

Roland Hausser

Foundations of Computational Linguistics

Human-Computer
Communication
in Natural Language

Third Edition

 Springer

Foundations of Computational Linguistics

Roland Hausser

Foundations of Computational Linguistics

Human-Computer Communication
in Natural Language

Third Edition

Roland Hauser
Abteilung für Computerlinguistik (retired)
Friedrich-Alexander-Universität
Erlangen-Nürnberg
Erlangen, Germany

ISBN 978-3-642-41430-5

ISBN 978-3-642-41431-2 (eBook)

DOI 10.1007/978-3-642-41431-2

Springer Heidelberg New York Dordrecht London

Library of Congress Control Number: 2013956746

© Springer-Verlag Berlin Heidelberg 1999, 2001, 2014

This work is subject to copyright. All rights are reserved by the Publisher, whether the whole or part of the material is concerned, specifically the rights of translation, reprinting, reuse of illustrations, recitation, broadcasting, reproduction on microfilms or in any other physical way, and transmission or information storage and retrieval, electronic adaptation, computer software, or by similar or dissimilar methodology now known or hereafter developed. Exempted from this legal reservation are brief excerpts in connection with reviews or scholarly analysis or material supplied specifically for the purpose of being entered and executed on a computer system, for exclusive use by the purchaser of the work. Duplication of this publication or parts thereof is permitted only under the provisions of the Copyright Law of the Publisher's location, in its current version, and permission for use must always be obtained from Springer. Permissions for use may be obtained through RightsLink at the Copyright Clearance Center. Violations are liable to prosecution under the respective Copyright Law.

The use of general descriptive names, registered names, trademarks, service marks, etc. in this publication does not imply, even in the absence of a specific statement, that such names are exempt from the relevant protective laws and regulations and therefore free for general use.

While the advice and information in this book are believed to be true and accurate at the date of publication, neither the authors nor the editors nor the publisher can accept any legal responsibility for any errors or omissions that may be made. The publisher makes no warranty, express or implied, with respect to the material contained herein.

Printed on acid-free paper

Springer is part of Springer Science+Business Media (www.springer.com)

Preface to the Third Edition

The Third Edition has been gently modernized, leaving the overview of traditional, theoretical, and computational linguistics, analytic philosophy of language, and mathematical complexity theory with their historical backgrounds intact. The format of the empirical analyses of English and German syntax and semantics has been adapted to current practice. Sects. 22.3–24.5 have been rewritten to focus more sharply on the construction of a talking robot.

Because the Second Edition is quite long, one aim of the present edition has been to shorten the number of pages. Little is lost, however, because later publications make up for a reduction of about 10 percent. The broad, leisurely style induced by the book's origin as material for a four semester introductory lecture course, however, has been preserved.

Special thanks to Ronan Nugent, Senior Editor of Computer Science at Springer, for his continued support.

Erlangen, Germany
June 2013

Roland Hausser

Remark on the Use of *Risci*

Science books often use a menagerie of tables, graphs, examples, definitions, theorems, lemmata, schemata, etc., which are either (i) not numbered or (ii) use different counters for each kind. For the author, (i) makes it difficult to refer to a particular item. For the reader, (ii) makes it difficult to find the item.

As a remedy, we use a uniform container called *riscus*¹ to hold any of such items. A *riscus* begins with a numbered head line and ends when normal text resumes. In this book, the number of a *riscus* consists of three parts, e.g., 3.2.1, which specify the location in terms of the chapter within the text, the section within the chapter, and the *riscus* within the section.

Independently of the *riscus* numbering, the head line may classify and number the item contained, as in the following fictitious example:

3.2.1 THEOREM 9: THE POLYNOMIAL COMPLEXITY OF CF LANGUAGES

Content of the theorem followed by the proof.

Elsewhere in the text, the above *riscus* may be referred to with something like As shown in 3.2.1, ... or In Theorem 9 (3.2.1) In this way, the *riscus* numbering allows the author to be brief yet precise because the reader may be guided to a more detailed presentation given in another place. The method works regardless of whether the book is a hard copy or an eBook.

The numbering of a *riscus* may be provided automatically, for example, by using the \LaTeX `\subsection` command for *riscus* headings. The *riscus* numbering, like that of chapters and sections, may be hooked up with automatic cross-referencing, for example, by adding `\usepackage{hyperref}` to the \LaTeX preamble of the text.

As a result, the reader of an eBook may click on a three part number to be transported automatically to the *riscus* in question, and may get back to the earlier position with a simple command like `cntrl [`. The comfort of accessing a *riscus* directly, instead of having to leaf through the chapter and the section as in a hard copy, is perhaps not the least advantage added by the transition from a hard copy to an eBook.

Erlangen, Germany

Roland Hausser

¹ *riscus* *<i>m*, is the Latin word for suitcase. A *riscus* was made of woven willow branches covered with fur.

Preface to the First Edition

The central task of a future-oriented computational linguistics is the development of cognitive machines which humans can freely talk with in their respective natural language. In the long run, this task will ensure the development of a functional theory of language, an objective method of verification, and a wide range of applications.

Natural communication requires not only verbal processing, but also non-verbal recognition and action. Therefore the content of this textbook is organized as a theory of language for the construction of talking robots. The main topic is the *mechanism of natural language communication* in both the speaker and the hearer.

The content is divided into the following parts:

- I. Theory of Language
- II. Theory of Grammar
- III. Morphology and Syntax
- IV. Semantics and Pragmatics

Each part consists of 6 chapters. Each of the 24 chapters consists of 5 sections. A total of 797 exercises help in reviewing key ideas and important problems.

Part I begins with current applications of computational linguistics. Then it describes a new theory of language, the functioning of which is illustrated by the robot CURIOUS. This theory is referred to with the acronym SLIM, which stands for *Surface compositional Linear Internal Matching*. It includes a cognitive foundation of semantic primitives, a theory of signs, a structural delineation of the components syntax, semantics, and pragmatics, as well as their functional integration in the speaker's utterance and the hearer's interpretation. The presentation refers to other contemporary theories of language, especially those of Chomsky and Grice, as well as to the classic theories of Frege, Peirce, de Saussure, Bühler, and Shannon and Weaver, explaining their formal and methodological foundations as well as their historical background and motivations.

Part II presents the theory of *formal grammar* and its methodological, mathematical, and computational roles in the description of natural languages. A description of Categorical (C) Grammar and Phrase Structure (PS) Grammar is combined with an introduction to the basic notions and linguistic motivation of generative grammar. Further topics are the declarative vs. procedural aspects of parsing and gener-

ation, type transparency, as well as the relation between formalisms and complexity classes. It is shown that the principle of possible *substitutions* causes empirical and mathematical problems for the description of natural language. As an alternative, the principle of possible *continuations* is formalized as LA Grammar. LA stands for the left-associative derivation order which models the time-linear nature of language. Applications of LA Grammar to relevant artificial languages show that its hierarchy of formal languages is orthogonal to that of PS Grammar. Within the LA hierarchy, natural language is in the lowest complexity class, namely the class of C1 languages which parse in linear time.

Part III describes the *morphology* and *syntax* of natural language. A general description of the notions word, word form, morpheme, and allomorph, the morphological processes of inflection, derivation, and composition, as well as the different possible methods of automatic word form recognition is followed by the morphological analysis of English within the framework of LA Grammar. Then the syntactic principles of valency, agreement, and word order are explained within the left-associative approach. LA Grammars for English and German are developed by systematically extending a small initial system to handle more and more constructions such as the fixed vs. free word order of English and German, respectively, the structure of complex noun phrases and complex verbs, interrogatives, subordinate clauses, etc. These analyses are presented in the form of explicit grammars and sample derivations.

Part IV describes the *semantics* and *pragmatics* of natural language. The general description of language interpretation begins by comparing three different types of semantics, namely those of logical languages, programming languages, and natural languages. Based on Tarski's foundation of logical semantics and his reconstruction of the Epimenides paradox, the possibility of applying logical semantics to natural language is investigated. Alternative analyses of intensional contexts, propositional attitudes, and the phenomenon of vagueness illustrate that different types of semantics are based on different ontologies which greatly influence the empirical results. It is shown how a semantic interpretation may cause an increase in complexity and how this is to be avoided within the SLIM theory of language. The last two chapters, 23 and 24, analyze the interpretation by the hearer and the conceptualization by the speaker as a time-linear navigation through a database called *word bank*. A word bank allows the storage of arbitrary propositional content and is implemented as a highly restricted specialization of a classic (i.e., record-based) network database. The autonomous navigation through a word bank is controlled by the explicit rules of suitable LA Grammars.

As supplementary reading the *Survey of the State of the Art in Human Language Technology*, Ron Cole (ed.) 1998 is recommended. This book contains about 90 contributions by different specialists giving detailed snapshots of their research in language theory and technology.

Contents

Part I. Theory of Language

1. Computational Analysis of Natural Language	3
1.1 Human-Computer Communication	4
1.2 Language Sciences and Their Components	7
1.3 Methods and Applications of Computational Linguistics	12
1.4 Electronic Medium in Recognition and Synthesis	13
1.5 Second Gutenberg Revolution	16
<i>Exercises</i>	22
2. Smart vs. Solid Solutions	25
2.1 Indexing and Retrieval in Textual Databases	25
2.2 Using Grammatical Knowledge	29
2.3 Smart vs. Solid Solutions in Computational Linguistics	31
2.4 Beginnings of Machine Translation	33
2.5 Machine Translation Today	37
<i>Exercises</i>	42
3. Cognitive Foundations of Semantics	45
3.1 Prototype of Communication	45
3.2 From Perception to Recognition	47
3.3 Iconicity of Formal Concepts	50
3.4 Context Propositions	56
3.5 Recognition and Action	59
<i>Exercises</i>	62
4. Language Communication	65
4.1 Adding Language	65
4.2 Modeling Reference	68
4.3 Using Literal Meaning	71
4.4 Frege's Principle	73
4.5 Surface Compositionality	76
<i>Exercises</i>	84

5. Using Language Signs on Suitable Contexts	87
5.1 Bühler’s Organon Model	87
5.2 Pragmatics of Tools and Pragmatics of Words	89
5.3 Finding the Correct Subcontext	91
5.4 Language Production and Interpretation	94
5.5 Thought as the Motor of Spontaneous Production	97
<i>Exercises</i>	99
6. Structure and Functioning of Signs	103
6.1 Reference Mechanisms of Different Sign Kinds	103
6.2 Internal Structure of Symbols and Indexicals	107
6.3 Repeating Reference	110
6.4 Exceptional Properties of Icon and Name	114
6.5 Pictures, Pictograms, and Letters	118
<i>Exercises</i>	121

Part II. Theory of Grammar

7. Formal Grammar	127
7.1 Language as a Subset of the Free Monoid	127
7.2 Methodological Reasons for Formal Grammar	132
7.3 Adequacy of Formal Grammars	133
7.4 Formalism of C Grammar	134
7.5 C Grammar for Natural Language	138
<i>Exercises</i>	141
8. Language Hierarchies and Complexity	143
8.1 Formalism of PS Grammar	144
8.2 Language Classes and Computational Complexity	146
8.3 Generative Capacity and Formal Language Classes	149
8.4 PS Grammar for Natural Language	154
8.5 Constituent Structure Paradox	159
<i>Exercises</i>	163
9. Basic Notions of Parsing	167
9.1 Declarative and Procedural Aspects of Parsing	167
9.2 Fitting Grammar onto Language	169
9.3 Type Transparency Between Grammar and Parser	174
9.4 Input-Output Equivalence with the Speaker-Hearer	180
9.5 Desiderata of Grammar for Achieving Convergence	183
<i>Exercises</i>	185
10. Left-Associative Grammar (LAG)	189
10.1 Rule Kinds and Derivation Order	189
10.2 Formalism of LA Grammar	193

10.3 Time-Linear Analysis 196

10.4 Absolute Type Transparency of LA Grammar 198

10.5 LA Grammar for Natural Language 201

Exercises 207

11. Hierarchy of LA Grammar 209

11.1 Generative Capacity of Unrestricted LAG 209

11.2 LA Hierarchy of A, B, and C LAGs 212

11.3 Ambiguity in LA Grammar 215

11.4 Complexity of Grammars and Automata 218

11.5 Subhierarchy of C1, C2, and C3 LAGs 221

Exercises 227

12. LA and PS Hierarchies in Comparison 231

12.1 Language Classes of LA and PS Grammar 231

12.2 Subset Relations in the Two Hierarchies 233

12.3 Nonequivalence of the LA and PS Hierarchy 235

12.4 Comparing the Lower LA and PS Classes 237

12.5 Linear Complexity of Natural Language 239

Exercises 245

Part III. Morphology and Syntax

13. Words and Morphemes 249

13.1 Words and Word Forms 249

13.2 Segmentation and Concatenation 254

13.3 Morphemes and Allomorphs 257

13.4 Categorization and Lemmatization 258

13.5 Methods of Automatic Word Form Recognition 262

Exercises 266

14. Word Form Recognition in LA Morph 269

14.1 Allo Rules 269

14.2 Phenomena of Allomorphy 273

14.3 Left-Associative Segmentation into Allomorphs 279

14.4 Combi Rules 282

14.5 Concatenation Patterns 285

Exercises 289

15. Corpus Analysis 291

15.1 Implementation and Application of Grammar Systems 291

15.2 Subtheoretical Variants 295

15.3 Building Corpora 297

15.4 Distribution of Word Forms 300

15.5 Statistical Tagging	304
<i>Exercises</i>	308
16. Basic Concepts of Syntax	311
16.1 Delimitation of Morphology and Syntax	311
16.2 Valency	314
16.3 Agreement	316
16.4 Free Word Order in German (<i>LA D1</i>)	319
16.5 Fixed Word Order in English (<i>LA E1</i>)	324
<i>Exercises</i>	326
17. LA Syntax for English	329
17.1 Complex Fillers in Pre- and Postverbal Position	329
17.2 English Field of Referents	334
17.3 Complex Verb Forms	336
17.4 Finite State Backbone of LA Syntax (<i>LA E2</i>)	339
17.5 Yes/No-Interrogatives (<i>LA E3</i>) and Grammatical Perplexity	343
<i>Exercises</i>	348
18. LA Syntax for German	351
18.1 Standard Procedure of Syntactic Analysis	351
18.2 German Field of Referents (<i>LA D2</i>)	354
18.3 Verbal Positions in English and German	359
18.4 Complex Verbs and Elementary Adverbs (<i>LA D3</i>)	362
18.5 Interrogatives and Subordinate Clauses (<i>LA D4</i>)	368
<i>Exercises</i>	374
<hr/>	
Part IV. Semantics and Pragmatics	
<hr/>	
19. Three Kinds of Semantics	379
19.1 Basic Structure of Semantic Interpretation	379
19.2 Logical, Programming, and Natural Languages	381
19.3 Functioning of Logical Semantics	383
19.4 Metalanguage-Based or Procedural Semantics?	388
19.5 Tarski's Problem for Natural Language Semantics	391
<i>Exercises</i>	395
20. Truth, Meaning, and Ontology	397
20.1 Analysis of Meaning in Logical Semantics	397
20.2 Intension and Extension	400
20.3 Propositional Attitudes	403
20.4 Four Basic Ontologies	407
20.5 Sorites Paradox and the Treatment of Vagueness	410
<i>Exercises</i>	415

21. Absolute and Contingent Propositions	417
21.1 Absolute and Contingent Truth	417
21.2 Epimenides in a [+sense, +constructive] System	421
21.3 Frege’s Principle as Homomorphism	424
21.4 Time-Linear Syntax with Homomorphic Semantics	428
21.5 Complexity of Natural Language Semantics	431
<i>Exercises</i>	434
22. Database Semantics	437
22.1 Database Metaphor of Natural Communication	437
22.2 Descriptive Aporia and Embarrassment of Riches	440
22.3 Combining Categorical Operation and Semantic Interpretation	444
22.4 Reference as Pattern Matching	447
22.5 Repercussion of the Semantic Interpretation on the DBS Syntax	449
<i>Exercises</i>	452
23. Semantic Relations of Structure	455
23.1 Coding Content at the Elementary, Phrasal, and Clausal Levels	455
23.2 Storing Content in a Word Bank	457
23.3 Representing the Semantic Relations of Structure Graphically	461
23.4 Paraphrase and the Universalist/Relativist Dichotomy	464
23.5 The Ten SLIM States of Cognition	467
<i>Exercises</i>	473
24. Conclusion	475
24.1 Hear Mode	475
24.2 Speak Mode	476
24.3 Questions and Answers	479
24.4 Autonomous Control	483
24.5 Coherence	486
<i>Exercises</i>	488
Bibliography	491
Name Index	503
Subject Index	507

Abbreviations Referring to Preceding and Subsequent Work

Preceding Work

- SCG Hausser, R. (1984) *Surface Compositional Grammar*, München: Wilhelm Fink Verlag, pp. 274
- NEWCAT Hausser, R. (1986) *NEWCAT: Natural Language Parsing Using Left-Associative Grammar*, Lecture Notes in Computer Science, Vol. 231, Springer, pp. 540
- CoL Hausser, R. (1989a) *Computation of Language: An Essay on Syntax, Semantics and Pragmatics in Natural Man-Machine Communication*, Symbolic Computation: Artificial Intelligence, Springer, pp. 425
- TCS Hausser, R. (1992) “Complexity in Left-Associative Grammar,” *Theoretical Computer Science* 106.2:283–308

Subsequent Work

- AIJ Hausser, R. (2001c) “Database Semantics for Natural Language.” *Artificial Intelligence* 130.1:27–74
- NLC Hausser, R. (2006) *A Computational Model of Natural Language Communication – Interpretation, Inferencing, and Production in Database Semantics*, Springer, pp. 365
- L&I Hausser, R. (2010) “Language Production Based on Autonomous Control – A Content-Addressable Memory for a Model of Cognition,” *Language and Information* 11:5–31
- CLaTR Hausser, R. (2011) *Computational Linguistics and Talking Robots, Processing Content in Database Semantics*, Springer, pp. 286

Introduction

I. BASIC GOAL OF COMPUTATIONAL LINGUISTICS

Transmitting information by means of a natural language like Chinese, English, or German is a real and well structured procedure. This becomes evident when we attempt to communicate with people who speak a foreign language. Even if the information we want to convey is completely clear to us, we will not be understood by our hearers if we fail to use their language adequately. Conversely, even if our foreign language partners use their language as they always do, and without any problems, we will not understand them.

The goal of computational linguistics is to reproduce the natural transmission of information by modeling the speaker's production and the hearer's interpretation on a suitable type of computer. This amounts to the construction of autonomous cognitive machines (robots) which can freely communicate in natural language.

The development of speaking robots is not a matter of fiction, but a real scientific task. Remarkably, however, theories of language have so far avoided a functional modeling of the natural communication mechanism, concentrating instead on peripheral aspects such as methodology (behaviorism), innate ideas (nativism), and scientific truth (model theory).

II. TURING TEST

The task of modeling the mechanism of natural communication on the computer was described in 1950 by ALAN TURING (1912–1954) in the form of an 'imitation game' known today as the Turing test. In this game, a human interrogator is asked to question a male and a female partner in another room via a teleprinter in order to determine which answer was given by the man and which by the woman.¹ The people running the test count how often the interrogator classifies his communication partners correctly and how often he is wrong.

Then one of the two humans is replaced by a computer. The computer passes the Turing test if the man or the woman replaced by the computer is simulated so well that the guesses of the human interrogator are just as often right and wrong as with

¹ As a corresponding situation in real life consider an author wondering about the gender of the anonymous copyeditor courteously provided by the publisher.

the earlier natural partner. In this way Turing wanted to replace the question “Can machines think?” with the question “Are there imaginable digital computers which would do well in the imitation game?”

III. ELIZA PROGRAM

In its original intention, the Turing test requires the construction of an artificial cognitive agent with a language behavior so natural that it cannot be distinguished from that of a human. This presupposes complete coverage of the language data and of the communicative functions in real time. At the same time, the test tries to avoid all aspects not directly involved in human language behavior.²

However, the Turing test does not specify what cognitive structure the artificial agent should have in order to succeed in the imitation game. It is therefore possible to misinterpret the aim of the Turing test as fooling the human interrogator rather than providing a functional model of communication on the computer. This was shown by the Eliza program of Weizenbaum (1965).

The Eliza program simulates a psychiatrist encouraging the human interrogator, now in the role of a patient, to talk more and more about him- or herself. Eliza works with a list of words. Whenever one of these words is typed by the human interrogator/patient, Eliza inserts it into one of several prefabricated sentence templates. For example, when the word *mother* is used by the human, Eliza uses the template *Tell me more about your ___* to generate the sentence *Tell me more about your mother*.

Because of the way in which Eliza works, we know that Eliza has no understanding of the dialog with the interrogator/patient. Thus, the construction of Eliza is not a model of communication. If we regard the dialog between Eliza/psychiatrist and the interrogator/patient as a modified Turing test, however, the Eliza program is successful insofar as the interrogator/patient *feels* him- or herself understood and therefore does not distinguish between a human and an artificial communication partner in the role of a psychiatrist.

The Eliza program is the prototype of a *smart* solution (Sect. 2.3) in that it exploits the restrictions of a highly specialized application to achieve a maximal effect with a minimum of effort. The goal of computational linguistics, however, is a *solid* solution in science: it must (i) explain the mechanism of natural communication theoretically and (ii) verify the theory with an implementation (software machine) which may be loaded with the language-dependent lexicon and compositional operations of any natural language. The speak and the hear mode of the implementation must work in any practical application for which free natural language communication between humans and machines is desired (15.1.3).

² As an example of such an aspect, Turing (1950), p. 434, mentions the artificial recreation of human skin.

IV. MODELING NATURAL COMMUNICATION

Designing a talking robot provides an excellent opportunity for systematically developing the basic notions as well as the philosophical, mathematical, grammatical, methodological, and programming aspects of computational linguistics. This is because modeling the mechanism of natural communication requires

- a theory of language which explains the natural transfer of information in a way that is functionally coherent and mathematically explicit,
- a description of language data which is empirically complete for all components of this theory of language, i.e., the lexicon, the morphology, the syntax, and the semantics, as well as the pragmatics and the representation of the internal context, and
- a degree of precision in the description of these components which supports a straightforward computational implementation running in real time.

Fulfilling these requirements will take hard, systematic, goal-oriented work, but it will be worth the effort.

For theory development, the construction of talking robots is of interest because an electronically implemented model of communication may be tested both in terms of the language behavior observed externally, and internally via direct access to its cognitive states via the service channel. The work towards realizing unrestricted human-computer communication in natural language is facilitated by the fact that the functional model may be developed incrementally, beginning with a simplified, but fully general system to which additional functions as well as additional natural languages are added step by step.

For practical purposes, unrestricted communication with computers and robots in natural languages will make the interaction with these machines maximally user friendly and permit new, powerful ways of information processing. The use of artificial programming languages may then be limited to specialists developing and servicing the machines.

V. USING PARSERS

Computational linguistics analyzes natural languages automatically in terms of software programs called parsers. The use of parsers influences the theoretical viewpoint of linguistic research, distribution of funds, and everyday research practice as follows:

- *Competition*

Competing theories of grammar are measured with respect to the new standard of how well they are suited for efficient parsing and how well they fit into a theory of language designed to model the mechanism of natural language communication.

- *Funding*

Adequate parsers for different languages are needed for an unlimited range of practical applications, which has a major impact on the inflow of funds for research, development, and teaching in this particular area of the humanities.

– *Verification*

Programming grammars as parsers allows testing their empirical adequacy automatically on arbitrarily large amounts of real data in the areas of word form recognition and synthesis, syntactic analysis and generation, and the semantic-pragmatic interpretation in both the speak and the hear mode.

The verification of a theory of language and grammar by means of testing an electronic model in real applications is a new approach which clearly differs from the methods of traditional linguistics, psychology, philosophy, and mathematical logic.

VI. THEORETICAL LEVELS OF ABSTRACTION

So far there are no electronic systems which model the functioning of natural communication so successfully that one can talk with them more or less freely. Furthermore, researchers do not agree on how the mechanism of natural communication really works. One may therefore question whether achieving a functional model of natural communication is possible in principle. I would like to answer this question with an analogy³ from the recent history of science.

Today's situation in computational linguistics resembles the development of mechanical flight before 1903.⁴ For hundreds of years humans had observed sparrows and other birds in order to understand how they fly. Their goal was to become airborne in a similar manner. It turned out, however, that flapping wings did not work for humans. This was taken by some as a basis for declaring human flight impossible in principle, in accordance with the pious cliché "If God had intended humans to fly, He would have given them wings."⁵

³ See also CoL, p. 317.

⁴ In 1903, the brothers Orville and Wilbur Wright succeeded with the first manned motorized flight.

⁵ Irrational reasons against a modeling of natural communication reside in the subconscious fear of creating artificial beings resembling humans and having superhuman powers. Such *homunculi*, which occur in the earliest of mythologies, are regarded widely as violating a tabu. The tabu of Doppelgänger-similarity is described in Girard (1972).

Besides dark versions of homunculi, such as the cabalistically inspired Golem and the electrically initialized creature of the surgeon Dr. Victor Frankenstein, the literature provides also more lighthearted variants. Examples are the piano-playing doll automata of the 18th century, based on the anatomical and physical knowledge of their time, and the mechanical beauty singing and dancing in *The Tales of Hoffmann*. More recent is the robot C3PO in George Lucas' film *Star Wars*, which represents a positive view of human-like robots.

Today human air travel is commonplace. Furthermore, we now know that a sparrow remains air-borne in accordance with the same aero-dynamic principles as a jumbo jet. Thus, there is a certain level of abstraction at which the flights of sparrows and jumbo jets function in the same way.

Similarly, the modeling of natural communication requires an abstract theory which applies to human and artificial cognitive machines alike. Thereby, one naturally runs the risk of setting the level of abstraction either too low or too high. As in the case of flying, the crucial problem is finding the correct level of abstraction.

A level of abstraction which is too low is exemplified by closed signal systems such as vending machines. Such machines are inappropriate as a theoretical model because they fail to capture the diversity of natural language use, i.e., the characteristic property that one and the same expression may be used meaningfully in different contexts.

A level of abstraction which is too high, on the other hand, is exemplified by naive anthropomorphic expectations. For example, a notion of ‘proper understanding’ which requires that the computational system be subtly amused when scanning *Finnegans Wake* is as far off the mark as a notion of ‘proper flying’ which requires mating and breeding behavior from a jumbo jet.⁶

VII. ANALYZING HUMAN COGNITION

The history of mechanical flight shows how a natural process (bird flight) poses a conceptually simple and obvious problem to science. Despite great efforts it was unsolvable for a long time. In the end, the solution turned out to be a highly abstract mathematical theory. In addition to being a successful foundation of mechanical flight, this theory is able to explain the functioning of natural flight as well.

This is why the abstract theory of aero-dynamics has led to a new appreciation of nature. Once the development of biplanes, turboprops, and jets resulted in a better theoretical and practical understanding of the principles of flight, interest was refocused again on the natural flight of animals in order to grasp their wonderful efficiency and power. This in turn led to major improvements in artificial flight, resulting in less noisy and more fuel-efficient airplanes.

Applied to computational linguistics, this analogy shows that our abstract and technological approach does not imply a lack of interest in the human language capacity. On the contrary, investigating the particular properties of natural language communication by humans is meaningful only *after* the mechanism of natural language communication has been understood in principle, modeled computationally, and proven successful in concrete applications on massive amounts of data.

⁶ Though this may seem reasonable from the viewpoint of sparrows.

VIII. INTERNAL AND EXTERNAL TRUTHS

In science we may distinguish between internal and external truths. Internal truths are conceptual models, developed and used by scientists to explain certain phenomena, and held true by relevant parts of society for limited periods of time. Examples are the Ptolemaic (geocentric) view of planetary motion or Bohr's model of the atom.

External truths are the bare facts of external reality which exist irrespective of whether or not there are cognitive agents to appreciate them. These facts may be measured more or less accurately, and explained using conceptual models.

Because conceptual models of science have been known to change radically in the course of history, internal truths must be viewed as *hypotheses*. They are justified mainly by the degree to which they are able to systematically describe and explain large amounts of real data.

Especially in the natural sciences, internal truths have improved dramatically over the last five centuries. This is shown by an increasingly close fit between theoretical predictions and data, as well as a theoretical consolidation exhibited in the form of greater mathematical precision and greater functional coherence of the conceptual (sub)models.

In contrast, contemporary linguistics is characterized by a lack of theoretical consolidation, as shown by the many disparate theories of language⁷ and the overwhelming variety of competing theories of grammar.⁸ As in the natural sciences, however, there is external truth also in linguistics. It may be approximated by completeness of empirical data coverage and functional modeling.

IX. LINGUISTIC VERIFICATION

The relation between internal and external truth is established by means of a *verification method*. The verification method of the natural sciences is the repeatability of experiments. This means that, given the same initial conditions, the same measurements must result again and again.

On the one hand, this method is not without problems because experimental data may be interpreted in different ways and may thus support different, even conflicting, hypotheses. On the other hand, the requirements of this method are so minimal that by now no self-respecting theory of natural science can afford to reject it. Therefore

⁷ Examples are nativism, behaviorism, structuralism, speech act theory, model theory, as well as Givón's (1985) iconicity, Lieb's (1992) neostructuralism, and Halliday's (1985) systemic approach.

⁸ Known by acronyms such as TG (with its different manifestations ST, EST, REST, and GB), LFG, GPSG, HPSG, CG, CCG, CUG, FUG, UCG, etc. These theories of grammar concentrate mostly on an initial foundation of internal truths such as 'psychological reality,' 'innate knowledge,' 'explanatory adequacy,' 'universals,' 'principles,' etc., based on suitably selected examples. Cf. Sect. 9.5.

the repeatability of experiments has managed to channel the competing forces in the natural sciences in a constructive manner.

Another aspect of achieving scientific truth has developed in the tradition of mathematical logic. This is the principle of formal consistency, as realized in the method of axiomatization and the rule-based derivation of theorems.

Taken by itself the quasi-mechanical reconstruction of mathematical intuition in the form of axiom systems is separate from the facts of scientific measurements. As the logical foundation of natural science theories, however, the method of axiomatization has proven to be a helpful complement to the principle of repeatable experiments.

In linguistics, corresponding methods of verification have sorely been missing. To make up for this shortcoming there have been repeated attempts to remodel linguistics into either a natural science or a branch of mathematical logic. Such attempts are bound to fail, however, for the following reasons:

- The principle of repeatable experiments can only be applied under precisely defined conditions suitable for measuring. The method of experiments is not suitable for linguistics because the objects of description are *conventions* which have developed over the course of centuries and exist as the intuitions (“Sprachgefühl”) of the native speaker-hearer.
- The method of axiomatization can only be applied to theories which have consolidated on a high level of abstraction, such as Newtonian mechanics, thermodynamics, or the theory of relativity. In today’s linguistics, there is neither the required consolidation of theory nor completeness of data coverage. Therefore, any attempt at axiomatization in current linguistics is bound to be empirically vacuous.

Happily, there is no necessity to borrow from the neighboring sciences in order to arrive at a methodological foundation of linguistics. Instead, theories of language and grammar are to be implemented as electronic models which are tested automatically on arbitrarily large amounts of real data as well as in real applications of spontaneous human-computer communication. This method of verifying or falsifying linguistic theories objectively is specific to computational linguistics and may be viewed as the counterpart of the repeatability of experiments in the natural sciences.

X. EMPIRICAL DATA AND THEIR THEORETICAL FRAMEWORK

The methodology of computational linguistics presupposes a theory of language which defines the goals of empirical analysis and provides the framework into which components are to be embedded without conflict or redundancy. The development of such a framework can be extraordinarily difficult, as witnessed again and again in the history of science.

For example, in the beginning of astronomy scientists wrestled for centuries with the problem of providing a functional framework to explain the measurements that had been made of planetary motion and to make correct predictions based on such a framework. It was comparatively recently that Copernicus (1473–1543), Kepler

(1571–1630) and Newton (1642–1727) first succeeded with a description which was both empirically precise and functionally simple. This, however, required the overthrow of clusters of belief held to be true for millennia.

The revolution affected the *structural hypothesis* (transition from geo- to heliocentrism), the *functional explanation* (transition from crystal spheres to gravitation in space), and the *mathematical model* (transition from a complicated system of epicycles to the form of ellipses). Furthermore, the new system of astronomy was constructed at a level of abstraction where the dropping of an apple and the trajectory of the moon are explained as instantiations of one and the same set of general principles.

In linguistics, a corresponding scientific revolution has long been overdue. Even though the empirical data and the goals of their theoretical description are no less clear in linguistics than in astronomy, linguistics has not achieved a comparable consolidation in the form of a comprehensive, verifiable, functional theory of language.⁹

XI. PRINCIPLES OF THE SLIM THEORY OF LANGUAGE

The analysis of natural communication should be structured in terms of methodological, empirical, ontological, and functional principles of the most general kind. The SLIM theory of language presented in this book is based on *Surface compositional*, *Linear*, *Internal Matching*. These principles are defined as follows.

1. *Surface compositional* (methodological principle)

Syntactic-semantic composition assembles only concrete word forms, excluding all operations known to increase computational complexity to exponential or undecidable, such as using zero-elements, identity mappings, or transformations.

2. *Linear* (empirical principle)

Interpretation and production of utterances are based on a strictly time-linear derivation order.

3. *Internal* (ontological principle)

Interpretation and production of utterances are analyzed as cognitive procedures located inside the speaker-hearer.

4. *Matching* (functional principle)

Referring with language to past, current, or future objects and events, relations, and properties is modeled in terms of pattern matching between language meaning and a context, defined as content in a speaker-hearer internal database.

These principles originate in widely different areas (methodology, ontology, etc.), but within the SLIM theory of language they interact very closely. For example, the functional principle of (4) matching can only be implemented on a computer if the overall

⁹ From a history of science point of view, the fragmentation of today's linguistics resembles the state of astronomy before Copernicus, Kepler, and Newton.

system is handled ontologically as (3) an internal procedure of the cognitive agent. Furthermore, the methodological principle of (1) surface compositionality and the empirical principle of (2) time-linearity can be realized within a functional mechanism of communication only if the overall theory is based on internal matching (3,4).

In addition to what its letters stand for, the acronym SLIM is motivated as a word with a meaning like *slender*. This is so because detailed mathematical and computational investigations have proven SLIM to be efficient in the areas of syntax, semantics, and pragmatics – both relatively in comparison to existing alternatives, and absolutely in accordance with the formal principles of mathematical complexity theory.

XII. CHALLENGES AND SOLUTIONS

The SLIM theory of language is defined on a level of abstraction at which the mechanism of natural language communication in humans and in suitably constructed cognitive machines is explained in terms of the same principles of surface compositional, linear, internal matching.¹⁰ This is an important precondition for unrestricted human-computer communication in natural language. Its realization requires general and efficient solutions in the following areas.

First, SLIM must model the hearer's *understanding* of natural language. This process is realized as the automatic reading-in of propositions into a database and – most importantly – determining the correct place for their storage and retrieval. The semantic primitives are defined as the artificial or natural agent's basic recognition and action procedures (and not as truth conditions).

Second, SLIM must model the speaker's *production* of natural language. This includes determining the content to be expressed in language, traditionally called *conceptualization*. In its simplest form, it is realized as an autonomous navigation through the propositions of the agent-internal database. Thereby speech production is handled as a direct reflection (internal matching) of the navigation path in line with the motto: *speech is verbalized thought*.

Third, SLIM must derive blueprints for action by means of *inferences*. The autonomous control of the agent's actions is driven by the principle of maintaining the agent in a state of balance vis-à-vis a continuously changing external and internal environment. Inferencing also plays an important role in the pragmatics of natural language, both in the hear and the speak mode.

¹⁰ Moreover, the *structural hypothesis* of the SLIM theory of language is a regular, strictly time-linear derivation order – in contrast to grammar systems based on constituent structure. The *functional explanation* of SLIM is designed to model the mechanism of natural communication as a speaking robot – and not some tacit language knowledge innate in the speaker-hearer which excludes language use (performance). The *mathematical model* of SLIM is the continuation-based algorithm of LA-grammar, and not the substitution-based algorithms of the last 50 years.

CONCLUDING REMARK

In summary, the vision of unrestricted natural language communication between humans and machines is like the vision of motorized flight a 110 years ago: largely solved theoretically, but not yet realized in practical systems. At this point, all it will take to really succeed in computational linguistics is a well-directed, concentrated, sustained effort in cooperation with robotics, artificial intelligence, and psychology.