischen Kontexten wesentlich interessanter als die bereits genannten Prädikate. Das liegt vorallem daran, dass eine Heirat stets zwei Personen betrifft und bei einer Aktiväußerung (X heiratet Y) beide im Satz vorkommen müssen. Zwar ist es bei dem passiven Ausdruck des "Verheiratet Seins" nicht unbedingt notwendig, dass der Ehepartner genannt wird, doch ist dies meist der Fall.

Sowohl Aktiv- wie auch Passivkonstruktionen wurden für die Entwicklung einer Grammatik (siehe Abbildungg 10.12) bedacht, deren Ziel es ist, verheiratete Personen im Text aufzuspüren, sowie gewisse Informationen über diese Ehe herauszufiltern. Dabei wurden



**Abbildung 10.12:** Graph zur Erkennung von Verbalphrasen mit dem Verb "to marry so." in seiner Aktiv- und Passivform, sowie seiner direkten Synonyme - married.grf

die folgenden Verben zum Thema "Heiraten" oder "Verheiratet Sein" ausgewählt:

- to become man and wife
- to join in marriage
- to marry so.
- to be married to so.
- to get married to so.
- to plight one's troth to so.

- to pledge one's troth to so.
- to lead so. to the altar
- to take so. to wife/husband
- to wed so.
- to be wedded to so.

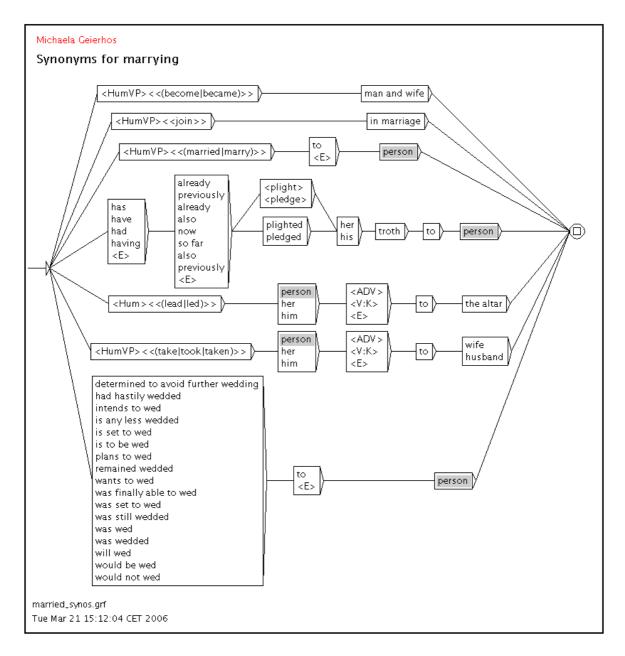
Die verschiedenen Verbformen dieser Prädikate werden zusammen mit ihren Argumenten im Graphen married\_synos.grf in Abbildung 10.13 auf Seite 122 behandelt.

Dagegen befasst sich der Hauptgraph married.grf mit der Spezifizierung von folgenden Angaben bezüglich einer Ehe:

- 1. Wer wurde geheiratet?
- 2. <u>Wann</u> wurde die Ehe geschlossen?
- 3. Wo fand die Hochzeit statt?
- 4. Wie lange liegt die Eheschließung schon <u>zurück</u>?
- 5. Wie viele Monate oder Jahre waren sie verheiratet?
- 6. Um die wievielte Ehe handelt es sich?

Allerdings wird die Frage "Wer ist <u>mit wem</u> verheiratet?" wieder vom Transduktor merger\_married.grf (ohne Abbildung) geklärt.

Die folgende Konkordanz zeigt, welche Phrasen beispielsweise von dieser Grammatik erkannt werden.



**Abbildung 10.13:** Graph zur Erkennung von Verben mit der Bedeutung von "to marry, to be married" - married\_synos.grf

Der Automat married.grf greift als erster von den hier vorgestellten Transduktoren auf den Graphen rel\_time\_span.grf zu, dessen Abbildung auf Seite 123 zu finden ist. So würde sich dieser Subgraph auf das Identifizieren von relativen Zeitspannen konzentrieren. Das heißt nichts anderes, als dass Zeiträume im Korpus gesucht werden, die auf einen Zeitpunkt im Text Bezug nehmen, der eventuell das Erscheinungsdatum des Artikels ist, das akutelle Jahr oder ein Datum, dass einige Passagen zuvor genannt wurde. Auf diese Weise würden auch Ausdrücke wie 5 years ago, 2 month later, days ahead, December ago, etc. von der Grammatik erfasst werden.

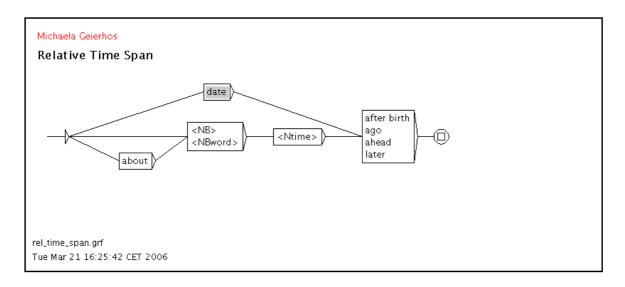


Abbildung 10.14: Graph zur Erkennung von relativen Zeiträumen - rel\_time\_span.grf

### 10.5 Die Scheidung: "to be divorced"

Eine Heirat ist leider auch die Voraussetzung für eine Scheidung. Auch syntaktisch gesehen ähneln diese beiden Antonyme sehr, was die Grammatik in Abbildung 10.15 auf Seite 124 sehr deutlich aufzeigt.

So beschäftigt sie sich unter anderem mit der Feststellung nach wie vielen Ehejahren eine Beziehung geschieden wurde, wann dies geschehen ist, und eventuell noch wo die Scheidung stattgefunden hat.

Eine Konkordanz zum Graphen divorced.grf könnte beispielsweise wie folgt aussehen:

```
three months after Cristina had broken up with <Person><PersonName>

Phillip<\PersonName><\Person>

Today she is divorced and hoping to marry Yitzhak Rabin's assassin bass player Dave Pegg is getting divorced

John Hughes, 51, is divorced with three children.{S}

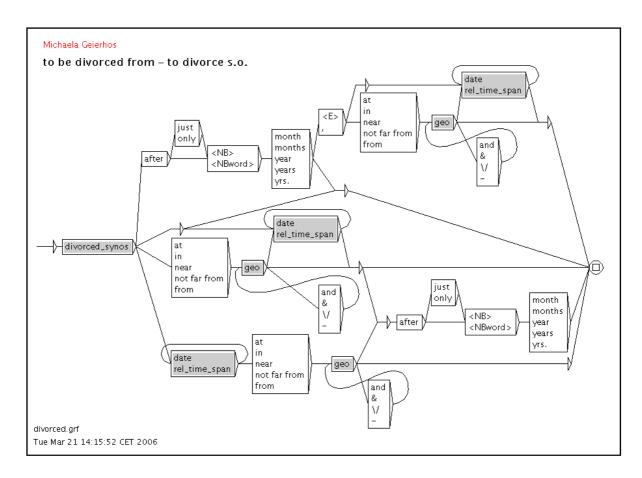
The Ohio congressman, who has been divorced twice

11 days after having been divorced, she married Captain Thomas George Symonds

Babb
```

Um jedoch zu vermeiden, dass sich eine Firma von einem ihrer Mitarbeiter "scheiden" lässt, wie es z.B. in folgendem Satz der Fall ist,

kommt der endliche Automat merger\_divorced.grf zum Einsatz, welcher sicherstellt, dass nur Verbalphrasen mit menschlichen Nominalphrasen im Subjekt gefunden werden.

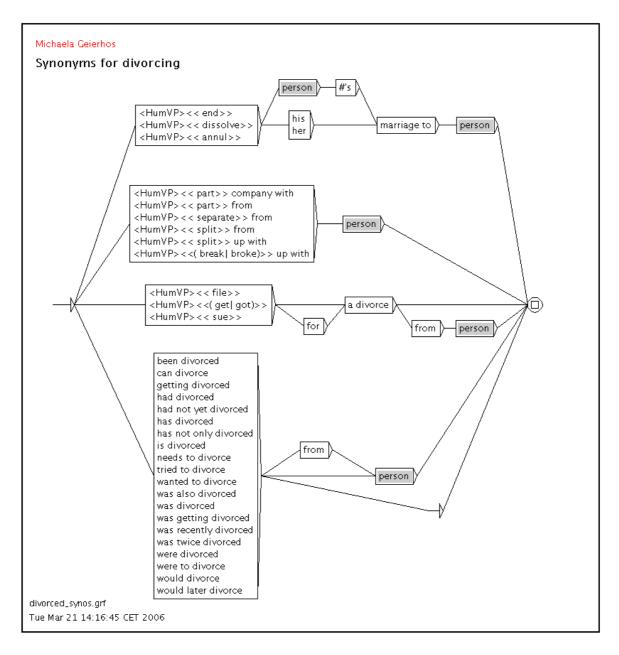


**Abbildung 10.15:** Graph zur Erkennung von Verbalphrasen mit dem Verb to be divorced und seiner direkten Synonyme - divorced.grf

Insgesamt werden 14 Verbkonstruktionen mittels dieser Grammatik im Korpus gefunden, welche im Subgraphen divorced\_synos.grf (siehe Abbildung 10.16 auf Seite 125) ausführlich beschrieben werden. Dabei wird sowohl auf die Aktiv- als auch auf die Passivverwendung dieser Prädikate eingegangen.

- to divorce so.
- to be divored from so.
- to file for a divorce from so.
- to sue for a divorce from so.
- to get a divorce from so.
- to part from so.
- to part company with so.

- to separate from so.
- to split from so.
- to split up with so.
- to break up with so.
- to end one's marriage to so.
- to dissolve one's marriage to so.
- to annul one's marriage to so.



**Abbildung 10.16:** Graph zur Erkennung von Verben mit der Bedeutung von "to be divorced" - divorced\_synos.grf

Überdies ist der Begriff der Synonymie hier relativ weit gefasst. Zwar entsprechen die meisten dieser Ausdrücke in ihrer Bedeutung einer Scheidung, doch können manche von ihnen auch nur eine Trennung oder das Ende einer Beziehung ohne Trauschein darstellen. Obwohl natürlich der Begriff der "Scheidung" auch den Aspekt des "getrennt Lebens" beinhaltet, sagt er doch noch etwas über die rechtliche Grundlage einer Partnerschaft aus. Doch dieser Faktor muss hier nicht berücksichtigt werden, und die syntaktische und semantische Ähnlichkeit der Verben reicht aus, sie in einem Graphen zusammenzufassen.

### 10.6 Der Tod: "to die"

Am Ende dieser Reihe von ausgewählten persönlichen Relationen stehen die Verben, welche den Tod eines Menschen in Worte fassen.

- to breath one's last
- to decease
- to depart one's life
- to die (off)
- to expire
- to lay down one's life

- to lose one's life
- to meet one's death
- to meet one's end
- to pass away
- to perish

Ausgehend von der folgenden Konkordanz, welche hauptsächlich nur darüber Aufschluss gibt, wann eine Person verstorben ist, oder wie viele Jahre ihr Tod schon zurückliegt, soll eine umfassende Grammatik entwickelt werden, welche weitere Informationen über den Tod des jeweiligen Menschen herausfindet.

```
Derek Jarman, who died 10 years ago <Date> today<\Date> Prepackaged Software Motion: least because he died 16 years ago, suffering a heart attack after falling off a hors Michael De-la-Noy died <Date> 12 August 2002<\Date> BY MICHAEL DE-LA-NOY Bloomsbury: n June 6th, 1923; died <Date> 9th February<\Date> , 2004 Europe Ireland Western Europe fashion district, died <Date> Feb. 1,<\Date> it was reported in London.{S} Arts Entering abortion law, died <Date> Feb. 11<\Date> after a battle with prostate cancer, faster a battle
```

Abbildung 10.17: Konkordanz zum Graphen died. grf

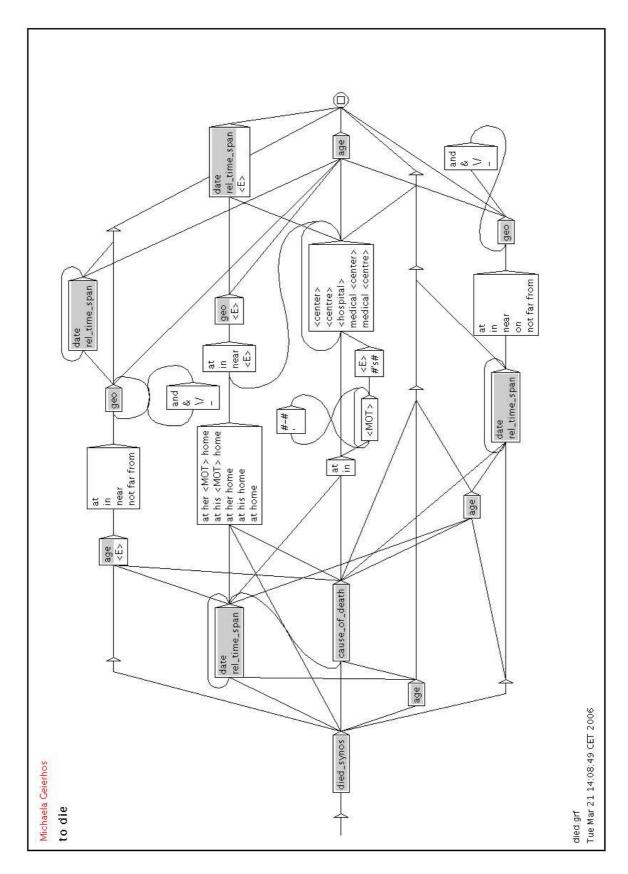
Doch gibt es noch andere Ergänzungen, welche im Umfeld des Ereignisses "Sterben" auftreten. Diese Fülle an Fakten, welche mit dem Ableben einer Person einhergeht, versucht der Graph auf Seite 127 in den Griff zu bekommen.

Dabei sollte das Alter beim jeweiligen Todeszeitpunkt nicht außer Acht gelassen werden, welches meist in Phrasen des folgenden Typs ausgedrückt wird:

```
Prunella Clough died <Date>in 1999<\Date> at the age of 80
and died <Date>last month<\Date> at age 83

Michael Dixon, who has died aged 71, was one of the Financial Times's longest-serving columnists

Adolf Mahr died in <GEO>Bonn<\GEO> in 1951, aged 64.{S}
```



 ${f Abbildung}$  10.18: Graph zur Erkennung von Verbalphrasen mit dem Verbto die - died.grf

Ein weiterer - nicht unbedeutender - Faktor ist die Todesursache.

Diese kann durchaus vielfältig sein, weil sie sich von diversen tödlichen Krankheiten und Drogenmissbrauch, über die unterschiedlichsten Unfälle mit Todesfolge, bis hin zu simplem Herzversagen erstrecken kann.

Aus diesem Grund wurde ein eigener Graph erstellt, welcher nur die Aufgabe hat, Todesursachen im Korpus zu erkennen (siehe Abbildung 10.19).

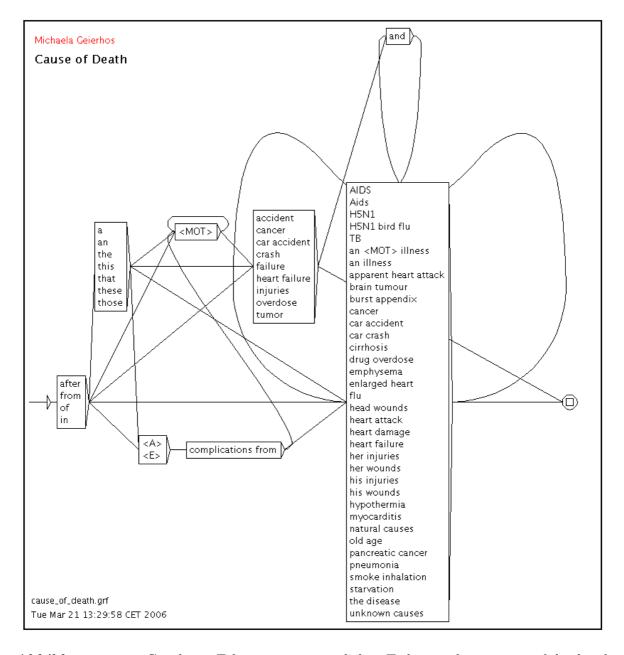


Abbildung 10.19: Graph zur Erkennung von möglichen Todesursachen - cause\_of\_death.grf

So würde der endliche Automat  $cause\_of\_death.grf$  beispielsweise folgende Todesursachen im Text aufspüren:

```
Thomas Hickey died in a <u>cycling accident</u>
Milt Bernhart, who has died of <u>heart failure</u> aged 77
Kristen Pfaff, had died of a <u>heroin overdose</u>
```

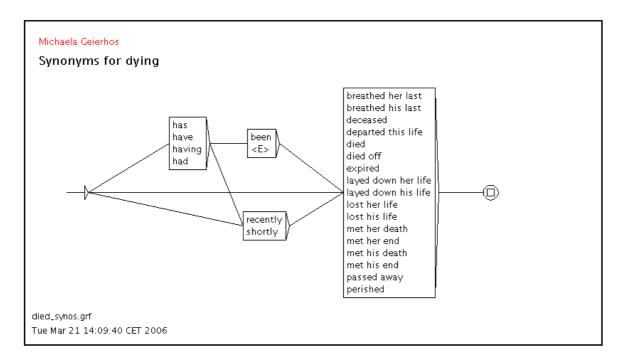
Zudem müssen in der Grammatik died. grf (siehe Seite 127) die verschiedenen Möglichkeiten erwähnt werden, welche ausdrücken, wo ein Mensch verstorben ist. Dabei sollte sie sich nicht nur auf Ortsangaben wie Städte oder Länder beschränken, sondern auch in Betracht ziehen, dass eine Person zu Hause oder im Krankenhaus sterben kann und dann keine genauere Ortsbestimmung darauf folgt.

```
Shirley Strickland de la Hunty has died at her <GEO>Perth<\GEO> home aged 78 Fischer died <Date>on Sunday<\Date> in a hospital in <GEO>Lugano<\GEO>
```

Der Subgraph died\_synos.grf in Abbildung 10.20 ist das Kernstück des Hauptgraphen died.grf, denn er übernimmt die Beschreibung der einzelnen Verbkonstruktionen. Alle Prädikate, die in ihrer Bedeutung dem Verb "to die" entsprechen, kommen im Vokabular dieses Automaten vor.

So erfasst dieser Transduktor auch folgende Verbalphrasen:

Colonel Ryszard Kuklinski, 74, passed away in a <GEO>Washington<\GEO> hospital Warren Zimmermann, who died of pancreatic cancer <DATE>last Tuesday<\DATE>

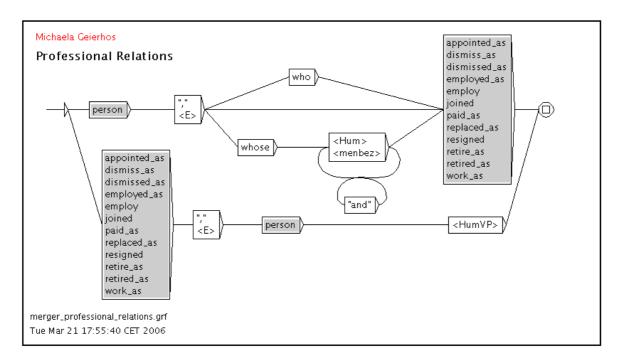


**Abbildung 10.20:** Graph zur Erkennung des Verbs  $to\ die$  und seiner direkten Synonyme -  $died\_synos.grf$ 

## 11 Grammatik beruflicher Relationen

Wie bereits zu Anfang in Abschnitt 1.2.2 angesprochen wurde, liegt der Schwerpunkt der hier vorgestellten linguistischen Untersuchungen verstärkt auf öffentlichen Relationen. Dabei stellen die beruflichen Relationen wohl die größte Untermenge der öffentlichen Beziehungen dar, denn die Möglichkeiten, ein Arbeitsverhältnis kombiniert mit einer Berufsbezeichnung wiederzugeben, sind vielfältig.

So versucht der endliche Automat aus Abbildung 11.1 ähnlich wie schon der Graph auf Seite 105 die einzelnen Grammatiken für ausgewählte Verbalphrasen auf recht kompakte Weise darzustellen. Auch hier werden Sätze mit derselben Struktur wie in der Grammatik für die persönlichen Relationen beschrieben.



**Abbildung 11.1:** Graph zur Erkennung von ausgewählten beruflichen Relationen -  $mer-ger\_professional\_relations.grf$ 

Zudem wurde sehr auf die Modularität dieses Transitionsnetzes geachtet, so dass jeder einzelnen Relation ein eigener Graph zugewiesen wurde. Dabei decken die Subgraphen nicht nur die Verben ab, welche schon im Namen der Automaten vorkommen, sondern auch noch deren Synonyme.

Auf diese Weise behandelt der Graph merger\_professional\_relations.grf 95 Verbkonstruktionen in verschiedenen Zeitformen und Variationen:

- to be adopted
- to be appointed
- to be chosen
- to be commissioned
- to be coopted
- to be designated
- to be elected
- to be engaged
- to be installed
- to be named
- to be nominated
- to be selected
- to be voted in
- to decapitate so.
- to decruit so.
- to disband so.
- to discard so.
- to discharge so.
- to dismiss so.
- to displace so.
- to eject so.
- to expel so.
- to fire so.
- to give so. one's notice
- to lay (off)
- to make so. redundant
- to oust so.

- to pay so. off
- to release so.
- to relieve so.
- to remove so.
- to throw so. out
- to be decapitated
- to be decruited
- to be disbanded
- to be discarded
- to be discharged
- to be dismissed
- to be displaced
- to be ejected
- to be expelled
- to be fired
- to be given one's notice
- to be laid (off)
- to be made redundant
- to be ousted
- to be paid off
- to be released
- to be relieved
- to be removed
- to be thrown out
- to enrol so.
- to employ so.
- to engage so.

- to enlist so.
- to hire so.
- to put so. on the payroll
- to recruit so.
- to sign (up) so.
- to take so. into employment
- to take on
- to retain so.
- to secure the services of so.
- to be enrolled
- to be employed
- to be engaged
- to be enlisted
- to be hired
- to be put on the payroll
- to be recruited
- to be signed (up)
- to be taken into employment
- to be taken on
- to be retained
- to have an employment
- to join

- to become a member of
- to draw salary
- to be paid
- to be replaced as
- to resign
- to quit
- to leave job
- to retire s.o.
- to stop s.o. working
- to stop work
- to be retired
- to give up work
- to be stopped working
- to stop work
- to reach retirement age
- to work
- to job
- to labour
- to labor
- to operate
- to serve
- to toil

Mit Hilfe dieser Prädikate erfasst die Grammatik der beruflichen Relationen die wichtigsten Beziehungen, welche zwischen einer Person, ihrem Beruf und einer Firma bestehen können. Darunter fallen einerseits Verbalphrasen, die den Anfang eines Arbeitsverhältnisses oder den Beginn einer neuen beruflichen Karriere ausdrücken. Andererseits dürfen auch die Verben nicht fehlen, welche die Art einer Beschäftigung wiedergeben und diejenigen, die das Ende eines Beschäftigungsverhältnisses in Worte fassen.

Des Weiteren besteht auch hier die Möglichkeit jeden der Subgraphen einzeln im Kontext des Graphen merger\_professional\_relations.grf (siehe Seite 130) aufzurufen, um sich ein Bild von der jeweiligen Relation im Satz zu verschaffen.

Die folgenden Graphen, welche hier nicht abgebildet sind, übernehmen diese Aufgabe:

- merger\_appointed\_as.grf
- $merger\_dismiss\_as.grf$
- merger\_dismissed\_as.grf
- merger\_employed\_as.grf
- merger\_employ.grf
- merger\_joined.grf

- merger\_paid\_as.grf
- merger\_replaced\_as.grf
- merger\_resigned.grf
- $merger\_retire\_as.grf$
- merger\_retired\_as.grf
- merger\_work\_as.grf

In den nächsten Abschnitten werden die verschiedenen Grammatiken präsentiert, welche die Verbalphrasen mit den bereits genannten Prädikaten beschreiben. Außerdem wurden im Zusammenhang mit den Verben entsprechende Automaten entwickelt, um auf die Struktur einer Berufsbezeichnung eingehen zu können.

### 11.1 Der Beginn eines Beschäftigungsverhältnisses

### 11.1.1 Die Ernennung: "to be appointed as"

Wenn jemand einmal eine gewisse Position in einem Unternehmen erlangt hat, so wird derjenige beispielsweise zum Vorsitzenden, Chef oder Direktor ernannt. Eine andere Möglichkeit ist die Ernennung zum Richter, Professor oder Kardinal, wobei diese Berufsgruppen in Wirtschaftsnachrichten nicht so zahlreich wie z.B. die Manager vertreten sind. Auch in der Politik werden gewählte Volksvertreter unter anderem als Gouverneur, Minister, Präsident usw. in ihr jeweiliges Amt eingeführt.

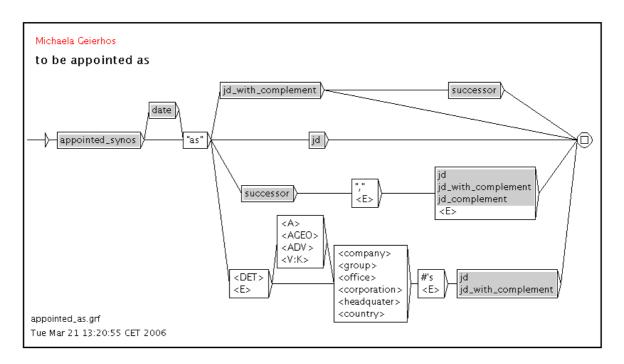
So bietet es sich an, Verben auszuwählen, welche den Aspekt der Berufung beinhalten, sich jedoch nicht auf eine bestimmte Domäne beschränken.

- to be adopted
- to be appointed
- to be chosen
- to be commissioned
- to be coopted
- to be designated
- to be elected

- to be engaged
- to be installed
- to be named
- to be nominated
- to be selected
- to be voted in

Diese Passivkonstruktionen werden vom Graphen appointed\_synos.grf in Abbildung 11.4 auf Seite 136 behandelt, welcher als Subgraph in der Grammatik zu "to be appointed as" zum Einsatz kommt.

Dabei ist die eigentliche Aufgabe der Grammatik aus Abbildung 11.2 die Erkennung von Verbalphrasen, in denen eine Person einen bestimmten Beruf ausübt oder eine gewisse Position erlangt.



**Abbildung 11.2:** Graph zur Erkennung von Verbalphrasen mit dem Verb "to be appointed as" und seiner direkten Synonyme - appointed\_as.grf

Im Umfeld von so genannten "Ernennungsrelationen" sind hilfreiche Informationen zu finden, welche die Umstände dieser Amtseinsetzung näher spezifizieren.

So wird meist im Zuge einer Amtseinführung angegeben, wessen Nachfolge nun angetreten wird. Um auf diese Textpassagen nicht verzichten zu müssen, wurde der Subgraph successor.grf (siehe Abbildung 11.3 auf Seite 135) entwickelt, welcher Personennamen im Kontext einer "Nachfolgerelation" aufspürt.

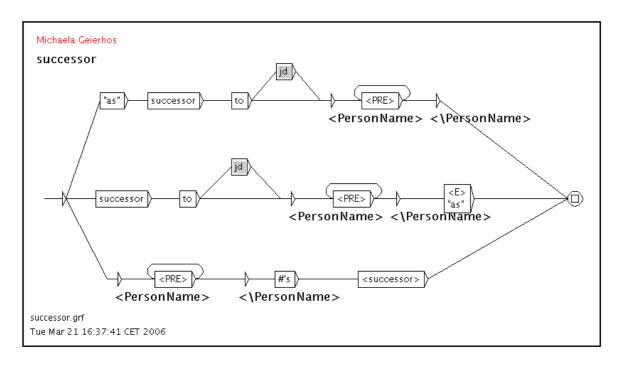


Abbildung 11.3: Graph zur Erkennung von Nachfolgern - successor.grf

Einen Vorgeschmack auf die volle Funktionalität des Graphen appointed\_as.grf gibt die folgende Konkordanz aus dem FT Korpus, wobei die Verberkennung vom Automaten appointed\_synos.grf (siehe Seite 136) vorgenommen wurde:

```
Arif Khan had been appointed as <PersonName> Hizb-ul Mojahedin<\PersonName>'s
     <JD>divisional commander<\JD> for <GEO>south Kashmir<\GEO> two years ago
after being appointed as <GEO>Netanya<\GEO>'s <JD>coach<\JD> <Date>on Sunday<\Date>
            following <PersonName>Eli Cohen<\PersonName>'s resignation
     Tun Mohamed Dzaiddin Abdullah has been appointed as <JD>chairman<\JD> of
                         <ORG>Deutsche Bank Malaysia<\ORG>
  Adrian Spencer Keane (41) has been appointed as <JD>Finance Director<\JD> and
        <JD>Company Secretary<\JD> with effect from <Date>1 May 2004<\Date>
   Charles Craven has been appointed as <JD>director<\JD> of <Sector>strategic
                                consulting<\Sector>
   Paolo Ceretti has been appointed as <JD>general manager<\JD> of <ORG>Italian
                        publishing group De Agostini.<\ORG>
   Abdol Wali Khan Zadran has been appointed as the <JD>head<\JD> of <GEO>Wazah
                   District<\GEO> in <GEO>Paktia Province<\GEO>
  Gerd Weiland, a lawyer from Hamburg, has been appointed as the <JD>provisional
                           insolvency administrator<\JD>
    Annie Bacon has been appointed as the company's <JD>vice president<\JD> of
                            <Sector>marketing<\Sector>
 Thomas P. Rice has been elected as the Company's <JD>Chief Executive Officer<\JD>
```

Hierbei wurden Berufsbezeichnungen an den verschiedensten Positionen im Text erkannt. Wie der Automat diese identifizieren konnte, wird im nächsten Abschnitt ersichtlich.

 $\langle JD \rangle (CEO) \langle JD \rangle$ 

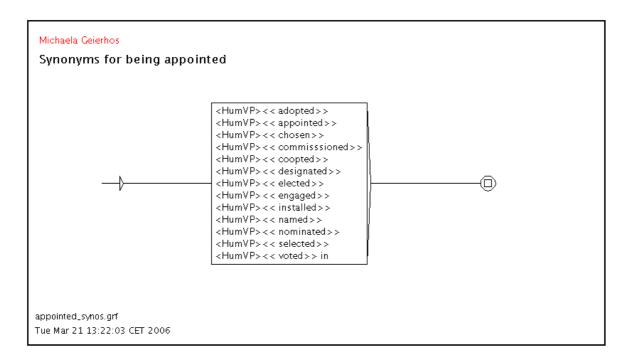


Abbildung 11.4: Graph zur Erkennung von Verben mit der Bedeutung von "to be appointed" - appointed\_synos.grf

#### 11.1.1.1 Grammatik der Berufsbezeichnungen

Wie die Beispielkonkordanz zur Nachfolgerelation auf Seite 134 gezeigt hat, sind Berufsbezeichner zwar auch als Attribute von Personennamen im Korpus zu finden, doch sollten sie an erster Stelle die Position angeben, für welche die jeweilige Person vorgesehen ist.

Um die syntaktische Variabilität von Berufsbezeichnern in den Griff zu bekommen, wurde eine Reihe von Automaten entworfen, welche sich diesem Problem annehmen. In Abbildung 11.5 auf Seite 137 ist der Hauptgraph jd.grf zu sehen, welcher weitere Subgraphen koordiniert.

Dieser Transduktor stützt sich nicht nur auf das in den Lexika kodierte Wissen über Berufe, sondern versucht auch den linken Kontext von Berufsbezeichnungen genauer zu beschreiben. Dabei geht er auf mögliche Adjektive bzw. Adverbien, sowie auf Nominalphrasen ein, welche die Beschäftigung einleiten und näher spezifizieren.

Phrasen, in denen diese Adjektive wichtige Eigenschaften über die Arbeit der betreffenden Person aussagen, könnten wie folgt aussehen:

```
Deputy Prime Minister Viktor Khristenko has been appointed as \langle JD \rangle acting prime minister \langle JD \rangle is expected to be appointed as \langle JD \rangle general production co-ordinator \langle JD \rangle the finance director of Pearson PLC, has been appointed as a \langle JD \rangle non-executive director \langle JD \rangle
```

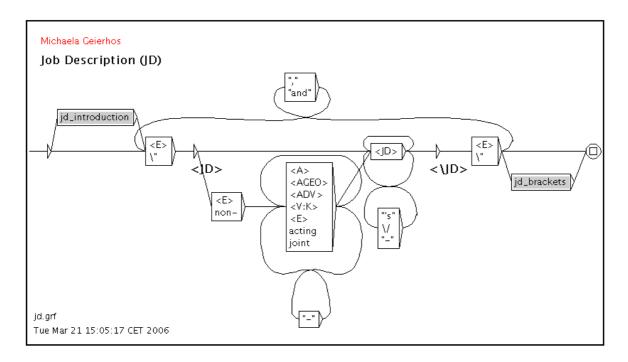


Abbildung 11.5: Graph zur Erkennung von Berufsbezeichnungen - jd. grf

Um diese attributiven Ergänzungen zu den jeweiligen Berufsbezeichnungen zu lokalisieren, reichte es über die Symbole <a>, <a>eachzo, <a>eachzo,

Dagegen ist es für die einleitenden Nominalphrasen notwendig, eine Grammatik anzugeben. Sie hat dann die Aufgabe Phrasen wie diese

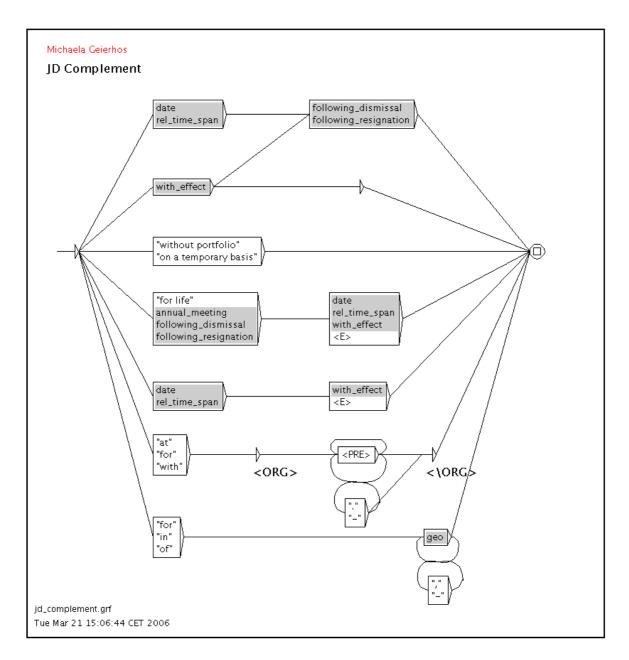
im Korpus zu erkennen. Eine entsprechende Abbildung zum Graphen  $jd\_introduction.grf$  wird im Anhang D auf Seite 171 unter dem Punkt D.1 gezeigt. Ebenso ist dort auch in Abbildung D.2 der Graph  $jd\_brackets.grf$  abgedruckt, der auf Firmenkürzel spezialisiert ist, welche in Klammern nach einem Organisationsnamen stehen.

Wie bereits in Abbildung 11.2 auf Seite 134 ersichtlich war, gibt es noch zwei weitere Grammatiken, welche den rechten Kontext des Verbes "to be appointed as", sowie seiner Synonyme beschreiben.

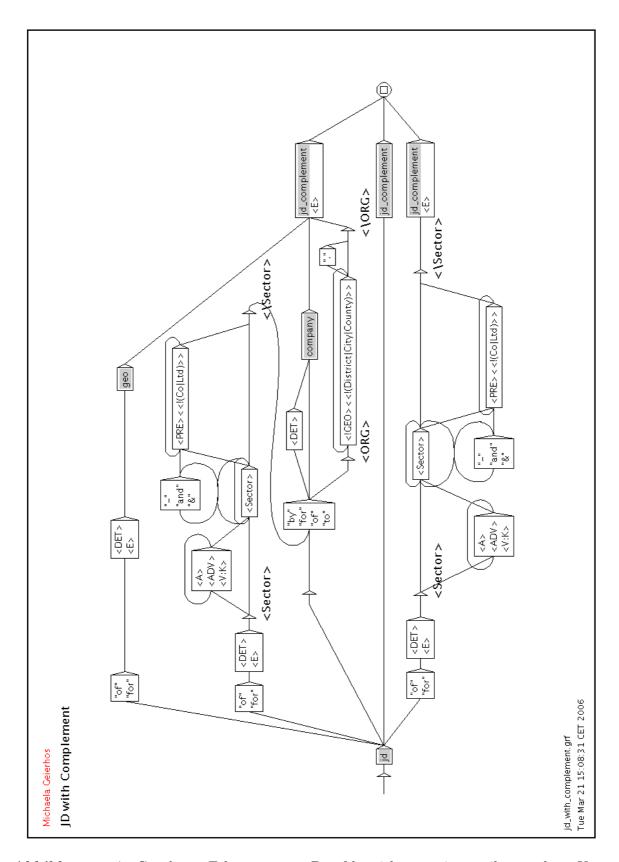
Dabei handelt es sich um den Automaten  $jd\_complement.grf$ , welcher den allgemeinen rechten Kontext einer Berufsbezeichnung spezifiziert und  $jd\_with\_complement.grf$ , der die Berufsbezeichnung gemeinsam mit ihrem attributiven rechten Kontext wiedergibt.

Die beiden Graphen werden auf den Seiten 138 und 139 dargestellt. Jedoch sind die dazugehörigen Subgraphen des Transduktors jd-complement.grf im Anhang D auf den Seiten 172 und 173 zu finden.

Dabei ist der Inhalt dieser Subgraphen fast selbsterklärend, da sie entweder Gründe (dismissal, resignation) angeben, warum der Posten neu besetzt wurde, oder ab wann die betreffende Person ihr Amt inne hat, oder ob die Ernennung eventuell auf der Jahresversammlung einer Firma bekannt gegeben wurde.

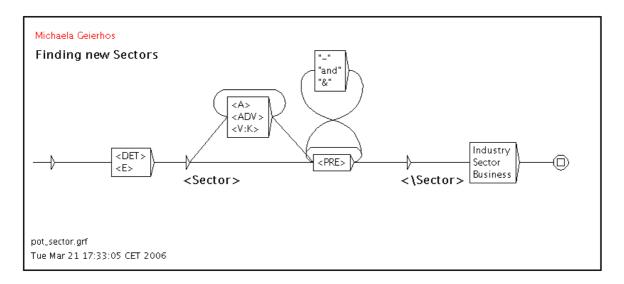


**Abbildung 11.6:** Graph zur Erkennung der Komplemente von Berufsbezeichnungen (rechter Kontext eines Berufsbezeichners) - jd-complement.grf



 ${\bf Abbildung\ 11.7:}\ {\bf Graph\ zur\ Erkennung\ von\ Berufsbezeichnern\ mitsamt\ ihrer\ rechten\ Kontexte\ -\ \it jd\_with\_complement.grf$ 

#### 11.1.1.2 Vervollständigung des Sektorenlexikons



**Abbildung 11.8:** Graph zur Erkennung potentieller Sektoren- und Branchenbezeichnungen - pot\_sector.grf

Bereits im Graphen  $jd\_with\_complement.grf$  auf Seite 139 fiel der Begriff des Sektors, der dort als Lexikonreferenz in Form von <Sector> eingesetzt, und später in der Ausgabe des Transduktors entsprechend annotiert wurde. In diesem Kontext wurden keine Variablen anstelle des Symbols <Sector> zugelassen, weil die vorangehende Präposition oft nicht ausreicht um eine Branche eindeutig zu identifizieren. Da manche Firmen Begriffe wie Marketing, Design oder Food enthalten, könnten Firmennamen irrtümlicherweise zu Sektorennamen und einem unbedeutenden Rest aufgespalten werden.

Aus diesem Grund wurde der endliche Automat aus Abbildung 11.8 erstellt, welcher auf Schlüsselbegriffe im rechten Umfeld von potentiellen Branchenbegriffen baut, um qualitativ gute Ergebnisse bei der Suche nach neuen Wörterbucheinträgen zu erzielen.

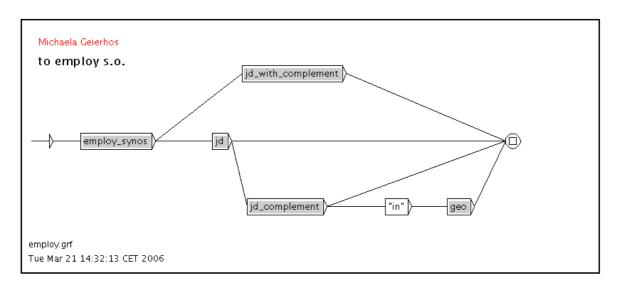
Dieser Graph findet nun folgende Kandidaten für neue Branchenbegriffe, die noch nicht in Lexikon Sector-.dic (siehe Abschnitt 5.2.7.2 auf Seite 64) enthalten sind.

Natürlich hätte man auch andere Indikatoren, wie z.B. "Director of", verwenden können, um neue Sektorennamen zu erhalten. Doch wenn nicht mehr Kontext spezifiziert wird, sind beide Ansätze für Fehler äußerst anfällig und bedürfen einer manuellen Korrektur, bevor Begriffe automatisch aus dieser Liste extrahiert werden.

### 11.1.2 Die Einstellung: "to employ so."

Nachdem die Verbkonstruktion "to be appointed" hier sehr ausführlich behandelt wurde, und sie in Kapitel 12 für die Evaluation des Systems noch einmal aufgegriffen wird, müssen die nachfolgenden Prädikate nicht mehr so intensiv behandelt werden. Das liegt unter anderem auch daran, dass sie alle auf die Berufsbezeichnergraphen (siehe Abschnitt 11.1.1.1 und Anhang D) zugreifen.

So fällt beispielsweise die Tatsache, dass eine Firma oder eine Person jemanden einstellt, der für sie arbeitet, auch in die Kategorie der beruflichen Relationen. Hierbei handelt es sich sogar um eine Prädikatsbeziehung, die zwei Menschen betreffen kann, und ist somit für diese Arbeit von besonderem Interesse.

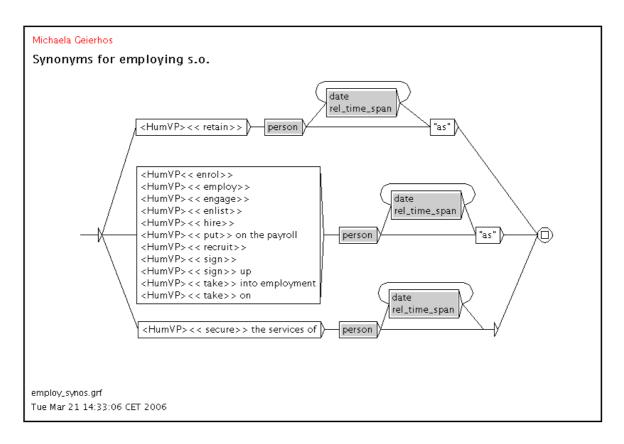


**Abbildung 11.9:** Graph zur Erkennung von Verbalphrasen mit dem Verb "to employ s.o." und seiner direkten Synonyme - employ.grf

Wenn man den Automaten in Abbildung 11.9 betrachtet, ist es leicht mit dem Vorwissen aus dem letzten Abschnitt seine Funktionsweise zu verstehen. Welche Verben er jedoch behandelt, wird nicht auf den ersten Blick ersichtlich, denn dafür ist der Subgraph *employ\_synos.grf* auf Seite 142 zuständig, welcher folgende Verbkonstruktionen umfasst:

- to enrol so.
- to employ so.
- to engage so.
- to enlist so.
- to hire so.
- to put so. on the payroll

- to recruit so.
- to sign (up) so.
- to take so. into employment
- to take on
- to retain so.
- to secure the services of so.



**Abbildung 11.10:** Graph zur Erkennung von Verben mit der Bedeutung von "to employ s.o." - employ\_synos.grf

Eine entsprechende Konkordanz zur Grammatik *employ.grf* (siehe Abbildung 11.9) könnte beispielsweise so aussehen:

```
it had hired <Person> <PersonName>Stewart Lawrie<\PersonName> <\Person> as
      <JD>director<\JD> of UK and international accounts from Proctor & Gamble
    announced that it has hired <Person> <PersonName>Frank Fiorillo<\PersonName>
         <\Person> as <JD>vice president<\JD> of worldwide customer support
         is pleased to announce that it has hired <Person> <PersonName>Jeri
              Silverman<\PersonName> <\Person> as a <JD>consultant<\JD>
  Evergreen Investments has hired <Person> <PersonName>Louis Membrino<\PersonName>
         <\Person> as a <JD>Director<\JD> in Evergreen Consultant Relations
Mr Stewart has recruited <Person> <PersonName>Lynne Peacock<\PersonName> <\Person> as
                <JD>business development director<\JD> at NAB Europe
  Hansen has retained <Person> <PersonName>Colin Charvis<\PersonName> <\Person> as
                                  <JD>captain<\JD>
has signed <Person> on <JD>legendary boxer<\JD> <PersonName>Muhammad Ali<\PersonName>
                          </Person> as a <JD>spokesman<\JD>
 as he was recently hired to replace the <Person>Irishman Johnny Murtagh<\Person> as
            the GEO>Aga<\GEO>'s GEO>hgirst-choice jockey\D> in Britain
```

Würde man stattdessen den Graphen merger\_employ.grf (ohne Abbildung) auf den Text anwenden, so sollten in der Konkordanz nur Sätze mit menschlichem Subjekt vorkommen und Firmen könnten auf diese Weise herausgefiltert werden.

### 11.1.3 Der Firmeneintritt: "to join"

Im Englischen sind wohl die folgenden Verbkonstruktionen die gebräuchlichsten Möglichkeiten, um den Eintritt einer Person in eine Firma zu beschreiben:

- to join
- to become a member of

Aufgrund ihrer minimalen Anzahl werden sie direkt in der Grammatik join.grf auf Seite 144 behandelt. Diese lokalisiert unter anderem die jeweiligen Unternehmen, zu denen die betreffende Person gegangen ist. Hierfür wird der Automat company.grf (siehe Seite 87) eingesetzt, welcher Firmennamen auf sehr umfassende Weise beschreibt. Außerdem werden noch weitere Informationen, wie der Sitz der Firma mit Hilfe des Graphen geo.grf (siehe Seite 93), das Einstellungsdatum mit date.grf (siehe Seite 96), oder der Beruf der betreffenden Person (siehe jd.grf auf Seite 137), sowie Beschäftigungszeit durch den Automaten rel\_time\_span.grf (siehe Seite 123) ermittelt.

Der Transduktor join. grf würde die Verbalphrasen wie folgt annotieren:

```
Jim Sweeney will also be joining <ORG>AmeriQuest<\ORG> as <JD>Vice President<\JD>
Andrew Bird will be joining <ORG>Hot Group<\ORG> as <JD>head<\JD> of e-marketing
    is joining the <ORG>Selfridges board<\ORG> as <JD>deputy chairman<\JD>
        and join <ORG>Reuters<\ORG> as <JD>non-executive chairman<\JD>
        But Mr Stewart, who joined <ORG>BSkyB<\ORG> <Date> in 1996<\Date>
He joined <ORG>Caledonian Insurance Services Limited<\ORG> <Date>in November<\Date>
Trevor Williams FCIS (42) joined <ORG>Imperial Tobacco<\ORG> <Date>in 1996<\Date>
Mr Rothschild, who joined <ORG>Tower Hotels<\ORG> <Date>in 1992<\Date>
```

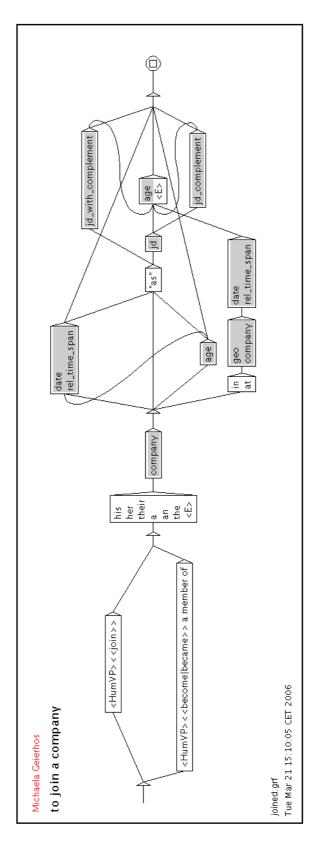
Möchte man jedoch erfahren, wie der Name der Person lautet, welche in die mit <ORG> markierte Firma eingetreten ist, muss der Graph merger\_join.grf (ohne Abbildung) aufgerufen werden, welcher dafür zuständig ist, ebenfalls das menschliche Subjekt im Satz zu kennzeichnen.

### 11.2 Die Ausübung des Berufes

Wurde eine Arbeitsstelle einmal angetreten, beginnt in der Regel die Zeit, in der ein Beruf ausgeübt wird. Die Tatsache, dass eine Person bei einem Unternehmen angestellt ist und somit dort arbeitet, kann auf die verschiedensten Arten ausgedrückt werden. Dabei muss es sich nicht nur um Aktiväußerungen handeln, in denen auf die Art der Tätigkeit des Arbeitnehmers eingegangen wird, es können auch Passivkonstruktionen sein, welche das aktuelle Beschäftigungsverhältnis eines Menschen beschreiben.

## 11.2.1 Das Beschäftigungsverhältnis: "to be employed"

Das Prädikat "to be employed" ist ein typisches Beispiel für so eine Passivphrase. Bereits in Abschnitt 11.1.2 wurde das Verb to employ in seiner Aktivform eingeführt. Dabei ging es um den Akt des Einstellens - eine Firma oder eine Person stellt einen Arbeitssuchenden bei sich ein. Bei dieser Variante handelte es sich um den Beginn eines



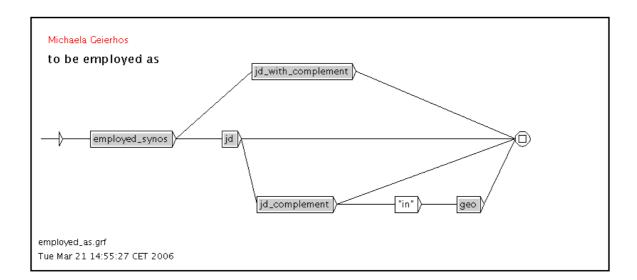
**Abbildung 11.11:** Graph zur Erkennung von Verbalphrasen mit dem Verb "to join" und seiner Paraphrasierung "to become a member of" - joined.grf

Arbeitsverhältnisses, wobei die passive Ausdrucksweise die Tatsache wiedergibt, dass jemand gerade bei einer Organisation beschäftigt ist. Das heißt nichts anderes, als dass derjenige schon die Phase des Firmeneintritts hinter sich gebracht hat und nun voll und ganz im Arbeitsleben steht.

Deshalb mussten lediglich die Verben des Aktivgraphen to employ für diese Grammatik in ihre jeweilige Passivkonstruktion überführt werden und konnten so im Automaten employed\_synos.grf (siehe Abbildung 11.13 auf Seite 146) kodiert und von employed\_as.grf (siehe Abbildung 11.12) verwendet werden.

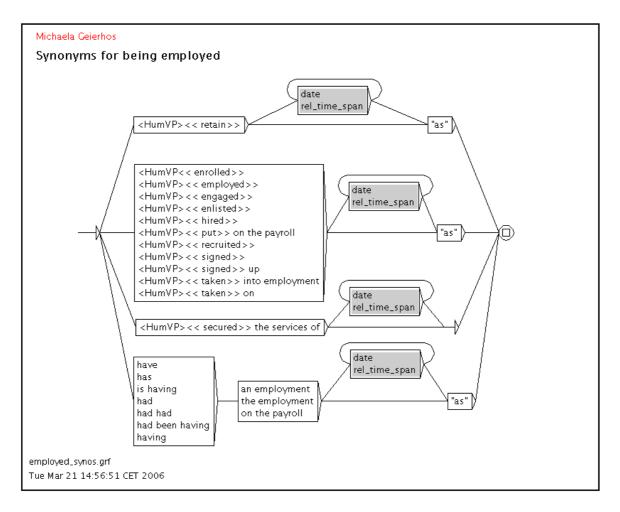
- to be enrolled
- to be employed
- to be engaged
- to be enlisted
- to be hired
- to be put on the payroll

- to be recruited
- to be signed (up)
- to be taken into employment
- to be taken on
- to be retained
- to have an employment



**Abbildung 11.12:** Graph zur Erkennung von Verbalphrasen mit dem Verb "to be employed as" und seiner direkten Synonyme -  $employed\_as.grf$ 

Um jedoch herauszufinden, wer die beschäftigte Person ist, muss allerdings der Graph  $merger\_employed\_as.grf$  (ohne Abbildung) zum Einsatz kommen, welcher den ganzen Satz und nicht nur die Verbalphrase analysiert.



**Abbildung 11.13:** Graph zur Erkennung von Verben mit der Bedeutung von "to be employed as" - employed\_synos.grf

Eine mögliche Konkordanz des Automaten *employed\_as.grf* (siehe Abbildung 11.13) würde folgendermaßen aussehen:

### 11.2.2 Die Bezahlung: "to be paid as"

Wenn jemand in einem Unternehmen beschäftigt ist, wird derjenige in der Regel auch für seine Arbeit finanziell entschädigt. Also kann davon ausgegangen werden, dass die Bezahlung als Gegenleistung für eine bereits verrichtete bzw. noch fortführende Beschäftigung erfolgt. Diese Tatsache wird in der Regel im Passiv ausgedrückt, wobei der Name des Arbeitgebers meist in diesem Zusammenhang nicht fällt. Um die Relation angemessen beschreiben zu können, wurden die folgenden beiden Prädikate dafür in Betracht gezogen.

- to draw salary
- to be paid

Im Graph der Abbildung 11.14 wurde nur die Präposition "as" im Anschluss an das jeweilige Verb zugelassen, auf welche eine Berufsbezeichnung folgen muss. Somit sollten nur Sätze in der Konkordanz enthalten sein, welche inhaltlich darauf Bezug nehmen, dass eine Person für eine bestimmte Tätigkeit (ihren Beruf) entlohnt wird.

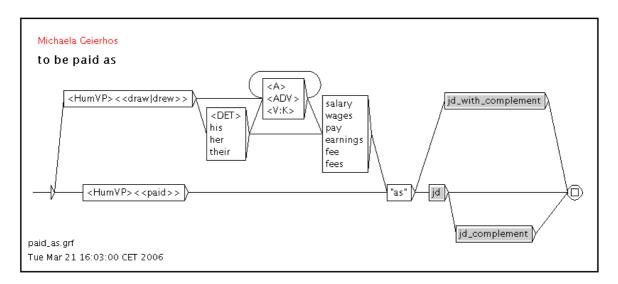
Corman, who is paid as a <JD>consultant<\JD>, holds 650,000 options

Dabei wird die Suche nach den entsprechenden Berufsbezeichnern und ihren Kontexten wieder von den Transduktoren jd.grf (siehe Seite 137),  $jd\_complement.grf$  (siehe Seite 138) und  $jd\_with\_complement.grf$  (siehe Seite 139) übernommen.

Außerdem werden bei der Floskel "to draw salary as" noch weitere Variationen zugelassen, indem noch zusätzlich die Begriffe

wages, pay, earnings, fee, fees

als Synonyme von "salary" im Vokabular des Graphen zugelassen werden.



**Abbildung 11.14:** Graph zur Erkennung von Verbalphrasen mit dem Verb "to be paid as" und seiner Paraphrasierung "to draw salary as" - paid\_as.grf

### 11.2.3 Die Tätigkeit: "to work as"

Wenn es darum geht, das Betätigungsfeld einer Person näher zu spezifizieren, bietet sich an, einfach auf die Relation "to work" zurückzugreifen. Auch ist offensichtlich, dass sich diese Person in einem momentanen Beschäftigungsverhältnis befinden muss.

Natürlich gibt es weitere Verben, welche den gleichen Sachverhalt ausdrücken:

• to work

• to operate

• to job

• to serve

• to labour

• to labor

• to toil

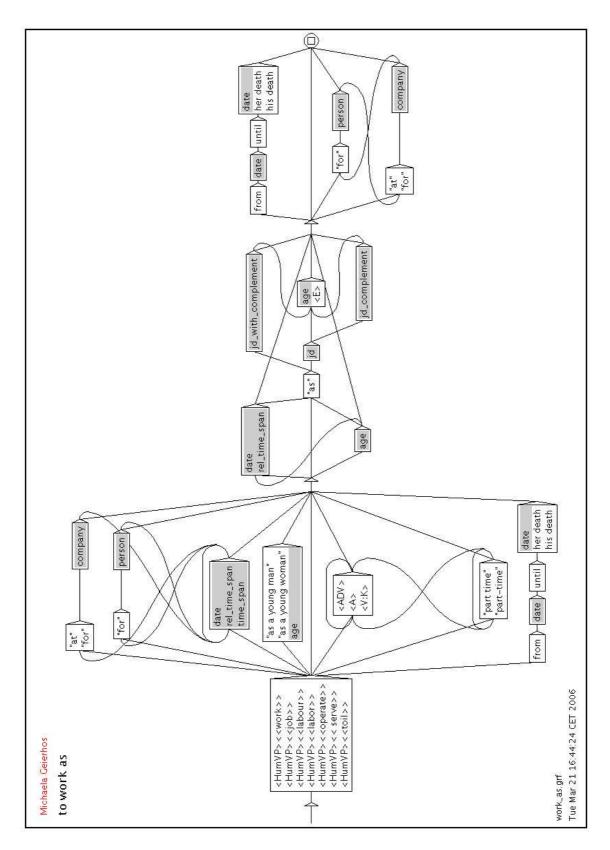
Zudem werden im Umfeld einer Arbeitsbeziehung wichtige Fragen geklärt, die sowohl den Arbeitgeber als auch den Arbeitnehmer betreffen:

- 1. Für wen arbeitet die betreffende Person? Wer ist ihr Arbeitgeber?
  - a) Handelt es sich hierbei um eine Organisation?
  - b) Oder wird sie von einer Privatperson beschäftigt?
- 2. <u>Seit wann</u> ist die betreffende Person für ihren Arbeitgeber tätig?
- 3. Wie lange arbeitet die betreffende Person schon bei diesem Arbeitgeber?
- 4. <u>Von wann bis wann</u> hat die betreffende Person dort gearbeitet?
- 5. Als was ist die betreffende Person dort beschäftigt? Welche Position hat sie inne?
- 6. <u>In welcher Branche</u> hat die betreffende Person gearbeitet?
- 7. Wie ist die betreffende Person beschäftigt? Vollzeit? Teilzeit?
- 8. Welches Alter hatte die betreffende Person, als sie dieser Tätigkeit nachging?

Jedoch wird der Arbeitnehmer selbst nicht in der Grammatik ermittelt, die sich mit den oben genannten Fragen beschäftigt, sondern im endlichen Automaten merger\_work\_as.grf (ohne Abbildung).

Dagegen ist der Graph, welcher für die Erkennung der Verbalphrase "to work as" zuständig ist, auf der nächsten Seite abgedruckt.

Die dazugehörige Konkordanz ist auf Seite 150 zu finden und illustriert, welche Informationen der Transduktor im Text erkannt und annotiert hat.



**Abbildung 11.15:** Graph zur Erkennung von Verbalphrasen mit dem Verb "to work as" und seiner direkten Synonyme -  $work_-as.grf$ 

```
Chief Operating Officer, has agreed to serve as <JD>interim Director<\JD> of Network Marketing.{S}
      President Jean-Bertrand Aristide can continue to serve effectively as <GEO>Haiti<\GEO>'s <JD>leader<\JD>.{S}
         Gary L. Primes, who continues to serve as <JD>CIO<\JD> in addition to his recent appointments as COO
                  and hired Gutman, who had been working as <JD>CEO<\JD> of 300-employee company Liraz
                     John Blue, who had served as <JD>Acting CEO<\JD> from <Date>March 2003<\Date>
                  David Brocklehurst, had served as <JD>finance director<\JD> for several ad agencies
          who had served as <JD>justice secretary<\JD> during the administration of deposed President Joseph
              bestselling author and historian who had served as <JD>librarian<\JD> of <ORG>Congress<\ORG>
                         had served as <JD>police chief<\JD> for <GEO>the Kushiro region<\GEO>
           John Lewis Ashcroft had worked as a <JD>professional trapper</JD> near Guyra, at Kangaroo Camp. (S)
         Wells Rich Greene, and has also served as <JD>Assistant Vice President<\JD>, Marketing Communications
                         He has also served as <JD>President<\JD> of <ORG>Midway Airlines<\ORG>
                  Col Mohammad Esa has been working as the <JD>acting head<\JD> of the department.{S}
 Schneider has served as <JD>chief financial officer<\JD> and <JD>principal<\JD> of <ORG>Leonard Green & Partners<\ORG>
                    Sam Skinner has served as <JD>Co-Chairman<\JD> of <ORG>Hopkins and Sutter<\ORG>
             He has served as <JD>editor<\JD> of the English edition of Haaretz since its founding in 1997
  and has served as <JD>Executive Vice President<\JD> and <JD>Chief Financial Officer<\JD> at <ORG>Riverview Community
                                                       Bank<\ORG>
a member of the board of directors and has served as <JD head<\JD> of <Sector>all commercial banking activities<\Sector>
```

Abbildung 11.16: Konkordanz zum Graphen worked\_as.grf

topped by Ernie Zampese, 67, who will serve as a <JD>consultant<\JD> to the team's offense.{S} to Vice President of Finance and will serve as <JD>interim Chief Financial Officer<\JD> until a decision is made Boniface Mukhwana Mutali was working as an <JD>untrained teacher<\JD> at <ORG>Emusire High School<\ORG> in Vihiga

### 11.3 Das Ende eines Arbeitsverhältnisses

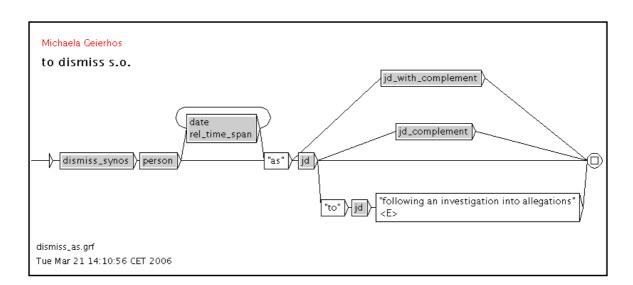
Jedes Arbeitsverhältnis wird irgendwann aufgelöst. Dabei kann es sich um ein abruptes Ende handeln, welches beispielsweise durch eine Kündigung ausgelöst wird, oder man ist sich des Termins bewusst, weil das Rentenalter erreicht wurde, ein besseres Johangebot vorliegt, oder man freiwillig aus dem Berufsleben ausscheiden will.

### 11.3.1 Die Entlassung: "to dismiss so." bzw. "to be dismissed"

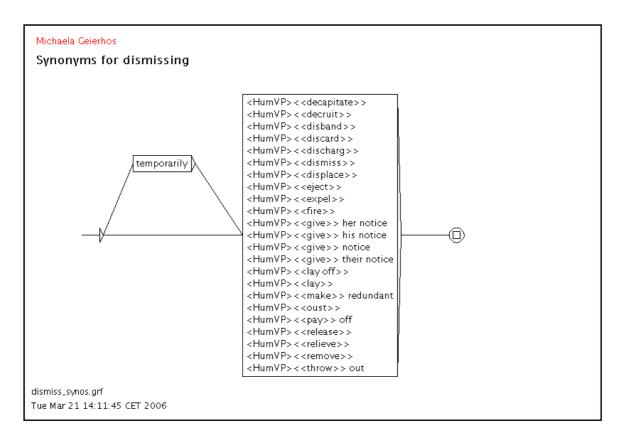
Im Englischen bringen folgende Verbkonstruktionen zum Ausdruck, dass ein Arbeitsverhältnis gekündigt wurde:

- to decapitate so.
- to decruit so.
- to disband so.
- to discard so.
- to discharge so.
- to dismiss so.
- to displace so.
- to eject so.
- to expel so.
- to fire so.

- to give so. one's notice
- to lay (off)
- to make so. redundant
- to oust so.
- to pay so. off
- to release so.
- to relieve so.
- to remove so.
- to throw so. out



**Abbildung 11.17:** Graph zur Erkennung von Verbalphrasen mit dem Verb "to dismiss s.o." und seiner direkten Synonyme - dismiss\_as.grf



**Abbildung 11.18:** Graph zur Erkennung von Verben mit der Bedeutung von "to dismiss s.o." - dismiss\_synos.grf

All diese Prädikate wurden in die Grammatik zur Erkennung von "Entlassungsergeignissen" aufgenommen. Allerdings handelt es sich bei diesen Verben um Aktivformen, d.h. dass sie Relationen wiedergeben, in denen eine Firma oder eine Person jemanden aus einem Beschäftigungsverhältnis entlässt.

Während der Subgraph dismiss\_synos.grf aus Abbildung 11.18 die einzelnen Verbkonstruktionen behandelt, analysiert der Hauptgraph dismiss\_as.grf auf Seite 151 den rechten Kontext dieser Prädikate. Hierbei wird obligatorisch zuerst die zu entlassende Person genannt, bevor möglicherweise Angaben über die Position gemacht werden, welche sie zuvor beansprucht hat.

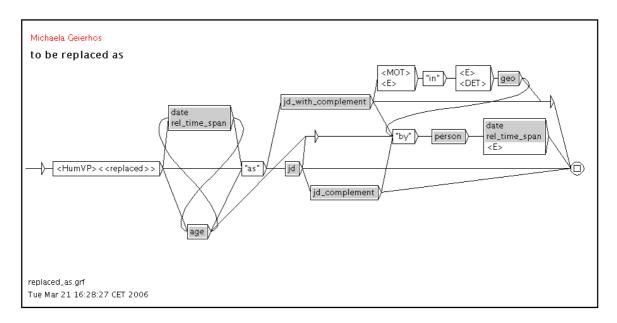
Natürlich lässt sich der gleiche Sachverhalt auch in einer passiven Äußerung darstellen. In diesem Fall nimmt die entlassene Person die Stelle es Subjekts im Satz ein, und die Information, wer sie gekündigt hat, fällt dann meist weg.

Im Wesentlichen verhalten sich die endlichen Automaten für diesen Typ Satz ganz analog zu der eben vorgestellten Phrasenstruktur und bedürfen deshalb keiner weiteren Erklärung. Die entsprechenden Graphen, Passivkonstruktionen der Verben und eine Beispielkonkordanz sind im Anhang E auf Seite 174 zu finden.

Dabei ist offensichtlich, dass sich die nachstehende Konkordanz und die aus Abbildung E.3 stark ähneln, was aber nicht verwunderlich sein dürfte.

#### 11.3.2 Die Nachfolge: "to be replaced as"

In großen Unternehmen ist es durchaus keine Seltenheit, dass von Zeit zu Zeit ein Machtwechsel stattfindet. Dabei liest man oft als Schlagzeile, dass wieder ein Generationenwechsel vollzogen wurde. Doch im Grunde wird hier nur eine Nachfolge für einen Firmenposten geregelt. Immer wenn jemand ein Amt abgibt, egal ob freiwillig oder ungewollt, muss dessen Stelle neu besetzt werden. Diese eben beschriebene Relation wird durch das Prädikat "to be replaced" ausgedrückt.



**Abbildung 11.19:** Graph zur Erkennung von Verbalphrasen mit dem Verb "to be replaced as" -  $replaced\_as.grf$ 

Welche Faktoren bei einer derartigen Neubesetzung einer Arbeitsstelle berücksichtigt werden müssen, versucht die Grammatik aus Abbildung 11.19 in den Griff zu bekommen. Denn einerseits sollte erkannt werden, welcher Posten neu zu vergeben ist, was durch den Automaten zur Lokalisierung von Berufsbezeichnungen (siehe Seite 137) er-

folgt. Andererseits dürfen Informationen, wie der Zeitpunkt des Wechsels, der Name der betroffenen Firma, und wer der bereits bestimmte Nachfolger ist, auf keinen Fall übergangen werden.

Wie die folgende Konkordanz zeigt, werden all diese Dinge im Kontext des Verbes to be replaced" berücksichtigt:

### 11.3.3 Die Abdankung: "to resign as"

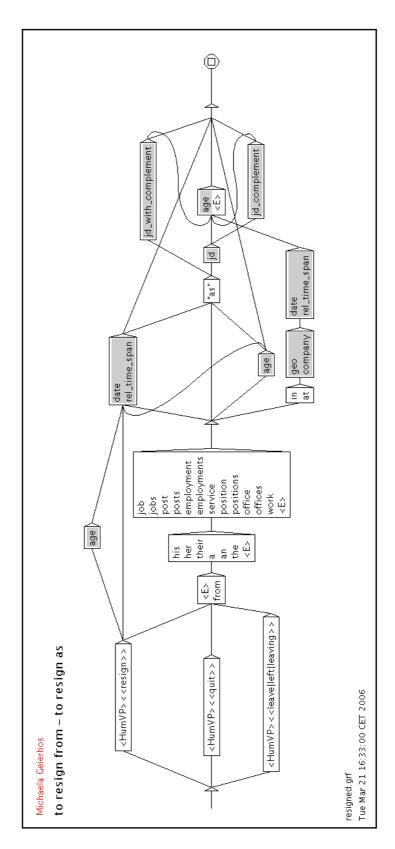
Wenn jemand darüber nachdenkt, abzudanken, dann ist dies vergleichbar mit einer Entlassung auf eigenen Wunsch. Um dieses Ereignis besser in Worte fassen zu können, wird die dazu entwickelte Grammatik aus Abbildung 11.20 auf der nächsten Seite folgende Prädikate behandeln:

- to resign
- to quit
- to leave job

Zudem werden bei der letztgenannten Verbkonstruktion im Vokabular des Graphen noch weitere Synonym- und Morphemvarianten von "job" zugelassen. Darunter fallen die Begriffe jobs, post, posts, employment, employments, service, position, positions, office, offices und work.

Auf diese Weise erfasst der Transduktor resigned.grf folgende Beispielphrasen:

```
Rauf Denktas said he had considered resigning as a <JD>negotiator<\JD> his daughter Lea Rose, 23, who had quit her job as a <JD>housemaid<\JD> in <GEO>Manila<\GEO> Chancellor Gerhard Schroder has resigned as <JD>chairman<\JD> of the <ORG>Social Democratic Party<\ORG> <ORG>(SPD)<\ORG> Martin Stewart is to quit as <ORG>BSkyB<\ORG>'s <JD>chief financial officer<\JD>
```



**Abbildung 11.20:** Graph zur Erkennung von Verbalphrasen mit dem Verb "to resign (as/from)" und seiner direkten Synonyme - resigned.grf

### 11.3.4 Die Pensionierung: "to retire so." bzw. "to be retired"

Bei der Pensionierung scheidet zwar eine Person aus Altersgründen aus dem Amt bzw. aus dem Berufsleben aus, doch geschieht dies in der Regel ebenfalls auf freiwilliger Basis.

Um die Tatsache zu beschreiben, dass jemand in Rente geht oder bereits im Ruhestand ist, kommen diese Aktiv- und Passivkonstruktionen folgender englischer Verben in Frage:

- to retire s.o.
- to stop s.o. working
- to stop work
- to be retired

- to give up work
- to be stopped working
- to stop work
- to reach retirement age

Die aktiven Verbalphrasen werden nun von der Grammatik retire\_as.grf aus Abbildung 11.21 und die passiven Konstruktionen vom Graphen retired\_as.grf auf Seite 157 erfasst.

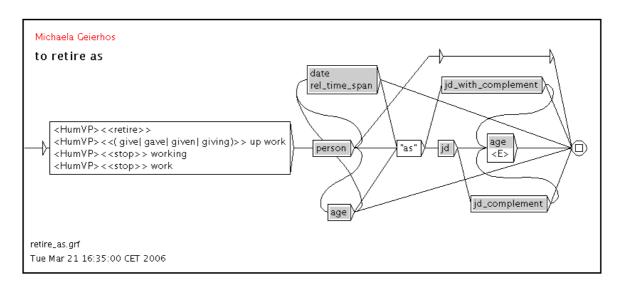
Da sich die beiden Automaten natürlich sehr stark gleichen, können sie hier zusammen behandelt werden, denn jeder von ihnen berücksichtigt das Rentenalter, eventuelle Angaben zum Beruf der jeweiligen Person, wann jemand in Rente gegangen ist, oder wie lange die Pensionierung schon zurückliegt.

Jedoch beschränken sich auch diese Graphen auf die Analyse von Verbalphrasen und geben keinen Aufschluss über ein mögliches Subjekt des Satzes. Hierfür sind die Transduktoren merger\_retire\_as.grf und merger\_retired\_as.grf (ohne Abbildung) zuständig, welche bestimmte Satzstrukturen untersuchen und außer der Verbalphrase noch das entspechende Subjekt annotieren.

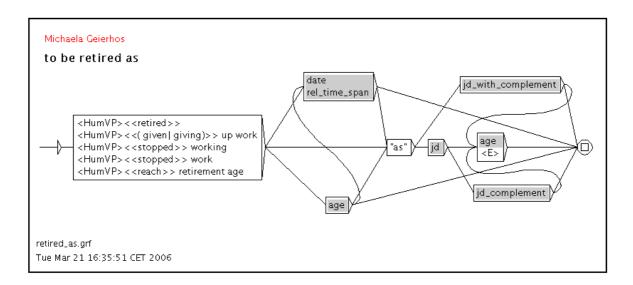
Die markierten Verbkonstruktionen mitsamt ihres rechten Kontextes könnten beispielsweise wie folgt aussehen:

Rocky Marciano, have retired as <JD>champion<\JD> and stayed retired Luzviminda Tancangco who retired <Date>on Feb. 2<\Date>
Denis Brosnan retired as <JD>Chairman<\JD> and <JD>Director<\JD> of the Group David Creary, who retired as <JD>chief operating officer<\JD> <Date>last week<\Date>
The 17-year-old son of retired Croatian <Person> General <PersonName>Vladimir Zagorac<\PersonName> <\Person> was kid

Minister Dan Naveh has appointed retired <Person> <ORG>Jerusalem District Court<\ORG> <JD>president<\JD> <PersonName>Vardi Zeiler<\PersonName> <\Person>



**Abbildung 11.21:** Graph zur Erkennung von Verbalphrasen mit dem Verb "to retire s.o." und seiner direkten Synonyme - retire\_as.grf



**Abbildung 11.22:** Graph zur Erkennung von Verbalphrasen mit dem Verb "to be retired as" und seiner direkten Synonyme - retired\_as.grf

## 12 Auswertung der Ergebnisse

### 12.1 Evaluationsmaße [97]

Vollständigkeit und Genauigkeit sind zwei Maße zur Beschreibung der Güte eines Suchergebnisses beim Information-Retrieval oder bei einer Recherche im Allgemeinen.

Für die Evaluierung eines Information-Retrieval-Systems sollten die beiden zusammenhängenden Maße gemeinsam betrachtet werden. In der Regel sinkt mit steigendem Recall (mehr Treffer) die Precision (mehr irrelevante Ergebnisse) und umgekehrt sinkt mit steigender Precision (weniger irrelevante Ergebnisse) der Recall (mehr relevante Dokumente, die nicht gefunden werden). Somit besteht eine negative Korrelation.

Stellt man das Verhältnis zwischen Recall und Precision in einem Diagramm dar, so wird der (höchste) Wert im Diagramm, an dem der Precision-Wert gleich dem Recall-Wert ist - also der Schnittpunkt des Precision-Recall-Diagramms mit der Identitätsfunktion - der Precision-Recall-Breakeven-Punkt genannt.

Für die Evaluierung des Information-Retrieval-Systems gibt es mit dem Fall-Out noch ein drittes Kriterium.

### 12.1.1 Precision bzw. Genauigkeit

Die Precision beschreibt die Genauigkeit eines Suchergebnisses. Sie ist definiert als der Anteil der gefundenen relevanten Dokumente von allen bei einer Suche gefundenen Dokumenten.

$$\label{eq:precision} \begin{aligned} \text{Precision} &= \frac{|\{\text{relevante Dokumente}\} \cap \{\text{gefundene Dokumente}\}|}{|\{\text{gefundene Dokumente}\}|} \end{aligned}$$

### 12.1.2 Recall bzw. Vollständigkeit

Der Recall beschreibt die Vollständigkeit eines Suchergebnisses. Er ist definiert als der Anteil bei einer Suche gefundenen relevanten Dokumente bzw. Datensätze an den relevanten Dokumenten der Grundgesamtheit.

$$\text{Recall} = \frac{|\{\text{relevante Dokumente}\} \cap \{\text{gefundene Dokumente}\}|}{|\{\text{relevante Dokumente}\}|}$$

#### 12.1.3 Fall-Out

Das Fall-Out beschreibt in negativer Weise die Güte des zu bewertenden Verfahrens, indem die Anzahl der gefundenen irrelevanten Dokumente durch die Gesamtanzahl irrelevanter Dokumente geteilt wird.

$$Fall-Out = \frac{|\{irrelevante\ Dokumente\} \cap \{gefundene\ Dokumente\}|}{|\{irrelevante\ Dokumente\}|}$$

### 12.2 Qualität des Systems

Um ein Verständnis dafür zu erlangen, wie präzise die im Laufe dieser Arbeit vorgestellten Automaten auf dem FT Korpus arbeiten, müssen die Suchergebnisse entsprechend ausgewertet werden.

Für die nun folgende Beispielauswertung der Treffermenge des Automaten person\_name.grf aus Abbildung 6.1 (siehe Seite 78) wird mit Hilfe der Evaluationsmaße Precision und Recall die Qualität dieser Grammatik in Bezug auf das Korpus ermittelt.

Hierfür ist es notwendig, die tatsächliche Anzahl an Vorkommen einer gesuchten Entität im Text zu kennen. Deshalb ist es kaum möglich diese Auswertung auf dem kompletten Financial Times Korpus durchzuführen, und es ist ratsam, sich nur auf einen Teil des Korpuses zu konzentrieren. Dieses Teilkorpus sollte nicht verstärkt zum Training der Automaten eingesetzt worden sein, da sonst kein repräsentatives Auswertungsergebnis erzielt werden könnte. Denn würde die Evaluation auf dem gleichen Text erfolgen, welcher schon während der Entwicklung der Grammatiken für Zwischentests eingesetzt wurde, wären herausragende Werte sicherlich eine direkte Konsequenz bei diesem Vorgehen.

Auch lässt sich plausibel erklären, warum diese Auswertung ausgerechnet mit dem endlichen Automaten zur Erkennung von Personennamen erfolgen soll. Wie bereits in Kapitel 6 diskutiert wurde, ist der Einsatz eines Automaten ohne größeres Kontextwissen ein sehr risikobehaftetes Unterfangen. Da der Personennamengraph viel weniger Kontextinformation als der Personengraph einsetzt und ohne die Kontexte der Verbalphrasenautomaten aufgerufen wird, müsste er verhältnismäßig schlechtere Ergebnisse erzielen.

Die folgende statistische Auswertung dieser Treffermenge soll nun folgendes beweisen: Wenn eine Grammatik auf qualitativ gute Lexikoneinträge setzen kann, fällt die vernachlässigte Kontextinformation nicht so stark in das Gewicht der Auswertung. Denn wäre die Wörterbuchbasis, welche der Automat nutzen konnte, wesentlich kleiner gewesen, könnten sich die Werte für Precision und Recall nicht mehr sehen lassen.

Tatsächliche Anzahl der Personennamen im Text	4267
Anzahl der gefundenen Personennamen im Text	4561
Anzahl der korrekt gefundenen Personennamen	3916
Recall	85,85%
Precision	91,77%

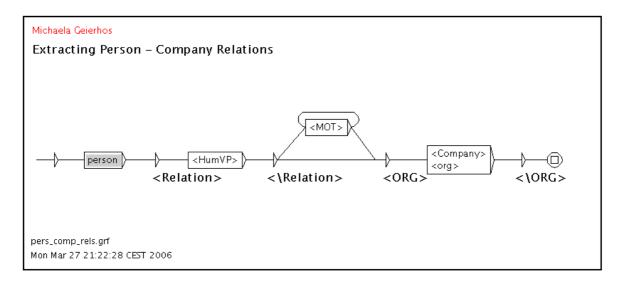
Tabelle 12.1: Statistische Auswertung der Extraktionsresultate

## 13 Anwendungen

Nachdem nun eine Reihe an verschiedenen ausgewählten Verbrelationen präsentiert wurde, ist es an der Zeit weitere Verbkonstruktionen zu ermitteln, welche mit den hier behandelten Entitäten wie z.B. den Personen und Organisationen einhergehen.

Die folgenden Graphen versuchen mit Hilfe einer groben Schematisierung des potentiellen Kontextes von personenbezogenen Prädikaten, diese auf relativ einfache und schnelle Weise automatisch aus den Korpora zu extrahieren.

# 13.1 Automatische Extraktion von Relationen zwischen Personen und Organisationen



**Abbildung 13.1:** Graph zur automatischen Erkennung weiterer Verbkonstruktionen, die eine Beziehung zwischen einer Person und einer Firma ausdrücken

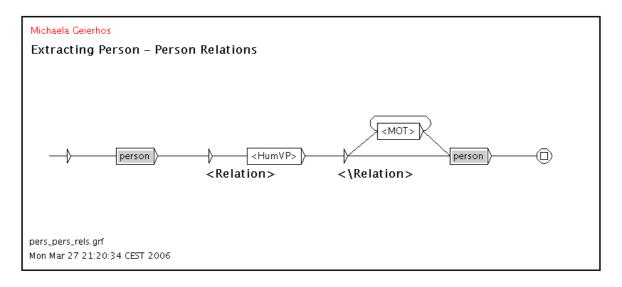
Im Laufe meiner Untersuchungen von biographischen Relationen stellte sich heraus, dass es ratsam ist, das Suchfenster nicht zu weit zu fassen, gerade wenn neue Informationen zwischen zwei Entitäten gewonnen werden sollen.

Zwar scheint der Graph aus Abbildung 13.1 auf den ersten Blick recht einfach aufgebaut zu sein, doch liefert er sehr gute Treffer bei der Erkennung potentieller Verbrelationen, welche im Zusammenhang mit Menschenbezeichnern und Firmennamen stehen.

Um die Effizienz des Automaten bei der Suche nach Kandidaten für Personen-Firmen-Relationen nicht allzu sehr zu beeinträchtigen, wurde auf den Einsatz des *company.grf*-Graphen verzichtet. Da es hierbei nicht um die genaue Lokalisierung von Organisationsnamen geht, reicht es mit den Symbolen <Company> und <org> aus den Lexika zu arbeiten und vorangehende Textpassagen mit dem "Platzhalter" <MOT>\* zu versehen.

Auf diese Weise lassen sich folgende Sequenzen im Korpus aufspüren und aufgrund ihrer Annotation automatisch weiterverarbeiten:

# 13.2 Automatische Extraktion von Relationen zwischen mindestens zwei Personen



**Abbildung 13.2:** Graph zur automatischen Erkennung weiterer Verbkonstruktionen, die eine Beziehung zwischen mindestens zwei Personen ausdrücken

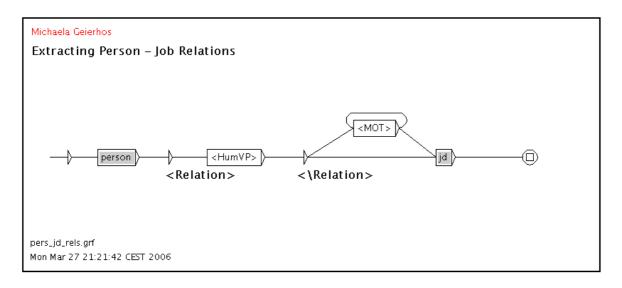
Das Gleiche gilt auch für die Beziehung zwischen zwei Menschen. Jedes Verb, welches an beiden Argumentpositionen einen Menschenbezeichner verlangt, sollte in einer Konkordanz dieses Graphen enthalten sein. Jedoch stellt dieser Automat - so wie schon sein Vorgänger - die Bedingung, dass die potentiellen Verben aus dem Lexikon der personenbezogenen Prädikate stammen müssen. Das ist unter anderem ein Grund dafür, warum tatsächlich nur Verben erkannt werden.

Leider hat diese Extraktionsmethode einen kleinen "Schönheitsfehler", da sie einen Zyklus in einer freien Variable wie <mot>MOT> zulässt. So kann es passieren, dass der Kontext zu stark ausgedehnt wird und der Bezug zum jeweiligen Verb verloren geht.

Wenn dies jedoch nicht der Fall ist, bekommt man ausgezeichnete Ergebnisse wie diese:

# 13.3 Automatische Extraktion von Relationen zwischen Personen und ihren Berufen

declared<\Relation> a terrorist by a <Person>Portuguese<\Person>



**Abbildung 13.3:** Graph zur automatischen Erkennung weiterer Verbkonstruktionen, die eine Beziehung zwischen einer Personen und ihrem Beruf ausdrücken

Nachdem hier die Thematik der Berufsbezeichnungen und den dazugehörigen Arbeitsverhältnissen ausgiebig behandelt wurde, bietet es sich noch an, einen weiteren Transduktor zu erstellen, welcher sich auf das Finden neuer Verbrelationen zwischen Menschenbezeichnern und ihren jeweiligen beruflichen Tätigkeiten spezialisiert.

So können damit beispielsweise folgende Treffer erzielt werden:

## 14 Zusammenfassung und Ausblick

Das Ziel dieser Arbeit war es Menschenbezeichner innerhalb biographischer Relationen in englischsprachigen Wirtschaftsnachrichten automatisch zu erkennen.

In diesem Zusammenhang wurde der Begriff der "biographischen Relation" ausführlich diskutiert, wobei dessen einzelne Facetten zum Vorschein kamen. Hierbei wurde ersichtlich, dass ein biographischer Kontext recht vielfältig sein kann und die Grammatikentwicklung zunächst ihre Tücken hat. Doch letztendlich besteht eine Möglichkeit, das textuelle Umfeld von Personen in einige Hauptkategorien zu unterteilen, und auf ihnen basierend erste linguistische Untersuchungen vorzunehmen.

Auf diese Weise können wiederum einzelne Verbkonstruktionen innerhalb dieser Gruppen, syntaktisch oder semantisch zusammengefasst, bzw. Ähnlichkeiten unter ihnen ausgemacht werden. Genauso verlief auch die Arbeit mit der Grammatik der beruflichen Relationen. Nach diversen Vorüberlegungen, welche Prädikate in die jeweilige Kategorie fallen, mussten diese semantisch klassifiziert, und im Anschluss an die strukturelle Analyse schematisiert und modularisiert werden. So konnten sich Teilgrammatiken in den verschiedensten Automaten wiederfinden, und wurden anhand diverser Kontexte konzeptionell und inhaltlich verbessert.

Für die hier ausgewählten Relationen ließ sich das Problem der automatischen Lokalisierung von Personen mit biographischen Kontexten relativ gut bewältigen. Außerdem liefern die entsprechenden Grammatiken zufriedenstellende Ergebnisse für die einzelnen Korpora.

Jedoch sollte man sich dessen bewusst sein, dass mit dieser Arbeit nur die "Spitze des Eisbergs" angekratzt wurde. Auch wenn das Feld der biographischen Relationen überschaubar ist, weil es eine endliche, nicht allzu große Menge von ihnen gibt, ist es ihr Kontext, welcher erst einmal analysiert und in Form von lokalen Grammatiken formalisiert werden muss. Je mehr man über diese Relationen im einzelnen auf syntaktischer und semantischer Ebene herausfindet, desto besser werden die Ergebnisse von "intelligenten" Informationsextraktionssystemen.

Mit dieser Arbeit wurde der erste Grundstein für sich weiterentwickelnde und umfassendere Grammatiken zur Erkennung von Menschenbezeichnern und ihrer biographischen Verbrelationen gelegt. Wobei sich das hier vorgestellte Konzept ohne größere Schwierigkeiten auf andere personenbezogene Prädikate, wie z.B. to stay with so., to live, to restore so. to life, to kill so., to laugh, to be moved to tears, to read sth., to cook sth., to sing sth., to kiss so., to love so., to visit so., to feat sth., to repare sth. übertragen lässt. Natürlich gibt es noch eine Reihe von Prädikaten, welche aus biographischer Sicht noch viel relevanter sind und ebenfalls untersucht werden sollten. Diese Verbkonstrukte können zukünftig Gegenstand weiterer linguistischer Studien sein und müssen hier nicht mehr aufgeführt werden.

# A Hilfe zum Tokenizer-Programm von Sebastian Nagel

#### tokenizer - a tokenizer with end-of-sentence detection

```
tokenizer OPTIONS [FILES]
options:
  -o <filen>
               output filename
  -L <lang>
               language in a specific charset, actually supported:
                 de german (iso-8859-1)
                 de-win german-win-cp1252 (cp1252)
                 de-u8 german-utf8 (rudimentary support for utf-8)
                 en english (iso-8859-1)
                 en-win english-win-cp1252 (cp1252)
                 en-u8 (utf-8)
                 ru russian (iso-8859-5)
                 ru-win russian-win-cp1251 (cp1251)
  -S
               enable end-of-sentence detection
  -E <str>
               specify EOS-mark (default: "<EOS />")
  -n
               treat a new line as EOS
  -N
               treat two or more new lines (paragraph break) as EOS
               combine continuation I: hyphenated words on line
  -с
                 breaks will be put together.
                 The hyphen is skipped.
  -C
               combine continuation II: same as above, but the
                 hyphen is preserved. This may be a good option
                 if you know that there are no hyphenated words,
                 but 'bindestrichwoerter' (like end-of-sentence)
                 in your text.
               detect www-adresses and treat them as one token
  -i | -1
               convert all tokens to lowercase (according to
                 language settings)
  -s
               single line mode: each token on a separate line
```

-X <chr> use <chr> as separator in single line mode instead <newline>. Original newlines are preserved because putting the whole input in one line isn't a good idea paragraph mode: -p two or more newlines are interpreted as a paragraph break, a single newline will not. All lines of one paragraph are collected in one line -P prints each sentence in a separate line print spaces: -xIn single line mode horizontal spaces will be printed as one space in a single line, vertikal spaces as two line breaks. In paragraph mode (including combination with -P) an additional newline is inserted between paragraphs -h | -? print this help and exit

Other arguments will be read as input filenames. If no input files are given, input is read from stdin.

If no output file is given, the tokenized text is written to stdout

WARNING: When the input contains long words or many following newlines tokenizer stops with "input buffer overflow". To avoid this use putzer (included in your package) with option -m <max-word-length> as filter!

tokenizer, v0.7, Sebastian Nagel (wastl@cis.uni-muenchen.de)

# B Übersicht aller Kategorien in den Wörterbüchern

### **B.1 Semantische Kategorien**

Abkürzung	Kategorietyp	wörtliche Bedeutung	Erläuterung
Abbrev	semantisch	Abbreviation	Abkürzung
ABourough	semantisch	Bourough Adjective	Adjektiv für einen Stadt- teil
ACity	semantisch	City Adjective	Adjektiv für eine Metropole, Stadt
AGEO	semantisch	Geographical Term Adjective	Adjektiv für ein Topo- nym, einen geographi- schen Begriff
ANation	semantisch	Nation Adjective	Adjektiv für ein Land
AProvince	semantisch	Province Adjective	Adjektiv für eine Provinz
AState	semantisch	State Adjective	Adjektiv für einen Bun- desstaat
AuProvinceCitizen	semantisch	Australian Province Citizen	Bewohner einer australischen Provinz
Borough	semantisch	Borough	Stadtteil, Stadtbezirk
CaProvinceCitizen	semantisch	Canadian Province Citizen	Bewohner einer kanadi- schen Provinz
CaProvince	semantisch	Canadian Province	Kanadische Provinz
Citizen	semantisch	Citizen	(Staats)Bürger
City	semantisch	City	Metropole, Stadt
Company	semantisch	Company Name	Organisationsname bzw. Firmenname
Continent	semantisch	Continent	Kontinent
County	semantisch	County	Grafschaft
Discipline	semantisch	Discipline	Fachbereich, Lehrbereich
Département	semantisch	Département	Département
FN	semantisch	First Name	Vorname
GEO	semantisch	Geographical Term	Toponym, geographi- scher Begriff
Hum	semantisch	Human	Menschenbezeichner
JD	semantisch	Job Descriptor	Berufsbezeichner

Abkürzung	Kategorietyp	wörtliche Bedeutung	Erläuterung
LN	semantisch	Long Name	Vollständiger Personen-
			name, bestehend aus
			Vor- und Nachname,
			evtl. mit Titel
Nation	semantisch	Nation	Land
NYCBourough	semantisch	New York City Bou-	Stadtteil von New York
		rough	City
NYCcitizen	semantisch	New York City Citizen	Bewohner von New York
			City
PR	semantisch	Proper Noun	Eigenname
Region	semantisch	Region	Region, Gebiet
Sector	semantisch	Sector	Sektor, Branche
SN	semantisch	Surname	Nachname
Title	semantisch	Title	Titel, wie z.B. akademi-
			sche Grade, aristokrati-
			sche Titel
Urbanite	semantisch	Urbanite	Stadtbewohner
USstate	semantisch	US State	US Bundesstaat
USstateCitizen	semantisch	US State Citizen	Bewohner eines US Bun-
			desstaates

## **B.2 Grammatikalische Kategorien**

Abkürzung	Kategorietyp	wörtliche Bedeutung	Erläuterung
A	grammatikalisch	Adjective	Adjektiv
N	grammatikalisch	Noun	Nomen
XA	grammatikalisch	Extended Adjective	lexikalische Einheit, wel-
			che die Funktion eines
			Adjektivs erfüllt
XN	grammatikalisch	Extended Noun	Mehrwortlexem

# C Syntaktische Variabilität am Beispiel von "Bill Gates"

- 1. B Gates
- 2. B. Gates
- 3. B Gates III
- 4. B Gates, III
- 5. B. Gates III
- 6. B. Gates, III
- 7. B Gates III, KBE
- 8. B. Gates III, KBE
- 9. B Gates, KBE
- 10. B. Gates, KBE
- 11. B H Gates
- 12. B. H. Gates
- 13. B H Gates III
- 14. B H Gates, III
- 15. B. H. Gates III
- 16. B. H. Gates, III
- 17. B H Gates III, KBE
- 18. B. H. Gates III, KBE
- 19. B H Gates, KBE
- 20. B. H. Gates, KBE

- 21. Bill Gates
- 22. Bill Gates III
- 23. Bill Gates, III
- 24. Bill Gates III, KBE
- 25. Bill Gates, KBE<sup>44</sup>
- 26. Bill Henry Gates
- 27. Bill Henry Gates III
- 28. Bill Henry Gates, III
- 29. Bill Henry Gates III, KBE
- 30. Bill Henry Gates, KBE
- 31. Bill H Gates
- 32. Bill H. Gates
- 33. Bill H Gates III
- 34. Bill H Gates, III
- 35. Bill H. Gates III
- 36. Bill H. Gates, III
- 37. Bill H Gates III, KBE
- 38. Bill H. Gates III, KBE
- 39. Bill H Gates, KBE
- 40. Bill H. Gates, KBE

 $<sup>^{44}</sup>Knights\ Commander\ of\ the\ British\ Empire$ 

- 41. Bill (William Henry) Gates
- 42. Bill (William Henry) Gates III
- 43. Bill (William Henry) Gates, III
- 44. Bill (William Henry) Gates III, KBE
- 45. Bill (William Henry) Gates, KBE
- 46. Bill (William H.) Gates
- 47. Bill (William H.) Gates III
- 48. Bill (William H.) Gates, III
- 49. Bill (William H.) Gates III, KBE
- 50. Bill (William H.) Gates, KBE
- 51. Gates
- 52. Gates, Bill (William H.)
- 53. Gates, Bill (William Henry)
- 54. Gates, Bill (William Henry) III
- 55. Gates, Bill (William Henry), III
- 56. Gates, Bill (William H.) III
- 57. Gates, Bill (William H.), III
- 58. Gates III
- 59. Gates III, KBE
- 60. Gates, William (Bill) H.
- 61. Gates, William (Bill) Henry
- 62. Gates, William (Bill) Henry III
- 63. Gates, William (Bill) Henry, III
- 64. Gates, William "Bill" H.
- 65. Gates, William "Bill" Henry
- 66. Gates, William "Bill" Henry III
- 67. Gates, William "Bill" Henry, III

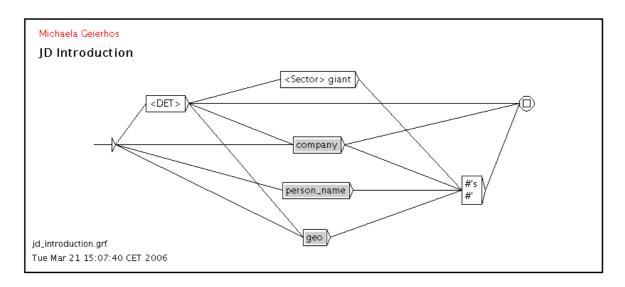
- 68. Gates, William H.
- 69. Gates, William Henry
- 70. Gates, William Henry III
- 71. Gates, William H, III
- 72. Gates, William H. III
- 73. Gates, William H., III
- 74. W H Gates
- 75. W. H. Gates
- 76. W H Gates III
- 77. W H Gates, III
- 78. W. H. Gates III
- 79. W. H. Gates, III
- 80. W H Gates III, KBE
- 81. W. H. Gates III, KBE
- 82. W H Gates, KBE
- 83. W. H. Gates, KBE
- 84. William (Bill) Henry Gates
- 85. William (Bill) Henry Gates III
- 86. William (Bill) Henry Gates, III
- 87. William (Bill) Henry Gates III, KBE
- 88. William (Bill) Henry Gates, KBE
- 89. William (Bill) H Gates
- 90. William (Bill) H. Gates
- 91. William (Bill) H Gates III
- 92. William (Bill) H Gates, III
- 93. William (Bill) H. Gates III
- 94. William (Bill) H. Gates, III

- 95. William (Bill) H Gates III, KBE
- 96. William (Bill) H. Gates III, KBE
- 97. William (Bill) H Gates, KBE
- 98. William (Bill) H. Gates, KBE
- 99. William Gates
- 100. William Gates III
- 101. William Gates, III
- 102. William Gates III, KBE<sup>45</sup>
- 103. William Gates, KBE
- 104. William "Bill" Henry Gates
- 105. William "Bill" Henry Gates III
- 106. William "Bill" Henry Gates, III
- 107. William "Bill" Henry Gates III, KBE
- 108. William "Bill" Henry Gates, KBE
- 109. William "Bill" H Gates
- 110. William "Bill" H. Gates
- 111. William "Bill" H Gates III
- 112. William "Bill" H Gates, III
- 113. William "Bill" H. Gates III
- 114. William "Bill" H. Gates, III

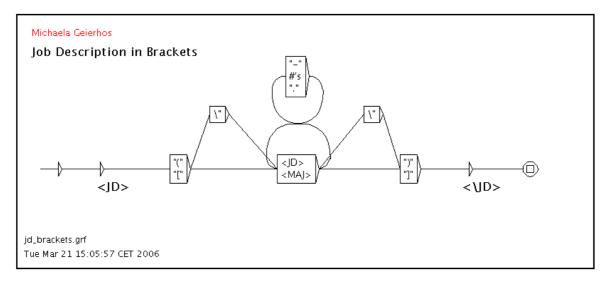
- 115. William "Bill" H Gates III, KBE
- 116. William "Bill" H. Gates III, KBE
- 117. William "Bill" H Gates, KBE
- 118. William "Bill" H. Gates, KBE
- 119. William Henry Gates
- 120. William Henry Gates III
- 121. William Henry Gates, III
- 122. William Henry Gates III, KBE
- 123. William Henry Gates, KBE
- 124. William H Gates
- 125. William H. Gates
- 126. William H Gates III
- 127. William H Gates, III
- 128. William H. Gates III
- 129. William H. Gates, III
- 130. William H Gates III, KBE
- 131. William H. Gates III, KBE
- 132. William H Gates, KBE
- 133. William H. Gates, KBE
- 134. ...

 $<sup>^{45} \</sup>mathrm{vgl.\,http://seattletimes.nwsource.com/html/editorialsopinion/2001845028\_billed28.html}$ 

# D Weitere Berufsbezeichnergraphen



**Abbildung D.1:** Graph zur Erkennung vorangestellter Berufsbezeichnerattribute -  $jd\_introduction.grf$ 



**Abbildung D.2:** Graph zur Erkennung von geklammerten Berufsbezeichnungen -  $jd\_brack$ ets.qrf

### Die Subgraphen des Automaten "jd\_complement.grf"

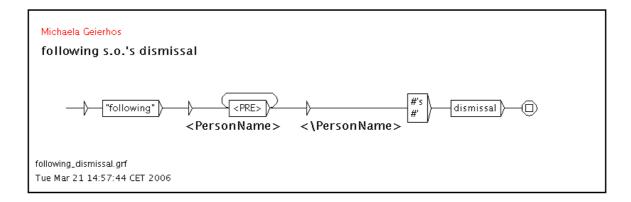
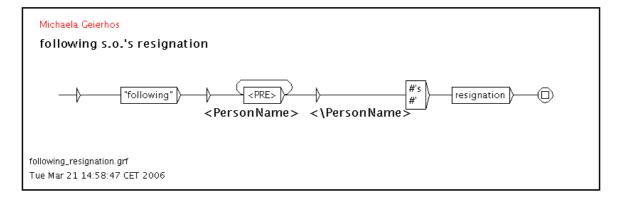
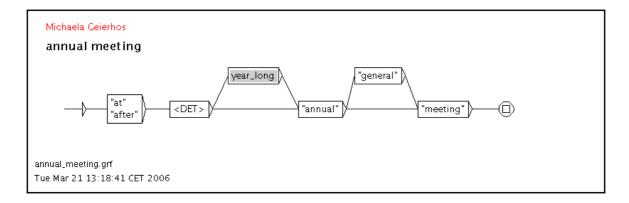


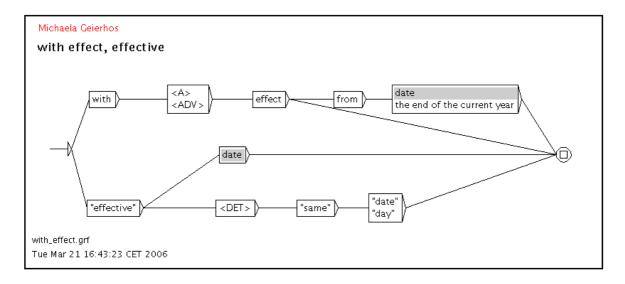
Abbildung D.3: Graph zur Erkennung von entlassenen Personen - following\_dismissal.grf



 $\textbf{Abbildung D.4:} \ \text{Graph zur Erkennung von abgedankten oder zur ückgetretenen Personen} - following\_resignation.grf$ 

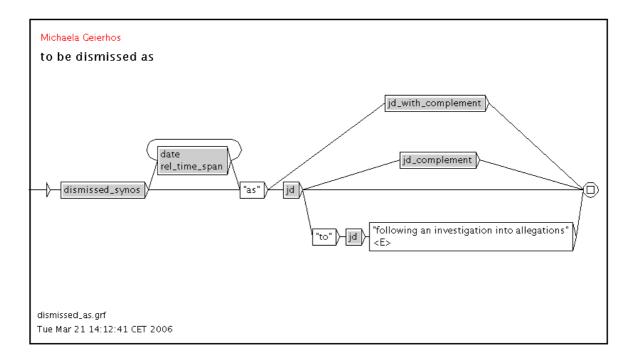


 $\textbf{Abbildung D.5:} \ \textbf{Graph zur Erkennung von Jahreshauptversammlungen -} \ annual\_meeting.grf$ 

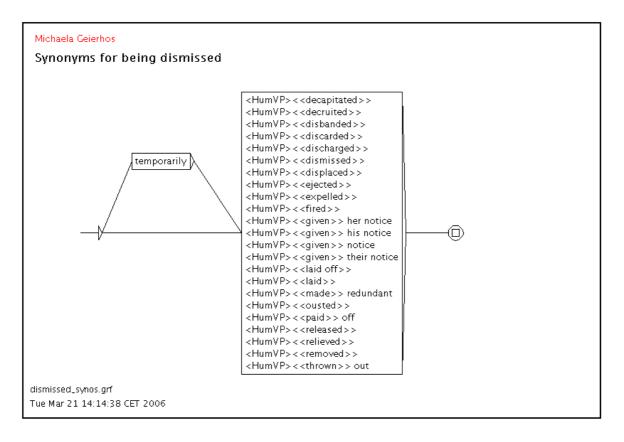


**Abbildung D.6:** Graph zur Erkennung von Zeitpunkten, ab denen etwas gültig ist -  $with\_ef$ -fect.grf

# E Graphen und Konkordanz zur Verbalphrase "to be dismissed"



**Abbildung E.1:** Graph zur Erkennung von Verbalphrasen mit dem Verb "to be dismissed as" und seiner direkten Synonyme - dismissed\_as.grf



**Abbildung E.2:** Graph zur Erkennung von Verben mit der Bedeutung von "to be dismissed as" - dismissed\_synos.grf

Dabei behandelt der Graph dismissed\_synos.grf folgende Passivverbkonstruktionen:

- to be decapitated
- to be decruited
- to be disbanded
- to be discarded
- to be discharged
- to be dismissed
- to be displaced
- to be ejected
- to be expelled
- to be fired

- to be given one's notice
- to be laid (off)
- to be made redundant
- to be ousted
- to be paid off
- to be released
- to be relieved
- $\bullet$  to be removed
- to be thrown out

```
176
```

Ephraim Mashaba, who was dismissed as <GEO>South Africa<\GEO>'s <JD>coach<\JD> just before the Nations Cup, has taken over at top division club Ann Wilson had been dismissed as <JD>managing director<\JD> of BBC Technology, its IT services division Vladimir Zelezny, who was dismissed as <ORG>Nova<\ORG>'s <JD>director-general<\JD> in 2003.{S} since being dismissed as the <JD>player-manager<\JD> of Bury.{S} Bob Francis was fired as <JD>coach<\JD> of the struggling Phoenix Coyotes on Tuesday almost four years after being fired as its <JD>basketball coach<\JD>.{S} a week after Michael Sears was fired as the aerospace giant's <JD> chief financial officer<\JD> after questions surfaced Michael Green, only weeks after being ousted as <JD>chairman<\JD> of Carlton, is planning to bid for the 330-screen UCI business has been vacant since Michael Green was ousted as <JD>chairman<\JD> - designate during a shareholder revolt in October Last October Michael Green was ousted as <JD>chairman<\JD> when he was ousted as <JD>chief executive<\JD> of Britain's second-largest hotel group, Thistle, in a torrent of publicity.{S} Iain Lumsden, who was ousted as <JD>Chief Executive<\JD> and replaced by his deputy Sandy Crombie earlier this year, received GBP 1.13 million in salary ROBIN SAUNDERS, who was ousted as <JD>head<\JD> of the German bank WestLB's principal finance business last year Messier was ousted as <JD>head<\JD> of Vivendi in the summer of 2002 Greg Dyke, ousted as <ORG>BBC<\ORG> <JD>director general<\JD> last week, has decided to rein in his public attacks on the Hutton report in favour of a more measured response. {S} Steve Case was removed as <JD>chairman<\JD> of AOL Time Warner Inc. at least partly because of pressure from Gordon Crawford, a media fund manager at Capital Research & Management Co. that Pol Maj-Gen Kosin Hinthao be removed as <JD>chief<\JD> of the Crime Suppression Division.{S} but were not legally removed as <JD>directors<\JD>.{S} Mr Horace Okumu, be removed as <JD>prosecutor<\JD> for allegedly violating their right to fair hearing.{S}

Abbildung E.3: Konkordanz zum Graphen dismissed\_as.qrf

### Literaturverzeichnis

- [1] E. AGICHTEIN AND L. GRAVANO, Snowball: Extracting relations from large plaintext collections, in Proceedings of the Fifth ACM International Conference on Digital Libraries, 2000, pp. 85–94. http://www.cs.columbia.edu/~eugene/papers/dl00.pdf.
- [2] H. ALANI, S. KIM, D. E. MILLARD, M. J. WEAL, P. H. LEWIS, W. HALL, AND N. R. SHADBOLT, Automatic extraction of knowledge from web documents, in Proceedings of 2nd International Semantic Web Conference Workshop on Human Language Technology for the Semantic Web abd Web Services, Sanibel Island, Florida, USA, 2003. http://eprints.ecs.soton.ac.uk/8194/01/Alani-HLT03-final.pdf.
- [3] R. Barzilay, Information Fusion for Multidocument Summarization: Paraphrasing and Generation, PhD thesis, Columbia University, 2003. http://www1.cs.columbia.edu/nlp/theses/regina\_barzilay.ps.
- [4] BIOGRAPHY.COM, 2005/2006. http://www.biography.com.
- [5] O. BLANC AND A. DISTER, Automates lexicaux avec structure de traits, in RE-CITAL 2004, Fès, 21 avril 2004. http://valibel.fltr.ucl.ac.be/Elements/BLANC\_DISTER-RECITAL\_2004.pdf.
- [6] I. Blank, Computerlinguistische Analyse mehrsprachiger Fachtexte, CIS-Bericht-98-109., cis-bericht-98-109, Centrum für Informations- und Sprachverarbeitung der Ludwig-Maximilians-Universität München, 1997. http://www.ifi.unizh.ch/~volk/TermDB/Blank\_Diss\_CIS-Bericht-98-109.pdf.
- [7] J. CALLAN, Human language technologies: Information extraction. Lecture Notes, Fall 2005. http://hartford.lti.cs.cmu.edu/classes/11-682/Lectures/16-InfoExtractionB.pdf.
- [8] CAREERBUILDER.COM, 2005. http://www.careerbuilder.com.
- [9] M. CIARAMITA AND Y. ALTUN, Named-Entity Recognition in Novel Domains with External Lexical Knowledge, in Workshop on Advances in Structured Learning for Text and Speech Processing, NIPS, 2005. http://www.cis.upenn.edu/~crammer/workshop\_material/ciaramita\_altun\_structlearn.pdf.

- [10] M. Constant, Description d'expressions numériques en français, in Revue Informatique et Statistique dans les Sciences humaines 36, Actes des troisièmes journées INTEX, A. Dister, ed., Liège, 2000, pp. 119–135.
- [11] M. CONSTANT, Grammaires locales pour l'analyse automatique de textes : Méthodes de construction et outils de gestion, PhD thesis, Université de Marne la Vallée, 2003. http://www-igm.univ-mlv.fr/~mconstan/papers/constant\_these.pdf.
- [12] M. Constant, *GRAAL*, une bibliothèque de graphes : mode d'emploi, in Cahiers de la MSH Ledoux 1, INTEX pour la linguistique et le traitement automatique des langues, C. Muller, J. Royeauté, and M. Silberztein, eds., Besançon, 2004, Presse Universitaire de Franche-Comté, pp. 321–330.
- [13] B. Courtois, *Dictionnaires électroniques DELAF anglais et français*, in Lexique, syntaxe et lexique-grammaire; syntax, lexis & lexicon-grammar, C. Leclère, Éric Laporte, M. Piot, and M. Silberztein, eds., John Benjamins Publishing Company, 2004, pp. 113–123.
- [14] A. D. Cruse, Lexical Semantics, Cambridge University Press, 1986.
- [15] A. CUCCHIARELLI AND P. VELARDI, Unsupervised named entity recognition using syntactic and semantic contextual evidence, Computational Linguistics, Volume 27 (2001). http://acl.ldc.upenn.edu/J/J01/J01-1005.pdf.
- [16] S. CUCERZAN AND D. YAROWSKY, Language independent named entity recognition combining morphological and contextual evidence, in Proceedings of the Joint SIGDAT Conference on EMNLP and VLC, 1999, pp. 90–99. http://acl.ldc.upenn.edu/W/W99/W99-0612.pdf.
- [17] L. Danlos and M. Gross, Building Electronic Dictionaries for Natural Language Processing, in Proceedings of the 2nd Symposium France-Japan, L. Kott, ed., Amsterdam: North Holland, 1988.
- [18] DoPL, Division of Professional Licensure List of Professions. Online-Verzeichnis, 2005. http://www.mass.gov/dpl/boards/dirprofs.htm.
- [19] DOT, Dictionary Of Occupational Titles. Online-Verzeichnis, 2005. http://www.occupationalinfo.org/.
- [20] G. Drosdowski, W. Müller, W. Scholze-Stubenrecht, and M. Wermke, eds., *DUDEN - Rechtschreibung der deutschen Sprache*, Dudenverlag, Mannheim, 1996. 21. Auflage. Neue Rechtschreibung.
- [21] P. Duboué and K. McKeown, Statistical Acquisition of Content Selection Rules for Natural Language Generation, in Proceedings of the 2003 Conference on Empirical Methods for Natural Language Processing (EMNLP 2003), Sapporo, Japan, July 2003. http://www.cs.columbia.edu/~pablo/publications/EMNLP2003selection.pdf.

- [22] P. Duboué, K. McKeown, and V. Hatzivassiloglou, *ProGenIE: Biographical descriptions for intelligence analysis*, in Proceedings of the NSF/NIJ Symposium on Intelligence and Security Informatics, vol. 2665 of Lecture Notes in Computer Science, Tucson, Arizona, USA, June 2003, Springer-Verlag, pp. 343–345. http://www.cs.columbia.edu/~pablo/publications/ISI2003biographies.pdf.
- [23] C. Fairon, Structures non-connexes. Grammaire des incises en français : description linguistique et outils informatiques, PhD thesis, Université Paris 7, 2000.
- [24] C. Fellbaum, ed., WordNet An Electronic Lexical Database, MIT Press, 1998. http://mitpress.mit.edu/catalog/item/default.asp?ttype=2&tid=8106.
- [25] N. Friburger, Reconnaissance automatique des nomes propres Application à la classification automatique de textes journalistiques, PhD thesis, Université François Rabelais, Tours, 2002.
- [26] N. Friburger and D. Maurel, Elaboration d'une cascade de transducteurs pour l'extraction des noms personnes dans les textes, in TALN 2001, Tours, 2-5 juillet 2001. http://www.up.univ-mrs.fr/veronis/Atala/TALN/pdf/art16\_p183\_192.pdf.
- [27] R. GAIZAUSKAS, Information extraction an information extraction perspective on text mining: Tasks, technologies and prototype applications. Euromap Text Mining Seminar, September 4, 2002. http://www.itri.brighton.ac.uk/projects/euromap/Text%20Mining%20Event/Rob\_Gaizauskas.pdf.
- [28] M. GEIERHOS, Einführung in WordNet 2.0. http://www.cis.uni-muenchen.de/~micha/old/WordNet.pdf, 3. Februar 2004.
- [29] M. GEIERHOS, DELA Wörterbücher: Der Umgang mit externen Ressourcen in Unitex. Was man beim Erstellen eigener Lexika beachten sollte? http://www.cis.uni-muenchen.de/~micha/old/Dela.pdf, 2. Mai 2005.
- [30] GOVERNMENT OF NEWFOUNDLAND AND LABRADOR, Canadian job classification. Online-Verzeichnis, 2005. http://www.fin.gov.nl.ca/fin/pensions/pdf/employer\_specs/job\_class.pdf.
- [31] M. Gross, On the relations between syntax and semantics, in Formal Semantics of Natural Language, E. L. Keenan, ed., Cambridge: Cambridge University Press, 1975, pp. 389–405.
- [32] M. Gross, Remarks on the separation between syntax and semantics, in Studies in Descriptive and Historical Linguistics, Festschrift for Winfred P. Lehmann, Amsterdam/Philadelphia, 1977, Benjamins, pp. 71–81.
- [33] M. Gross, *Taxonomy in syntax*, SMIL, Journal of Linguistic Calculus 1978:3-4, (1978), pp. 73–96. Stockholm: Skriptor.

- [34] M. Gross, On the failure of generative grammar, Language 55:4, (1979), pp. 859–885.
- [35] M. Gross, On structuring the lexicon, Quaderni di Semantica 4:1, (1983), pp. 107–120.
- [36] M. Gross, Lexicon-Grammar: The Representation of Compound Words, in COLING-1986 Proceedings, 1986, pp. 1–6.
- [37] M. GROSS, Methods and Tactics in the Construction of a Lexicon-Grammar, Linguistics in the Morning Calm 2, Selected papers from SICOL 1986, (1988), pp. 177–197. Seoul: Hanshin.
- [38] M. Gross, Linguistic representations and text analysis, Linguistic Unity and Linguistic Diversity in Europe, (1991), pp. 31–61. London: Academia Europaea.
- [39] M. Gross, The argument structure of elementary sentences, Language Research 28:4, (1992), pp. 699–716. Seoul National University.
- [40] M. Gross, Local grammars and their representation by finite automata, in Data, Description, Discourse, Papers on the English Language in honour of John McH Sinclair, M. Hoey, ed., Harper-Collins, London, 1993, pp. 26–38.
- [41] M. Gross, Constructing lexicon-grammars, in Computational Approaches to the Lexicon, Atkins and Zampolli, eds., Oxford Univ. Press, 1994, pp. 213–263.
- [42] M. Gross, *Lexicon grammar*, in Concise encyclopedia of syntactic theories, K. Brown and J. Miller, eds., Pergamon, 1996, pp. 244–258.
- [43] M. Gross, *The Construction of Local Grammars*, in Finite-State Language Processing, E. Roche and Y. Schabès, eds., Language, Speech, and Communication, Cambridge, Mass.: MIT Press, 1997, pp. 329–354.
- [44] M. Gross, Lemmatization of compound tenses in English, Lingvisticae Investigationes, XXII (1998-1999), pp. 71–122.
- [45] M. GROSS, A bootstrap method for constructing local grammars, in Contemporary Mathematics: Proceedings of the Symposium, University of Belgrad, Belgrad, 1999, pp. 229–250.
- [46] M. GROSS AND J. SENELLART, Nouvelles bases pour une approche statistique, in JADT98, Nice, France, 1998.
- [47] F. GUENTHNER, Electronic lexica and corpora research at cis, International Journal of Corpus Linguistics, 1 (1996), pp. 287–301.
- [48] GUIDE TO THE WORLD OF OCCUPATIONS, Alphabetical list of occupations, 2005. http://www.occupationsguide.cz/en/abecedni/abecedni.htm.

- [49] J. E. HOPCROFT, R. MOTWANI, AND J. D. ULLMAN, Einführung in die Automatentheorie, Formale Sprachen und Komplexitätstheorie, Addison-Wesley Pearson Studium, 2002.
- [50] LABOURMARKET. Online-Verzeichnis, 2005. http://www.labourmarket.co.nz/labourmarket.htm.
- [51] J. LACOMBE, List of occupations. Online-Verzeichnis, 2004. http://www.cpcug.org/user/jlacombe/terms.html.
- [52] S. LANGER, Selektionsklassen und Hyponymie im Lexikon, CIS-Bericht 96-94, Centrum für Informations- und Sprachverarbeitung, Ludwig-Maximillians-Universität Mnchen, 1996.
- [53] E. LAPORTE, Graphes paramétrés et lexique-grammaire, in Interface lexique-grammaire et lexiques syntaxiques et sémantiques, 12 mars 2005. http://www.atala.org/doc/JE\_050312/Lexsynt-Laporte.pdf.
- [54] P. Maier-Meyer, Lexikon und automatische Lemmatisierung, CIS-Bericht 95-84, Centrum für Informations- und Sprachverarbeitung, Ludwig-Maximillians-Universität München, 1995.
- [55] F. Mallchok, Automatic Recognition of Organization Names in English Business News, PhD thesis, Ludwig-Maximillians-Universität München, 2004.
- [56] I. Mani, Recent developments in text summarization, in Proceedings of the Tenth International Conference on Information and Knowledge Management, Atlanta, Georgia, USA, November 5-10 2001, pp. 529-531. http://complingone.georgetown.edu/~linguist/papers/cikm.pdf.
- [57] I. Mani, K. Concepcion, and L. V. Guilder, *Using summarization for automatic briefing generation*, in NAACL-ANLP 2000 Workshop: Automatic Summarization, The MITRE Corporation, 2000. http://acl.ldc.upenn.edu/W/W00/W00-0410.pdf.
- [58] G. S. Mann, Fine-grained proper noun ontologies for question answering, tech. rep., Department of Computer Science, John Hopkins University, Baltimore, Maryland, 2001. http://www.cs.ust.hk/~hltc/semanet02/pdf/mann.pdf.
- [59] MAPPLANET GMBH, MapPlanet.com, 2006. http://mapplanet.com/ix/.
- [60] A. McCallum and D. Jensen, A note on the unification of information extraction and data mining using conditional-probability, relational models, in IJ-CAI'03 Workshop on Learning Statistical Models from Relational Data, 2003. http://www.cs.umass.edu/~mccallum/papers/iedatamining-ijcaiws03.pdf.
- [61] A. McCallum and W. Li, Early Results for Named Entity Recognition with Conditional Random Fields, Feature Induction and Web-Enhanced Lexicons, in Seventh Conference on Natural Language Learning (CoNLL), 2003. http://www.cs.umass.edu/~mccallum/papers/mccallum-conll2003.pdf.

- [62] T. MIKOLAJEWSKI, Eine Untersuchung der Formen und Konstruktionen von Menschenbezeichnern für das Elektronische Lexikonsystem CISLEX, Studien zur Informations- und Sprachverarbeitung Band 7, Centrum für Informations- und Sprachverarbeitung, Ludwig-Maximillians-Universität Mnchen, 2003.
- [63] G. A. MILLER, R. BECKWITH, C. FELLBAUM, D. GROSS, AND K. MILLER, Introduction to WordNet: An On-line Lexical Database, Cognitive Science Laboratory, Princeton University, 1993. http://www.isi.edu/isd/kr/5papers.pdf.
- [64] E. MILLS, Google balances privacy, reach, CNET News.com, (2005). http://news.com.com/Google+balances+privacy%2C+reach/2100-1032\_3-5787483.html.
- [65] MOM, Ministry of Manpower List of Occupations. Online-Verzeichnis, 2003. http://www.mom.gov.sg/MOM/MRSD/Others/2003W\_OccList.pdf.
- [66] G. NAVARRO, R. BAEZA-YATES, AND J. M. ARCOVERDE, *Matchsimile: A flexible approximate matching tool for personal names searching*, in Proceedings SBBD '01, 2001, pp. 228-242. http://www.dcc.uchile.cl/~gnavarro/ps/sbbd01.ps.gz.
- [67] A. NENKOVA AND K. MCKEOWN, Improving the coherence of multi-document summaries: A corpus study for modeling the syntactic realization of entities, tech. rep., Columbia University, 2003. http://www.cs.columbia.edu/techreports/cucs-001-03.pdf.
- [68] A. NENKOVA AND K. MCKEOWN, References to named entities: A corpus study, in Proceedings of NAACL-HLT, 2003. http://acl.ldc.upenn.edu/N/N03/N03-2024.pdf.
- [69] C. Niu, W. Li, J. Ding, and R. K. Srihari, Bootstrapping for named entity tagging using concept-based seeds, in HLT-NAACL, 2003. http://acl.ldc.upenn.edu/N/N03/N03-2025.pdf.
- [70] OOH, Occupational outlook handbook. Online-Verzeichnis, 1998. http://www.umsl.edu/services/govdocs/ooh9899/1.htm.
- [71] OOH, Occupational outlook handbook. Online-Verzeichnis, 2000-2001. http://www.umsl.edu/services/govdocs/ooh20002001/1.htm.
- [72] S. PATEL AND J. SMARR, Automatic Classification of Previously Unseen Proper Noun Phrases into Semantic Categories Using an N-Gram Letter Model, in CS224N/Ling237 Final Projects, Stanford University, 2001. http://nlp.stanford.edu/courses/cs224n/2001/jsmarr/NGramWordClassifier.pdf.
- [73] S. Paumier, Some remarks on the application of a lexicon-grammar, Lingvisticæ Investigationes XXIV:2, (2001), pp. 245–256. Amsterdam/Philadelphia, John Benjamins.

- [74] S. PAUMIER, De la reconnaissance de formes linguistiques à l'analyse syntaxique, PhD thesis, Université de Marne-la-Vallée, 2003.
- [75] S. PAUMIER, Manuel d'utilisation d'Unitex, 2004. http://wwwigm.univmlv.fr/~unitex/.
- [76] PROSPECTS.AC.UK, 2005. http://www.prospects.ac.uk.
- [77] E. RILOFF AND R. JONES, Learning dictionaries for information extraction by multi-level bootstrapping, in Proceedings of the Sixteen National Conference on Artificial Intelligence (AAAI-99), 1999, pp. 1044-1049. http://www.cs.cmu.edu/afs/cs/project/cald/www/lis/jones\_learning\_dictionaries.ps.
- [78] E. Roche, Analyse syntaxique transformationnelle du français par transducteurs et lexique-grammaire, PhD thesis, Université Paris 7 Denis Diderot, janvier 1993.
- [79] J. ROTH, Der Stand der Kunst in der Eigennamen-Erkennung: Mit einem Fokus auf Produktenamen-Erkennung, Master's thesis, Universität Zürich, 2002. http://www.ifi.unizh.ch/cl/study/lizarbeiten/lizjeannetteroth.pdf.
- [80] B. Schiffman, I. Mani, and K. J. Concepcion, *Producing biographical sum-maries: Combining linguistic knowledge with corpus statistics*, in Proceedings of the 39<sup>th</sup> Annual Meeting of the Association for Computational Linguistics, Toulouse, France, 2001, pp. 450–457. http://acl.ldc.upenn.edu//P/P01/P01-1059.pdf.
- [81] J. Senellart, Locating noun phrases with finite state transducers, in Proceedings of the 17th International Conference on Computational Linguistics, Montréal, 1998, pp. 1212–1219. http://acl.ldc.upenn.edu/P/P98/P98-2198.pdf.
- [82] J. Senellart, Tools for locating noun phrases with finite state transducers, in The computational treatment of nominals. Proceedings of the Workshop, COLING-ACL'98, 1998, pp. 80–84. http://acl.ldc.upenn.edu/W/W98/W98-0611.pdf.
- [83] J. Senellart, Outils de reconnaissance d'expressions linguistiques complexes dans des grands corpus, PhD thesis, Université Paris 7 Denis Diderot, 1999.
- [84] M. Silberztein, Dictionnaire électroniques et analyse automatique de textes le système INTEX, 1993. Paris, Masson.
- [85] K. SPARCK-JONES, What might be in a summary?, in Information Retrieval '93: Von der Modellierung zur Anwendung, G. Knorz, J. Krause, and C. Womser-Hacker, eds., Universitätssverlag Konstanz, 1993, pp. 9-26. http://www.ftp.cl.cam.ac.uk/ftp/papers/ksj/ksj-whats-in-a-summary.ps.gz.
- [86] SpecialistInfo.com, Consultants, 2006. http://www.specialistinfo.com/directory.php.
- [87] L. SWEENEY, Finding Lists of People on the Web, ACM Computers and Society, 34 (2004). http://privacy.cs.cmu.edu/dataprivacy/projects/rosterfinder/rosterfinder2.pdf.

- [88] Systran, Document de référence, 2003. http://www.systransoft.com/company/investors/AnnualReport2003-FR.pdf.
- [89] H. Traboulsi, A local grammar for proper names, Master's thesis, University of Surrey, August 2004. http://portal.surrey.ac.uk/pls/portal/url/ITEM/F3E6A8EF87037385E0340003BA296BDE.
- [90] H. N. TRABOULSI, Towards the automatic acquisition of local grammars, in Proceedings of the 3<sup>rd</sup> Annual PhD Conference, Department of Computing, University of Surrey, UK, 2005, pp. 13–18. http://portal.surrey.ac.uk/pls/portal/url/ITEM/FAD86365EBCD5E15E0340003BA296BDE.
- [91] O. TSUR, M. DE RIJKE, AND K. SIMA'AN, Biographer: Biography questions as a restricted domain question answering task, in Proceedings ACL 2004 Workshop on Question Answering in Restricted Domains, 2004. http://www.science.uva.nl/~mdr/Publications/Files/acl2004-qa-ws.pdf.
- [92] U. S. DEPARTMENT OF LABOR, Office of Administrative Law Judges. Internet Law Library, März 2005. http://www.oalj.dol.gov/printdoc.htm?URL= %2Fpublic%2Fdot%2Frefrnc%2Fdotalpha.htm.
- [93] UBC, The University of British Columbia Library -1991 Standard Occupational Classification (SOC). Online-Verzeichnis, 1991. http://data.library.ubc.ca/datalib/restricted/other/statscan/soc/alph1.html.
- [94] UNIVERSITY OF SURREY, UK, Hayssam Traboulsi BSc, MSc, Februar 2006. http://portal.surrey.ac.uk/pls/portal/PORTAL.wwv\_media.show?p\_id= 211329&p\_settingssetid=726846&p\_settingssiteid=0&p\_siteid=798&p\_type=basetext&p\_textid=211330.
- [95] USDoL, U.S. Department of Labor Alphabetical List of SOC Occupations. Online-Verzeichnis, 2001/2002. http://www.bls.gov/oes/2001/oes\_alph.htm, http://www.bls.gov/oes/2002/oes\_alph.htm.
- [96] M. Waite, ed., Oxford Thesaurus of English, Oxford University Press, 2004.
- [97] WIKIPEDIA, The Free Encyclopedia, 2005/2006. http://en.wikipedia.org/wiki/Main\_Page.
- [98] WORDNET, An Electronic Lexical Database, Version 2.1. http://wordnet.princeton.edu/.
- [99] L. ZHOU, M. TICREA, AND E. HOVY, Multi-document biography summarization, in Proceedings of EMNLP, 2004, pp. 434-441. http://arxiv.org/pdf/cs.CL/0501078.
- [100] ZOOM INFORMATION INC., ZoomInfo People, Companies, Relationships, 2006. http://www.zoominfo.com/Search/.

## Abbildungsverzeichnis

3.1	HealthAndN.grf aus Gross, 1999 [45]	25
3.2	<i>VModToV.grf</i> aus Gross, 1998-1999 [44]	26
3.3	<i>Insert.grf</i> aus Gross, 1998-1999 [44]	26
3.4	Übersicht der Lemmatisierungsgraphen aus Gross, 1998-1999 [44] (Teil 1)	27
3.5	Übersicht der Lemmatisierungsgraphen aus Gross, 1998-1999 [44] (Teil 2)	28
3.6	MinisterOccupation.grf aus Senellart, 1998 [81]	32
3.7	FullName.grf aus Senellart, 1998 [81]	33
3.8	NounPhrases.grf aus Senellart, 1998 [81]	34
5.1	Auszug aus dem Lexikon FirstNamesdic	49
5.2	Auszug aus dem Lexikon LastNamesdic	50
5.3	Auszug aus dem Lexikon LongNamesBiosdic	52
5.4	Auszug aus dem Lexikon LongNamesFTdic	54
5.5	Auszug aus dem Lexikon LongNamesSpecialistInfodic	54
5.6	Auszug aus dem Lexikon LongNamesZoominfo[12345]dic	55
5.7	Auszug aus dem Lexikon LongNamesAuthorsdic	55
5.8	Auszug aus dem Lexikon <i>Titlesdic</i>	56
5.9	Auszug aus dem Lexikon MenbezWordnetdic	59
5.10	Auszug aus dem Lexikon <i>JDdic</i>	60
5.11	Auszug aus dem Lexikon Citizensdic	61
5.12	Auszug aus dem Lexikon $HumVP$ $dic$	62
5.13	Auszug aus dem Lexikon Disciplinesdic	63
5.14	Auszug aus dem Lexikon Sectordic	64
5.15	Auszug aus dem Lexikon orgbezdic	65
5.16	Auszug aus dem Lexikon orgdic	66
5.17	Auszug aus dem Lexikon Companiesdic	66
5.18	Auszug aus dem Lexikon org_adjdic	67
5.19	Auszug aus dem Lexikon context_beforedic	67
5.20	Auszug aus dem Lexikon Geos Wikipediadic	69
5.21	Auszug aus dem Lexikon GeosMapplanetdic	69
5.22	Auszug aus dem Lexikon Geos Wikipediadic	70
	Auszug aus dem Lexikon <i>USstatesdic</i>	70
5.24	Auszug aus dem Lexikon Geos Wikipediadic	71
	Auszug aus dem Lexikon <i>Monthdic</i>	72
	Auszug aus dem Lexikon MonthAbbrdic	72
	Auszug aus dem Lexikon DayOfWeekdic	73

	Auszug aus dem Lexikon $Ntimedic$	73 74
6.1 6.2 6.3 6.4 6.5 6.6	Graph zur Erkennung von Personennamen - person_name.grf Graph zur Erkennung potentieller Personennamen - pot_person_name.grf Graph zur Erkennung von Menschenbezeichnern - person.grf	78 80 82 83 84 85
7.1 7.2	Graph zur Erkennung von Firmennamen - $company.grf$ Graph zur Erkennung von Firmennamen in geklammerten Ausdrücken - $company\_brackets.grf$	87 88
7.3 7.4	Graph zur Erkennung von Rechtsformen und weiteren Firmenzusätzen - $company\_additions.grf$	89 90
8.1 8.2	Graph zur Erkennung von Toponymen - $geo.grf$	93 94
9.1 9.2 9.3	Graph zur Erkennung genauer Datumsangaben aus Gross, 1993 [40] Graph zur Erkennung von Datumsangaben - date.grf	95 96 97
9.3 9.4 9.5	Erster Pfad aus dem Graphen date.grf	97 97 98
9.6 9.7	Konkordanz zum zweiten Pfad aus dem Graphen $date.grf$	98 99 99
<ul><li>9.8</li><li>9.9</li><li>9.10</li></ul>	Konkordanz zum vierten Pfad aus dem Graphen $date.grf$	100 100
9.12	Konkordanz zum sechsten Pfad aus dem Graphen $date.grf$	
9.14	Zehnter Pfad aus dem Graphen $date.grf$	102
9.16	Graph zur Erkennung von numerischen Tagesangaben - $day.grf$	
9.18	Graph zur Erkennung von vierstelligen Jahreszahlen - $year\_long.grf$ Graph zur Erkennung von zweistelligen Jahreszahlen - $year\_short.grf$	104
	Graph zur Erkennung von ausgewählten persönlichen Relationen - merger_personal_relations.grf	105
	Graph zur Erkennung von Verbalphrasen mit dem Verb "to be born" und seiner Paraphrasierung "to see the light of day" - born.grf	

10.4 Graph zur Erkennung von Verbalphrasen mit dem Verb "to be raised (up)"	
und seiner direkten Synonyme - raised.grf	112
	113
10.6 Graph zur Erkennung von Verbalphrasen mit dem Verb "to graduate" und	
seiner direkten Synonyme - graduated.grf	114
10.7 Graph zur Erkennung von Altersangaben - age.grf	115
10.8 Graph zur Erkennung von Verben mit der Bedeutung von "to graduate"	
	116
10.9 Graph zur Erkennung von Abschlussbezeichnungen für den akademischen	
Grad des Bachelors - $bachelor.grf$	117
10.10Graph zur Erkennung von Abschlussbezeichnungen für den akademischen	
Grad des Masters - $master.grf$	118
10.11Graph zur Erkennung von Abschlussbezeichnungen für den akademischen	
Grad des Doktors - $doctor.grf$	119
10.12Graph zur Erkennung von Verbalphrasen mit dem Verb "to marry so." in	
seiner Aktiv- und Passivform, sowie seiner direkten Synonyme - married.grf 1	20
10.13Graph zur Erkennung von Verben mit der Bedeutung von "to marry, to	
be $married$ " - $married\_synos.grf$	122
$10.14 \text{Graph}$ zur Erkennung von relativen Zeiträumen - $\textit{rel\_time\_span.grf}$ 1	123
10.15Graph zur Erkennung von Verbalphrasen mit dem Verb to be divorced	
und seiner direkten Synonyme - divorced.grf	124
10.16Graph zur Erkennung von Verben mit der Bedeutung von "to be divorced"	
- $divorced\_synos.grf$	125
10.17 Konkordanz zum Graphen $\operatorname{died.grf}$	126
$10.18$ Graph zur Erkennung von Verbalphrasen mit dem Verb $to\ die$ - $died.grf$ $ 1$	127
$10.19 {\rm Graph\ zur\ Erkennung\ von\ m\"{o}glichen\ Todesursachen\ -\ {\it cause\_of\_death.grf}}  1$	128
10.20Graph zur Erkennung des Verbs to die und seiner direkten Synonyme -	
$died\_synos.grf$	129
11.1 Graph zur Erkennung von ausgewählten beruflichen Relationen - mer-	20
	L30
11.2 Graph zur Erkennung von Verbalphrasen mit dem Verb "to be appointed	124
	134
	135
11.4 Graph zur Erkennung von Verben mit der Bedeutung von "to be appoin-	126
	136
	L37
11.6 Graph zur Erkennung der Komplemente von Berufsbezeichnungen (rechten Ventaut eines Berufsbezeichnung)	120
ter Kontext eines Berufsbezeichners) - jd_complement.grf	.38
11.7 Graph zur Erkennung von Berufsbezeichnern mitsamt ihrer rechten Kon-	190
	L39
11.8 Graph zur Erkennung potentieller Sektoren- und Branchenbezeichnungen	10
1 0 0	L40
11.9 Graph zur Erkennung von Verbalphrasen mit dem Verb "to employ s.o." und seiner direkten Synonyme - employ.arf	L41
nna semer anekren avnonvine - <i>enthion.att</i>	41

11.10	OGraph zur Erkennung von Verben mit der Bedeutung von "to employ	
	$s.o.$ " - $employ\_synos.grf$	142
11.11	1 Graph zur Erkennung von Verbalphrasen mit dem Verb "to join" und	
	seiner Paraphrasierung "to become a member of" - joined.grf	144
11.12	2Graph zur Erkennung von Verbalphrasen mit dem Verb "to be employed	
	$as$ " und seiner direkten Synonyme - $employed\_as.grf$	145
11.13	BGraph zur Erkennung von Verben mit der Bedeutung von "to be employed	
	$as$ " - $employed\_synos.grf$	146
11.14	4Graph zur Erkennung von Verbalphrasen mit dem Verb "to be paid as"	
	und seiner Paraphrasierung "to draw salary as" - paid_as.grf	147
11.15	5Graph zur Erkennung von Verbalphrasen mit dem Verb "to work as" und	
		149
11.16	* * *	150
	7Graph zur Erkennung von Verbalphrasen mit dem Verb "to dismiss s.o."	
		151
11.18	8Graph zur Erkennung von Verben mit der Bedeutung von "to dismiss	
		152
11.19	9Graph zur Erkennung von Verbalphrasen mit dem Verb "to be replaced	-
		153
11.20	OGraph zur Erkennung von Verbalphrasen mit dem Verb "to resign (as/from)"	
	und seiner direkten Synonyme - $resigned.grf$	
11.21	1Graph zur Erkennung von Verbalphrasen mit dem Verb "to retire s.o."	
	und seiner direkten Synonyme - retire_as.grf	157
11.22	2Graph zur Erkennung von Verbalphrasen mit dem Verb "to be retired as"	10.
		157
	and somet are not spring the power as gry 1.1.1.1.1.1.1.1.1.1.1.1.1.1.1.1.1.1.1.	101
13.1	Graph zur automatischen Erkennung weiterer Verbkonstruktionen, die	
	eine Beziehung zwischen einer Person und einer Firma ausdrücken	160
13.2	Graph zur automatischen Erkennung weiterer Verbkonstruktionen, die	
	eine Beziehung zwischen mindestens zwei Personen ausdrücken	161
13.3	Graph zur automatischen Erkennung weiterer Verbkonstruktionen, die	
	eine Beziehung zwischen einer Personen und ihrem Beruf ausdrücken	162
D 4		
D.1	Graph zur Erkennung vorangestellter Berufsbezeichnerattribute - $jd_{-}intro$	
		171
D.2	Graph zur Erkennung von geklammerten Berufsbezeichnungen - $jdbrack$ -	
		171
D.3		172
D.4	Graph zur Erkennung von abgedankten oder zurückgetretenen Personen	
		172
D.5	Graph zur Erkennung von Jahreshauptversammlungen - $annual\_meeting.grf$	173
D.6	Graph zur Erkennung von Zeitpunkten, ab denen etwas gültig ist - with_ef-	
	fect.grf	173
T7. 1		
E.1	Graph zur Erkennung von Verbalphrasen mit dem Verb "to be dismissed	1 17 4
	as" und seiner direkten Synonyme - dismissed_as.grf	1/4

E.2	Graph zur Erkennung von Verben mit der Bedeutung von "to be dismissed	
	$as$ " - $dismissed\_synos.grf$	175
E.3	Konkordanz zum Graphen dismissed_as.qrf	176

## **Tabellenverzeichnis**

3.1	Statistische Auswertung der Extraktionsresultate aus Friburger (2001) [26]	39
5.1	Abkürzungen, die in FirstNamesdic verwendet werden	49
5.2	Abkürzungen, die in LastNamesdic verwendet werden	50
5.3	Abkürzungen, die in LongNames*dic verwendet werden	51
5.4	Abkürzungen, die in <i>Titlesdic</i> verwendet werden	56
5.5	Abkürzungen, die in MenbezWordnetdic verwendet werden	59
5.6	Abkürzungen, die in <i>JDdic</i> verwendet werden	60
5.7	Abkürzungen, die in <i>Citizensdic</i> verwendet werden	61
5.8	Abkürzungen, die in <i>Disciplinesdic</i> verwendet werden	63
5.9	Abkürzungen, die in Sectordic verwendet werden	64
5.10	Abkürzungen, die in <i>Companiesdic</i> verwendet werden	67
5.11	Abkürzungen, die in Geos*dic für Nomina verwendet werden	68
5.12	Abkürzungen, die in $Geos*dic$ für Adjektive verwendet werden	71
12.1	Statistische Auswertung der Extraktionsresultate	159

## Index

Master 117, 118 Akronym 42, 88 Aktivform 120, 124, 143, 145, 152, 156 Allgemeine Menschenbezeichner 81 Altersangabe 115 Anaphern 84 Anglizismus 44 Anredeform 53, 56, 57 Antonym 123 Antonymie 58 Arbeitgeber 147, 148 Arbeitsbeziehung 148 Arbeitsbeziehung 148 Arbeitsverhältnis 132, 145, 162 Aufwachsen 111 Australian Province Citizen 61, 166 Auswertung 158 Autobiographie 12 Automat 16  Biographische Persönlich Biography.com Biogra	e Relation . 12, 43, 46, 47, 1, 105, 120, 160, 163 e Relation . 130, 132, 141 he Relation 12, 62, 105, 126, 130 Relation 13, 62 he Relation 13, 62 he Relation
BCISLEX-EBeruf143Citizen	

City Adjective       71, 166         Clustering       36, 39         Company Name       67, 166         Continent       68, 166         County       68, 166	Evaluation
D	Fachbereich
Département	Fachbereichlexikon
DAG	Fachrichtung         116           Fairon, Cédrick         37
Datum	Fall-Out
Datumsangabe 72, 76, 99, 101	Financial Times
Datumsangaben 95	Firmeneintritt
Datumserkennung 72	Firmenkürzel
DELA	Firmenname 41, 65, 67, 92, 166
DELA Wörterbücher 20	Firmensitz
DELAC 21, 37	First Name
DELACF 21	Flexionsform
DELAF 21	Friburger, Nathalie
DELAF-L 22	FT Korpus 47, 111, 135, 159
DELAF-M 22	Full Name 51
DELAF-S	
DELAS 21, 37	<b>G</b>
Directed Acyclic Graph	Geburt
Disambiguierung 15, 16, 19, 79, 92	Geburtsdatum
Diskursanalyse	Geburtsname
Drogenmissbrauch	Geburtsort
E	Geburtstag
Ehejahre	Geographical Term
Ehepartner	Geographical Term Adjective 71, 166
Eigenname 11, 35, 45, 49–51, 53, 67,	Geographischer Begriff 91, 92
76, 80, 81, 84, 167	Gerichteter Azyklischer Graph 17
Eigennamenerkennung 10, 11, 40	Google
Einfache Wörter 21	Grafschaft 68, 70, 166
Einstellung 141, 143	grammatikalische Lexikoninformation
Einstellungsdatum 143	A 71
Einwohner	N 49-51, 56, 59-61, 63, 64, 67, 68
Einwohnerlexikon	XA 71
Elektronisches Lexikon	XN 56, 60, 63, 64, 67, 68
Eltern	Graph 16, 17
Entität 10, 11, 15, 42, 43, 45, 65, 76,	Gross, Maurice 19, 21, 24, 26, 42, 62, 95
86, 105, 159, 160 Entlassung	н
Ernennung	Harris, Zellig Sabbetai
Ernennungsrelation	Heirat
Erziehung 111	Herzversagen
	120

Himmelsrichtung 92	Krankheit 128
Holonymie 58	
Homograph 30	L
HPL 41	LADL 19, 21, 22
Human 49–51, 59–61, 166	Land 68, 69, 129, 167
Hyperonym 43	Lehrbereich
Hyperonymie	Lehreinrichtung 114, 116
Hyponym 35	Lemma 20
Hyponymie 58	Lemmaform
	Lemmatisierung 26, 62
I	Lexikalische Analyse
IE 10	Lexikon
Information Retrieval 44	Lexikon-Grammatik 22
Information Retrieval Methoden 30	Lexikoneintrag
Exact-Pattern-Algorithm 30	Lexikongrammatik 16, 19
Exact-String-Matching-Algorithm 30	Lexikonkodierung 22
Key-Word-Algorithm 30	Lexikonpriorität
Statistical Algorithm 30	Lokale Grammatik 15–19, 22, 23, 40,
Information-Retrieval-System 158	41, 43–45, 47, 75, 76, 95, 105,
Informationsextraktion 10, 44	111, 114, 120, 126, 132, 153,
INTEX 17, 19, 37	159, 163
	Lokativa
J	Long Name 51, 167
Jahresangabe	,
Jahresspanne	M
Jahresversammlung	Mallchok, Friederike . 40, 43, 46, 58, 66
Jahreszahl	Maurel, Denis
Jahreszeiten	Mehrwortlexem . 22, 56, 60, 63, 64, 67,
К	68, 167
Kanadische Provinz 68, 166	Menschenbezeichner 15, 22, 43, 46,
Kanonische Form	49–51, 59–62, 76, 84, 105, 120,
Kaskadierung	141, 143, 147, 152, 156,
Kategorietyp . 49–51, 56, 59–61, 63, 64,	160–163, 166
67, 68, 71, 167	Menschenbezeichnerlexika 58
grammatikalische Kategorie 49	Meronymie
semantische Kategorie 49	Metropole 68, 71, 166
Kindheit	Modularität
Klassifikation	Monatsabkürzung 72, 98
Komplexe Wörter	Monatsangabe
Konkordanz	Monatsname
Kontextfreie Grammatik	Morphologie
Kontinent	Mots Composés
Kontraktionsauflösung	Mots Simples
Korpusverarbeitung	MUC-7 11
Krankenhaus	Mustererkennung
1x1amxommaus 123	manufathaning

N	New York City Bourough 68, 167
Nachfolge 153, 154	New York City Citizen 61, 167
Nachfolgerelation	New York Times 41
Nachname 50, 80, 167	Nominalphrase 105, 136, 137
einfacher Nachname 50	Normalisierung 20, 47
komplexer Nachname 50	
Nachnamenlexikon 50	0
Nagel, Sebastian 46	ODL 41, 65
Named Entity 10–12	ONL 41, 65
Named-Entity-Recognition . 10, 11, 40,	Organisation 43, 86, 88, 90, 141, 143,
43, 44	147, 152, 153, 160
Namen der Lexika	Organisationsbeschreibungslexikon . 41,
Citizensdic 60	65
Companiesdic	Organisationsname . 40, 41, 65, 67, 76,
Disciplinesdic	79, 86, 88, 92, 137, 154, 160, 166
FirstNamesdic 49, 51	Organisationsnamenlexikon 41, 65
GeosMapplanetdic 68, 69	Ortsangabe
GeosSebastian+.dic 68	Ortsbezeichnung 41
GeosWikipediadic 68–70	Р
HumVPdic	•
JDdic 59	Pars Pro Toto
LastNamesdic 50	Passivform 120, 124, 143, 145, 147,
LongNamesAuthorsdic 51, 55	152, 156
LongNamesBiosdic 51, 52	Paumier, Sébastien
LongNamesFTdic	Pensionierung
LongNamesSpecialistInfodic 51, 53	Performanz
LongNamesZoominfo1dic 51, 54	Person
LongNamesZoominfo2dic 51, 54	Personenbezogene Prädikate 62, 160,
LongNamesZoominfo3dic 51, 54	161, 163
LongNamesZoominfo4dic 51, 54	Personenlexikon
LongNamesZoominfo5dic 51, 54	Personenname 51, 53, 55, 56, 76, 79,
MenbezWordnetdic	86, 88, 90, 143, 167
Monthdic	Personennamenlexika 51, 80
MonthAbbr.dic	Personensuchmaschine
Sectordic	Prädikat 107, 133
Titlesdic	Precision
USstatesdic	Prenom-prolex
context_beforedic	Prioritätsebene
orgdic	Produktname
org_adjdic	Prolintex 37
0 0	Proper Noun 49–51, 67, 167
0	Province Adjective 71, 166
Nation	R
Nation Adjective	
NER	Recall
Neubesetzung 153	Region 68, 70, 167

Rekursiver Aufruf 97	MonthAbbr 72
Relative Jahresangabe 99	Month 72
Relativsatz 105	NYCBourough 68, 167
Rentenalter	NYCcitizen 61, 167
Ressourcen	Nation 68, 167
Reuters Korpus 41, 44, 46	Ntime 73
Reuters News	PR 49–51, 67, 167
Roth, Jeannette 45	Region 68, 167
	SN 50, 167
S	Sector 63, 64, 167
Satzenderkennung 20, 37, 47	Title $56, 167$
Satzendmarkierung	USstateCitizen 61, 167
Scheidung 107, 123, 125	USstate 68, 167
Scheinname	Urbanite 61, 167
Schulabschluss	Senellart, Jean
Schule	SpecialistInfo.com 53
Sektor 63, 64, 140, 167	Sprachgebundenheit 44
Sektorenlexikon	Städtenamen 69
Semantik	Staatsbürger 61, 166
semantische Lexikoninformation	Stadt
ABourough 71, 166	Stadtbewohner 61, 167
ACity 71, 166	Stadtteil 68, 166
AGEO 71, 166	State Adjective 71, 166
ANation	Sterben 126
AProvince 71, 166	Straßenname
AState 71, 166	Subgraph 17
Abbrev 56, 166	Subsprache
AuProvinceCitizen 61, 166	Suchfenster 160
Borough 68, 166	Surname 50, 167
${\tt CaProvinceCitizen}\ \dots \dots\ 61,\ 166$	Sweeney, Latanya 55
CaProvince 68, 166	Synekdoche 77
Citizen 61, 166	Synonym
City 68, 166	Synonymie 58, 125
Company 67, 166	Synset 58
Continent	Syntaktische Variabilität 14, 76, 86
County 68, 166	Syntax 17, 19, 22
${\tt DayOfWeek}~\dots~73$	
Departement 68, 166	Т
Discipline 63, 166	Tätigkeit
FN 49, 166	Tagesangabe 98, 101
GEO 68, 166	Teilformenlexikon
HumVP 62	Temporalialexikon 72
Hum $49-51, 59-61, 166$	Textkodierung 19
JD 60, 166	Titel 53, 56, 167
LN 51, 167	Adelstitel 57

akademischer Grad	. 53, 56	Verbalphrasenlexikon 62
aristokratischer Titel	. 56, 57	Verben
militärischer Rang	53, 56	einfache Verben 62
Titellexikon	56	komplexe Verben 62
to be appointed	133, 141	Verifikation 75
to be born	107, 109	Vollform 77
to be dismissed	151	Vollständigkeit
to be divorced	123	Vorname
to be employed	143	einfacher Vorname 49
to be married	120	komplexer Vorname 49
to be paid	147	Vornamenlexikon 49
to be raised (up)	107, 111	
to be replaced	153	W
to die	126	Wörterbuch
to dismiss so		Wörterbuch-Lookup 91
to employ so	141	Wörterbucheintrag 21, 22, 43, 51, 52,
to graduate	114	54, 56, 75, 86, 90, 140, 159
to join		Wall Street Journal 41
to resign		Wikipedia 49, 50, 56, 59, 60, 63
to work		Wirtschaftsnachrichten 15, 44, 46, 47
Tod	107, 126	Wochentag
Todesalter	,	Wochentaglexikon 73
Todesdauer		WordNet 58
Todesort		Z
Todesursache		
Todeszeitpunkt	,	Zeitangabe
Tokenisierung 20		Zeitbestimmung
Tokenizer-Programm		Zeitbezogene Nomina
Toponym 11, 37, 43, 45, 68, 69		Zeitspanne 100
79, 91, 92, 166	, . , ,	ZoomInfo.com 54
Toponymlexikon	68	
Transduktor 16, 23, 36, 3'		
Transitionsnetz		
Trennung	,	
U		
Unfall	128	
unique beginners		
Unitex 17, 19, 23, 3'		
US Bundesstaat 68		
US State 68.		
US State Citizen	61, 167	
V		
Verbalphrase 106, 107, 109, 1	132–134	
143, 145, 148, 156	,	