

**COMPOSITIONALITY
AND
SYNTACTIC GENERALIZATIONS**

Jan Odijk

COMPOSITIONALITY
AND
SYNTACTIC GENERALIZATIONS

PROEFSCHRIFT

TER VERKRIJGING VAN DE GRAAD VAN DOCTOR
AAN DE KATHOLIEKE UNIVERSITEIT BRABANT,
OP GEZAG VAN DE
RECTOR MAGNIFICUS, PROF. DR. L.F.W. DE KLERK,
IN HET OPENBAAR TE VERDEDIGEN
TEN OVERSTAAN VAN EEN
DOOR HET COLLEGE VAN DEKANEN AANGEWEEZEN COMMISSIE
IN DE AULA VAN DE UNIVERSITEIT OP
VRIJDAG 12 NOVEMBER 1993
TE 14.15 UUR

DOOR

Johannes Engelbertus Josephus Maria Odijk

geboren te Schiedam

Promotores:

Prof. Dr. H.C. van Riemsdijk

Prof. Ir. S.P.J. Landsbergen

The work described in this thesis has been carried out at the Philips Research Laboratories in Eindhoven, the Netherlands, as part of the Philips Research program.

Acknowledgements

This thesis is one of the results of research conducted in the Rosetta machine translation project which was carried out at Philips Research Laboratories. I am grateful to Jan Landsbergen and Henk van Riemsdijk for their supervision. Furthermore, I would like to thank Jan for giving me the opportunity of working at Philips and for supplying the opportunities to write theses in general and this one in particular, and Henk for waiting so long for my original thesis, finally receiving something completely different from what he must have expected.

I would also like to thank all members from the Rosetta team, and the students who visited us. Some of them require special mentioning. Agnes Mijnhout initiated the research on predicate-argument relations, on which we report in chapter 4 of this book, and Angeliek van Hout made an initial study of the problems which R-pronouns pose in the Rosetta framework. My room mates André Schenk and later Joep Rous had a very positive influence on me and my work through the stimulating discussions we had on our work and other topics.

And I would like to thank Lisette Appelo. The often intense discussions with her and her comments on earlier versions of this book improved it considerably.

I am obliged to Lex Augusteijn and Wijbrand Siedenburg for their invaluable T_EXnical and PostScript support in preparing the final document.

Verder wil ik mijn ouders bedanken, met name mijn vader, die deze promotie helaas niet meer mee kan maken. Zij hebben mij altijd volop de gelegenheid gegeven en gestimuleerd te studeren.

Tot slot bedank ik Margriet. Zij bleef er altijd in geloven (en erop aandringen) dat ik zou promoveren. En het gaat er dus toch van komen.

Jan Odijk

Contents

1	Introduction	1
1.1	Formulation of the Problem	1
1.2	Syntax and Semantics	2
1.3	Theoretical and Application-Oriented Grammars	4
1.3.1	Artifacts	5
1.3.2	Descriptive Adequacy	7
1.3.3	Broad Coverage	8
1.3.4	Formal Differences	9
1.4	Organization of this Book	9
2	The Grammatical Framework	11
2.1	Compositional Grammars	11
2.2	S-trees	12
2.3	Free M-grammars	13
2.3.1	Description	13
2.3.2	Example M-grammar and Example Derivation	15
2.4	Controlled M-grammars	20
2.5	Translation Relation	22
2.6	Translation System	30
2.6.1	Description of the System	30
2.6.2	Example Translation	32
2.7	Some Specific Assumptions	35
2.8	Concluding Remarks	39
3	Compositionality and Syntactic Generalizations	41
3.1	Introduction	41
3.2	Auxiliaries and Inversion in English	42
3.2.1	The Problem	42
3.2.2	A Solution	45
3.3	Mood in Dutch	47
3.4	Order variants	49
3.5	Wh-movement	51
3.6	Generic Sentences	53

3.7	A Different Kind of Modularization	56
3.8	Concluding Remarks	60
4	Predicate-Argument Relations	63
4.1	Introduction	63
4.2	Predicate-Argument Relations and Compositionality	63
4.3	Arguments	65
4.3.1	Argument-Ordering Convention	68
4.3.2	External and Internal Arguments	70
4.3.3	Attribute-Value Pairs to Specify Arguments	72
4.4	Covert Arguments	75
4.5	Bound Adjuncts and Adverbials	79
4.6	Small Clauses	82
4.7	Systematic Relations between Predicates	92
4.8	Concluding Remarks	95
5	Some Constructions	97
5.1	Introduction	97
5.2	Passivization	98
5.3	Verb Second in Dutch	104
5.4	Unbounded Dependencies	109
5.5	Crossing Dependencies in Dutch	113
5.5.1	Outline of the Analysis	114
5.5.2	Verb-Raising Transformations	116
5.5.3	Pruning	118
5.5.4	Surface Grammar	120
5.5.5	Concluding Remarks	122
6	R-pronouns	125
6.1	Introduction	125
6.2	Some Relevant Facts	126
6.3	The Functions of R-pronouns	127
6.3.1	Expletive Function	127
6.3.2	Prepositional R-pronouns	136
6.3.3	Locative R-pronouns	138
6.3.4	Quantificational R-pronouns	138
6.4	The Distribution of R-pronouns	139
6.4.1	General Discussion	139
6.4.2	Global Characterization of the Results	142
6.5	The Assumptions in More Detail	148
6.6	Illustration	156
6.7	Concluding Remarks	158

CONTENTS

v

7 Concluding Remarks	163
7.1 Conclusions	163
7.2 Topics for Further Research	165
Samenvatting (summary in Dutch)	179
Curriculum Vitae	183

Chapter 1

Introduction

1.1 Formulation of the Problem

The central problem of this dissertation is the question how syntactic generalizations can be adequately captured in a compositional framework.

This problem will be investigated within the *controlled M-Grammar* formalism. I will describe how a number of complex syntactic constructions have been dealt with in this formalism, which has been used in developing the Rosetta machine translation system. In particular, I will show that these syntactic constructions have been dealt with in a syntactically adequate manner in a framework which is *compositional* in nature, and where consequently the grammar has a strong semantic bias.

The syntactic generalizations that I am mainly interested in here relate to the fact that many constructions can be described most adequately by a conglomerate of construction-independent rules. This construction-independence of syntactic rules, and the relation between syntax and semantics will be discussed in more detail in section 1.2.

In addition, the research has been carried out in the context of research into machine translation, i.e. application-oriented research. Application-oriented research differs from purely theoretical research in a number of respects. Some of the differences will be presented in more detail in section 1.3.

The general conclusions of this study are (1) that the grammatical formalism used, *controlled M-grammar*, supplies — due to its compositional nature — a firm framework to deal with certain phenomena, especially when they relate fairly directly to semantics (e.g. predicate-argument relations); (2) that the framework makes it possible to incorporate analyses in which constructions are created by a conglomerate of construction-independent rules; (3) that it is possible to incorporate and extend syntactically adequate descriptions based on insights from theoretical linguistics into this compositional framework in a fairly direct manner.

It will, however, also become clear that many improvements of the framework, or of specific linguistic analyses within it, are still possible.

In the Rosetta translation system, grammars have been developed for the languages Dutch, English and Spanish. The general conclusions presented will therefore be established on the basis of examples from these languages which have actually been implemented, with an emphasis on Dutch.

1.2 Syntax and Semantics

The aim of this book is to show that a number of complex syntactic phenomena can be (and have been) dealt with in a compositional grammar by means of rules which are not construction-specific. The most adequate syntactic description of a construction often turns out to be a description in which the construction arises as a consequence of the interaction of a number of different syntactic rules, most of which play a role in the formation of a wide variety of other constructions as well and are therefore not specific to one construction. I have attempted to achieve such descriptions in the Rosetta machine translation system as often as possible.

The attempt to make rules of grammars construction-independent is one aspect of the so-called *Move α* program as pursued by Chomsky and other theoretical linguists. Another aspect of this program is the attempt to replace language-specific rules by language-independent principles. In fact, within this program, the thesis that there are no construction-specific and no language-specific rules, for a substantial core of the grammar, is pushed to its limits. As will become clear later, many analyses of constructions within the Rosetta system have been inspired by analyses of these constructions proposed within the *Move α* program, though the emphasis was on developing analyses which use construction-independent but language-specific rules. The issue of language-specificity will be briefly mentioned, and the conclusion is that improvements are possible. A way of dealing with this aspect is suggested, but the analyses to be presented are all formulated in terms of language-specific rules. A comparison with other possible analyses of the same facts from computationally oriented frameworks is made incidentally.

It is certainly not the case that I attempt to describe a direct implementation of a grammar constructed in accordance with the *Move α* program. On the contrary, an attempt was usually made to extract a number of key features of a specific analysis and to express these in the analysis developed. There are several reasons for not directly implementing such a system. First, the aim is to develop adequate syntactic analyses in a compositional framework. The compositionality of the framework imposes restrictions which do not hold in a pure *Move α* framework. Second, analyses developed within this framework tend to be extremely non-deterministic, which is undesirable in a computational environment. Third, the restriction ‘for a substantial core’, mentioned above, is crucial. Given the nature of the Rosetta project (see below), it was impossible to restrict oneself to a ‘substantial core’. In addition, the *Move α* program is a program which is still under development, and which still has many unresolved problems. In many analyses presented within this framework, language-specific rules also occur, for phenomena which have not yet been investigated, or which are not yet fully understood, or

for which a further reduction is perhaps not possible at all.

The attempt to describe constructions as the result of a number of interacting construction-independent rules within the controlled M-grammar framework has the flavor of a paradox. The controlled M-grammar framework is designed to describe translation-equivalent utterances of different languages ‘isomorphically’. The exact meaning of this notion will become clear below, when the adopted framework will be described in more detail (chapter 2), but it can be explained informally in the following manner. The basic idea is that a construction from one language corresponds to a construction from another language with the same meaning. Under the ‘isomorphic’ approach, the grammar of one language contains a rule which creates this construction and the grammar of the other language contains a different rule which creates the translation-equivalent construction. This appears to force one to describe these constructions by means of construction-specific rules, and to a certain extent this is actually the case. So, how can one devise analyses with construction-independent rules in this framework? The paradox is resolved in the following way: it is carefully determined which aspects of a construction are essential to the construction and its meaning or translation. For these aspects, a rule (or more than one rule) is written which is construction-specific, and which partakes in the isomorphy relation. All other aspects, however, are dealt with by rules which do not partake in the isomorphy relation (called *syntactic transformations*), and which need not be construction-specific. Chapter 3 will describe this in more detail for a number of constructions.

Syntactic transformations play a crucial role here in making it possible to adequately express syntactic generalizations. Transformations are important because they are rules which have the identity operation as their meaning operation, which makes it possible to capture syntactic generalizations, and they allow, like all rules of the grammar formalism adopted here, powerful structural operations on syntactic trees.

In many other computationally oriented syntactic formalisms, transformational operations are not available, so in these frameworks many phenomena for which well-motivated analyses in terms of transformational operations exist must be dealt with in a different manner. I will show that in most cases these alternative analyses are inferior to the analyses from the transformational frameworks, either because of observational inadequacy (not all relevant facts are covered) or because of descriptive inadequacy (existing regularities and generalizations are not adequately captured), or both.

Rules with transformational power which apply to syntactic trees have never been very popular in computational linguistics. Nevertheless, I claim that they are useful, and in fact necessary to describe language adequately. Most of the objections to the use of transformations in natural language processing systems (e.g. the problems for transformational parsing pointed out by [King, 1983]) are not applicable to the specific variant of transformations and their mode of application within the Rosetta framework. I will argue, based on the analysis of several syntactic phenomena implemented in the Rosetta machine translation system, that transformations provide an excellent operation to adequately describe (and pro-

cess) natural language, and that the analysis given in the Rosetta grammars is superior to non-transformational analyses.

1.3 Theoretical and Application-Oriented Grammars

The work for this dissertation was carried out in the Rosetta project. This project was a unique combination of basic research into machine translation and large-scale development of grammars. Within the project, a number of research themes were defined, but at the same time it has always been the intention to develop a large prototype system which could form the basis for real applications of the machine translation system and of its spin-offs (see [Odijk, 1992]). All analyses described in this dissertation have actually been implemented in the Rosetta3¹ machine translation system, unless explicitly indicated otherwise. The system was developed and tested in the period from 1985 to 1991. One of the consequences of this character of the project is that from a given point in time certain assumptions had to be fixed, and the analysis had to be performed within the limits existing at that point. This had advantages, but also disadvantages. A clear disadvantage is that it was impossible or very difficult to change certain assumptions, even if they turned out to complicate the system. Advantages, however, were that the relevant assumptions were relatively fixed, and that certain analyses were carried out to their limits within these restrictions, which made very clear what problems the restrictions imposed. Certain improvements of the system would probably not have been thought of as being relevant if the current system had not been applied this systematically to such a large array of constructions. This dissertation will contain a number of suggestions and proposals for possible improvements of the current grammars which have not yet been implemented.

The research carried out in the Rosetta project was application-oriented research. This kind of research differs from purely theoretical research in a number of ways. In particular, I would like to indicate here how the grammars (or grammatical descriptions) constructed in application-oriented research differ from grammars or grammatical descriptions constructed in purely theoretically oriented research. The major differences relate to the following points:

- Application-oriented grammars are artifacts, not intended as theories of any real existing mental or Platonic object.
- In application-oriented grammars one attempts to achieve descriptive adequacy, but observational adequacy must never be sacrificed. In theoretically oriented grammars *explanatory adequacy* is aimed for.
- Broad coverage is desirable for application-oriented grammars, but of less importance for theoretically oriented grammars.

¹The Rosetta3 system was preceded by two smaller systems, Rosetta1 and Rosetta2.

- Application-oriented grammars must be expressed in a formalism which makes it possible to actually parse and generate utterances. In grammars set up for purely theoretical purposes these considerations play no role.

Each of these points will be discussed in more detail.

1.3.1 Artifacts

The Rosetta grammars are artifacts. They are not intended as a theory about the internally represented grammar of human beings, nor as descriptions of a grammar as an abstract Platonic object. They solely serve as a grammar for the machine translation system, and their design is specifically adapted to machine translation. This does not mean that they cannot be used for other purposes (in fact, they can be and have been used for other purposes). The fact that the Rosetta grammars are artifacts, constructed specifically to serve a special purpose (machine translation), implies that completely different criteria of adequacy are relevant for the Rosetta grammars than for other grammars which are intended for other purposes. The adequacy of the Rosetta grammars is to be measured by their success as a part of the machine translation system. The activity of constructing grammars oriented to a specific purpose is a form of *linguistic engineering*, and this dissertation is about linguistic engineering.

It must be emphasized that this does not imply that the Rosetta grammars can be more sloppy (i.e. overaccept or overgenerate) than theoretical grammars. It has always been the intention to construct very precise grammars which describe the language correctly. That is not only more interesting and more challenging from a theoretical point of view, but also important for practical reasons. Designing precise grammars will make it possible to use these grammars in applications other than machine translation, and ‘spurious analyses’ (which can lead to ‘spurious ambiguities’ for well-formed input) are avoided. ‘Spurious analyses’ are empirically incorrect analyses assigned to an expression by a grammar which is not sufficiently precise (see below for examples). Avoiding spurious ambiguities forms a major practical reason why the grammars must be very good. The relevancy of precise grammars to avoid spurious ambiguities has been emphasized by [Flickinger et al., 1987] and [Sag, 1991], and I fully agree with the remarks made there. The same point has also been made by [Isabelle, 1989] and Van Noord([van Noord, 1993]:11) who also argue that constructing *reversible* grammars makes the task of constructing precise grammars easier. One might think that a grammar can be more sloppy for sentences which will never be offered for translation. However, it is very difficult to establish that certain sentences will never be offered to a system, unless the system is intended for a very restricted application. This can be illustrated by one of our own experiences, which again shows that it is absolutely vital to have very precise grammars, and that ignoring certain rules will almost always cause problems.

In Rosetta, it was decided to first construct sentence grammars, and to build systems that could handle text and relations between sentences only in a later

phase. Because of this, it was decided that the interpretation of non-reflexive pronouns would not be dealt with in the first phase, because it was assumed that rules dealing with those phenomena would be rules of text grammar.

It was also generally assumed that there was at least one rule that dealt with a part of the interpretation of non-reflexive pronouns in the sentence grammar, viz. an analog of principle B of the Binding Theory of [Chomsky, 1981]. This principle states that a pronoun cannot be bound by an antecedent within a certain local, structurally determined domain.² It was expected, however, that this principle required no implementation in the first phase, for the following reasons. First, for third-person pronouns, principle B does not exclude sentences, only certain interpretations of sentences, in particular relating to the interpretation of pronouns. Leaving out principle B for these pronouns would have no negative effects if pronominal interpretations were not dealt with. Second, principle B does exclude sentences in the case of first and second-person sentences, but it was assumed that such ungrammatical input would never be offered to the system. So, for this case principle B could be omitted as well, and since no other cases were to be considered, it was decided that principle B would not be implemented in the first phase.

It turned out, however, that this soon led to problems. The reasons for this were the following. First, in Dutch, reduced first and second-person *reflexive* pronouns have the same form as reduced first and second-person *non-reflexive* pronouns. As a consequence, in a sentence such as *ik scheer me* the reduced pronoun *me* could be analyzed either as a reflexive or as a non-reflexive pronoun. With *me* as a reflexive pronoun, the sentence is well-formed, and it is correctly translated into ‘*I am shaving*’. With *me* as a non-reflexive pronoun, the sentence is ill-formed, and it would be excluded by principle B if it were implemented. Since this was not the case, the system yielded an additional result (**I am shaving me*), which is neither a correct English sentence, nor an acceptable translation of the Dutch sentence.

Even more unexpectedly, the absence of principle B caused problems in a completely different type of sentence. Dutch has imperative constructions with overt subjects, e.g. the sentence *komt u binnen*, literally *come you inside*, is a sentence containing an imperative verb (*komt*) and an overt subject (*u*). The sentence is correctly translated by the system into *please enter* or *please come inside*. Again, however, the system yielded an additional result which was neither a correct English sentence nor acceptable as translation. Since the Dutch word *binnen* is also a post-position, the sentence given could also be analyzed as an imperative sentence without an overt subject, and with *u* as a complement to the post-position *binnen* (as in *komt het huis binnen* ‘*come into the house*’), resulting in the translation **come into you*.

Both problems mentioned were caused because principle B was not implemented. When these undesirable results were encountered, it was decided to implement a version of principle B.

This example clearly shows that abstracting from certain rules of grammar will

²This domain can be defined very precisely, but that is irrelevant in this context.

almost always lead to spurious parses. In Rosetta, however, there are additional reasons that the highest quality grammars are required.

It has always been the central idea behind the Rosetta system that it would be an interactive system in which the system, after analyzing the source text, and after consulting the user in the source language to resolve ambiguities that the system cannot resolve on its own, would generate a well-formed sentence in the target language. The kind of applications envisaged requires no knowledge by the user of the target language. Given this kind of application, the grammars must be very reliable in generation: the system must be able to generate a correct sentence of the target language without any additional help from the user or any other person.

1.3.2 Descriptive Adequacy

A second difference between theoretically oriented grammars and application-oriented grammars concerns the relevant level of adequacy on the scale set by [Chomsky, 1964, 28-9].

In theoretical linguistics, one attempts to achieve *explanatory adequacy*, i.e. one tries to describe phenomena in a way that relates these phenomena to rules and principles that are part of a general theory of grammars and of language. This goal is so important that very often observational adequacy is sacrificed (temporarily) to achieve it.

Explanatory adequacy is not a main concern if one intends to develop a machine translation system. Of course, wherever possible one makes use of descriptions which relate to known general properties of grammar, so that e.g. the description of related phenomena from different languages is more uniform, but it is not a purpose in itself to make descriptions that satisfy the criterion of explanatory adequacy.

Because the success of the grammar in the machine translation system is the only criterion that counts, *observational adequacy*, i.e. correctly describing the facts of the language and translating correctly, is essential.³

In addition, extensibility and modifiability of the grammar is essential, so that one attempts to reach the level of *descriptive adequacy*, i.e. a description that correctly captures all relevant generalizations and regularities within the language. One has to attempt to achieve this goal, because if one does not do so, the size and complexity of the grammar will become intolerable. However, this level of adequacy often cannot be achieved, and in many cases a description is required which attains at best the level of observational adequacy.

These differences between theoretically oriented grammars and the Rosetta grammars can be illustrated with the following example. Consider again the distribution of reflexive pronouns. It is well-known that a reflexive pronoun entertains a relation with some other phrase in the sentence (its antecedent), and that there are strong conditions on where an antecedent for a reflexive pronoun can occur

³I do not limit the notion of observational adequacy to the correct specification of a *finite* set of observed data, though Chomsky might have intended this.

relative to the pronoun. Many theoretical works have dealt with the issue of how to characterize precisely this relation between the antecedent and a reflexive pronoun in a unifying way. In the Rosetta grammars, it is not necessary to find a unified way of characterizing this relationship: if necessary, simply all possible configurations, or all classes of configurations (if there is an infinite number of possible relationships) can be spelled out, and this has indeed been done in this manner for the case under discussion. Of course, the number of grammar rules will then increase, but for the practical purposes of machine translation coverage of the relevant facts counts more than searching for a unified account of this relationship. This is not to say that looking for a unified characterization of this relationship should not be done. On the contrary, it should, but it should be done by theoretical linguists and not by developers of a machine translation system. If a convincing unifying characterization is found, these results can be incorporated into the Rosetta grammars (if the formalism allows it, and if not, the formalism might be adapted to enable the incorporation of these results).

1.3.3 Broad Coverage

As pointed out before, *broad coverage* is essential. There are two aspects to this notion.

First, certain phenomena occur in virtually every utterance (e.g. time, predicate-argument structure, etc.). Since a real working system is to be developed, one simply cannot abstract from such phenomena. If these phenomena are not taken into account, not a single utterance can be dealt with. Thus, several different linguistic phenomena must somehow be dealt with in the grammars of a machine translation system.

It is also not possible to abstract from certain facts and make a distinction between *core* and *peripheral* facts in the same way as in theoretical linguistics. In theoretical linguistics, such distinctions are determined by theoretical considerations: recalcitrant facts which appear to form direct evidence against perhaps far-reaching hypotheses can be set aside and ignored temporarily by declaring them *peripheral*. One cannot do that when constructing a machine translation system. Of course, a distinction between core and periphery can be made (and it actually has been made), but it must be made on the basis of completely different criteria.

Related to this is the fact that in a machine translation system real grammars must be written: theoretical linguists do not write grammars. They investigate certain constructions in some detail, abstracting away from other phenomena, and they concentrate on the relevancy that the construction investigated has for the general theory of grammar. The fact that phenomena can be studied in isolation, abstracting away from other phenomena, can also lead to a situation in which analyses for different phenomena are mutually inconsistent, even if these analyses assume the same framework to work in. Of course, there is nothing wrong with that (though the inconsistency should of course be resolved at a certain point), but it does imply a difference between theoretical linguists and linguistic engi-

neers: linguistic engineers do write grammars, and they must incorporate in these grammars mutually consistent descriptions of several different phenomena.

In a second sense of the term *broad coverage*, it means that many different constructions must be dealt with. Actually occurring texts are enormously varied; very many different constructions occur even in the simplest and most uniform texts.

This has an important further consequence: for research into machine translation, the time and the resources to investigate every phenomenon in detail are simply not available: perhaps some phenomena can be studied in detail, but most of this work must be left to theoretical linguistics.

Linguistic engineers must therefore base their work on existing analyses. Their major task is to try to extract the common properties from all the proposals made in the theoretical literature and in traditional grammars etc., make concrete decisions about the non-common properties, translate them into the terms of their own framework, and integrate this consistently with all other elements in the grammars used in the system. This dissertation is in part a report of an attempt that has been made to do exactly this, although for some phenomena new and original analyses are proposed.

Because of this, the quality of systems which process natural language making use of explicit grammars will be a function of the status of theoretical linguistics.

1.3.4 Formal Differences

Further differences between application-oriented grammars and theoretical grammars relate to more formal properties.

First, in a machine translation system both generation and analysis procedures are required. It would be undesirable if generation and analysis were based on different grammars. In the ideal case one would like to have one grammar from which both a generation and analysis procedure is derived automatically. This requires the grammar to be reversible.

Second, both generation and analysis must be effective, i.e. they must generate or analyze a sentence in finite (and preferably very short) time. These two properties (reversibility and effectivity) impose additional requirements on the grammar. These additional requirements are irrelevant for theoretical grammars: theoretical linguists who do not want to implement their system will not reject a theory of grammar if it happens not to be reversible or effective, but these are very important properties of the Rosetta grammars.

1.4 Organization of this Book

This dissertation is organized as follows. In chapter 2 the framework assumed in this book, the *controlled M-grammar framework*, is described in detail. A general discussion on the necessity of syntactic transformations in a compositional framework, their role in such a framework, and their consequences for the relation

between syntax and semantics is given in chapter 3. In chapter 4 the treatment of predicate-argument relations within the adopted framework is outlined. It is shown there that the compositional nature of the framework determines this treatment to a large degree. In chapter 5 the analysis of several specific constructions is dealt with in detail: the analysis of passive constructions (section 5.2), the analysis of Verb-Second phenomena in Dutch (section 5.3), the analysis of unbounded dependencies (section 5.4), and crossing dependencies in Dutch (section 5.5). A special chapter is devoted to the complex syntax of R-pronouns in Dutch (chapter 6). In these chapters, existing analyses of the constructions mentioned are adapted to the specific framework adopted and the requirements imposed, and partially new and original analyses are proposed for complex syntactic phenomena.

Chapter 7 states the conclusions and summarizes which additional improvements are needed and possible in the system.

Chapter 2

The Grammatical Framework

In this chapter, the framework adopted in this dissertation, the *controlled M-grammar* framework, will be described in detail. It will be preceded by an introduction of the *free M-grammar* framework, since controlled M-grammars are an extension of free M-grammars. More elaborate descriptions of these formalisms can be found in [Landsbergen, 1987] for free M-grammars, [Appelo et al., 1987] for controlled M-grammars and in [Rosetta, 1993].

First, the concept of *compositional grammar* is defined (section 2.1). Free M-grammars and controlled M-grammars are briefly characterized and compared. Section 2.2 introduces the concept of *S-tree*, which plays an important role in these frameworks. Next, *free M-grammars* are discussed in more detail (section 2.3). The framework actually adopted in this dissertation, *controlled M-grammar*, is sketched in section 2.4. In section 2.5 it will be outlined how this grammatical framework can be used to define a translation relation, and section 2.6 describes how it can form the basis of an actual machine translation system (section 2.6). Finally, a number of assumptions which have been made within the formalism proposed are enumerated. These are substantive elements which are not part of the formalism, but which have been assumed in the grammars of the Rosetta3 machine translation system. They are crucial for a correct understanding of the rest of this book.

2.1 Compositional Grammars

We speak of *compositionality of meaning* if the meaning of an expression is a function of the meanings of the parts of the expression and the way they are combined.

Compositional Grammars are grammars which incorporate compositionality of meaning in a very direct manner: form and meaning are explicitly related by

the so-called rule-to-rule approach. In this approach it is assumed that form and meaning are composed in tandem in the following manner: a compositional grammar consists of *basic expressions* and *syntactic rules*. A syntactic rule can turn a tuple of expressions (its arguments) into a new expression. Basic expressions are associated with *basic meanings*, and syntactic rules are associated with *meaning operations*. The form of an expression can be constructed by recursively applying the syntactic rules to their arguments, initially basic expressions. The meaning of an expression can be derived in tandem by recursively applying the meaning operations associated with the rules to their arguments, initially the basic meanings associated with the basic expressions.

M-grammars are a special kind of compositional grammars, which can be viewed as a computational variant of Montague grammar (see e.g. [Montague, 1974b]). Their properties will be discussed in more detail in section 2.3. They will be briefly characterized here by comparing them to *Montague Grammar*.

The major differences between M-grammar and Montague grammar are: (1) an M-grammar uses a special kind of syntactic objects (called *S-trees*) and a special kind of rules (called *M-rules*) operating on S-trees in its syntax, while Montague used rules operating on symbol strings ([Montague, 1974a]) or a form of categorial grammar ([Montague, 1974b]). (2) M-grammars are intended to be used in natural language processing systems, in particular (though not exclusively) in machine translation systems. For this reason M-grammars satisfy a number of conditions which are motivated from a computational point of view, in particular, M-grammars must define a relation which can form the basis both for analysis and generation of utterances. We will see below how this is achieved. See [Landsbergen, 1981] for an adaptation of Montague grammar to the requirements of parsing.

Controlled M-grammars are an extension of free M-grammars. As in free M-grammars, the syntactic objects they manipulate are S-trees, and the rules are M-rules, but controlled M-grammars differ from free M-grammars in the following respects: (1) controlled M-grammars consist of subgrammars; in a subgrammar, it is possible to control the application of rules: whether they are optional, obligatory, iterative, etc. (2) in controlled M-grammars a special kind of rules, *syntactic transformations*, is distinguished. They do not contribute to the meaning, in other words they have the identity operation as associated meaning operation. Controlled M-grammars are discussed in more detail in section 2.4.

The basic objects manipulated by rules of the free M-grammar framework and of the controlled M-grammar framework are S-trees. In section 2.2 I will first define this notion and describe how they are notated.

2.2 S-trees

In order to be able to describe the grammatical frameworks, first the notion of *S-tree* will be introduced. Informally, an S-tree is a labeled ordered tree familiar from theoretical linguistics. It is used to specify the syntactic structure of an

utterance. From a more formal point of view, an S-tree is a tuple $\langle N, S \rangle$, in which N is a *node*, and S is a sequence of tuples $\langle \langle r_1, t_1 \rangle, \dots, \langle r_n, t_n \rangle \rangle$ in which each r_i is a *grammatical relation*, and each t_i is an S-tree. A *node* is a tuple $\langle C, AV \rangle$, in which C is a *syntactic category* and AV is a set of attribute-value pairs.

The following notation is used for S-trees:

$$(1) \quad N[r_1/t_1, \dots, r_n/t_n]$$

Here N is a node, and $r_1/t_1, \dots, r_n/t_n$ represents a sequence of tuples $\langle r_1, t_1 \rangle$ to $\langle r_n, t_n \rangle$. Each r_i is a grammatical relation and each t_i an S-tree. If this sequence is empty we notate $N/$ or simply N .

Nodes are represented as follows: $C\{a_1:v_1, \dots, a_n:v_n\}$, where C is a syntactic category and where each $a_i:v_i$ is an attribute-value pair.

It is often useful to focus on certain aspects of S-trees. In such cases S-trees will be represented partially, by only showing the properties relevant in the context and leaving out all other details. S-trees will also often be represented by the familiar tree diagrams, since these are generally easier to read than the notation given above.

The syntactic categories and the grammatical relations are defined for each individual language by enumeration. For each syntactic category in each language the set of appropriate attributes is also specified by enumeration. The values for attributes in attribute-value pairs can be atomic, or they can be finite sets of atomic values. The possible values for each attribute are defined for each language by enumeration or by a set-constructor.

A *surface tree* is a special kind of S-tree with the property that its 'leaves' correspond to the words of the actual sentence. These leaves are small S-trees, called *lexical S-trees*, which have special syntactic categories, called *lexical categories* as the categories of their top nodes and which correspond to words.

2.3 Free M-grammars

2.3.1 Description

Free M-grammars consist of three components: (1) the syntactic component, (2) the semantic component, and (3) the morphological component.

- (1) The syntactic component of a free M-grammar defines a set of surface trees. It consists of basic expressions and rules. The basic expressions in a free M-grammar are S-trees, called *basic S-trees*. The rules of the syntactic component are compositional rules that take S-trees as their arguments and produce an S-tree as their result; they are called *M-rules*. An S-tree which is formed by the application of M-rules to basic expressions is called a *derived S-tree*. A surface tree is a special kind of derived S-tree, of which the category of the top node is taken from a special set, called *TOPCATS*.

The derivation of an expression is represented in a *syntactic derivation tree* (I will often use *D-tree* as an abbreviation for *derivation tree*). A syntactic D-tree is a tree which consists of names of basic expressions (as its leaves) and names of rules (as its non-leaves). It indicates how a surface tree has been derived.

- (2) The semantic component of a free M-grammar relates syntactic derivation trees to *semantic derivation trees*. A semantic derivation tree is a tree which consists of names of basic meanings (as its leaves) and of names of meaning operations (as its non-leaves). The semantic component can relate syntactic D-trees to semantic D-trees simply by relating basic expressions to basic meanings and syntactic rules to meaning operations. Corresponding syntactic and semantic D-trees have the same geometry, and differ only in the labels at their nodes.

The semantic component is rather rudimentary in character. For translation purposes it is sufficient to deal with *translational equivalence*. For this reason the semantic component need not be worked out in full detail. See [Rosetta, 1993] for additional discussion of this point.

- (3) The morphological component relates lexical S-trees occurring in surface trees to symbol strings. The morphological component will not be discussed here. See [Rosetta, 1993].

Before free M-grammars are illustrated by deriving a sentence with an example M-grammar, the notions of *syntactic variable*, *substitution rule* and *start rules* will be introduced.

A special type of element in S-trees, called a *syntactic variable*, plays an important role in the syntax of the Rosetta grammars. The example grammar contains them as well. A syntactic variable is a basic S-tree consisting of one node, having a specific category¹ and associated attribute-value pairs, and a special attribute called the *index*. The *index* attribute takes integers as values. Variables function as place-holders for the phrases that will be substituted for them later in the derivation. There is a set of special M-rules, called *substitution rules*. These are parameterized rules, taking an *index parameter* that can assume integer values, and two arguments, i.e. two S-trees. A substitution rule with index parameter i applies to an S-tree containing a variable with index i , and some other S-tree (the so-called *substituent*), and substitutes this latter S-tree for the variable with index i occurring in the former S-tree. So, the index parameter determines which variable the substituent must be substituted for. In derivation trees the notation $R_{subst,i}$ is used to indicate that the substitution rule R_{subst} substitutes for the variable with index i .

A derivation usually begins with a *start rule*, i.e. a rule that combines some argument-taking basic expression (a *predicate*) with a number (zero or more) of

¹In the real grammars there are several categories for syntactic variables, but in the example grammar only one is distinguished.

variables. There are start rules which combine a monadic predicate (e.g. *pass*) with one syntactic variable, start rules which combine a dyadic predicate (e.g. *see*) with two variables, etc.

2.3.2 Example M-grammar and Example Derivation

In order to clarify how a sentence can be derived by a free M-grammar, I will (partially) specify an example M-grammar G_E and derive the example sentence *a car passed slowly* with it. The example is intended for expository purposes only, and it is for this reason unrealistically simple. No linguistic claims are made with the example at all.

First, the syntactic categories and the grammatical relations of the example grammar are specified (in table 2.1). For each syntactic category, it is specified which attributes are appropriate, and what their possible values are. The possible values of an attribute are indicated by a name for a type, which is defined separately.

The grammar contains the M-rules listed in (2). These rules will not actually be formulated here. An informal characterization of their task is given below. The names of the meaning operations associated with these rules are also given.

Rule Name	Meaning Name
Rstart1	Lstart1
Radv	Ladv
Rmod	Lmod
Rsubst	Lsubst
Rpast	Lpast
Rmaindecl	Lmaindecl
RNP	LNP

The basic expressions of the example grammar are given in (3), together with the names of their associated basic meanings. The names of the basic expressions are equal to the values of the attribute *name*.

Basic Expression	Basic Meaning
BADV{name:THERE, boundness:stem}	B1
BN{name:CAR}	B2
BV{name:PASS}	B3
BADV{name:LY, boundness:suffix}	B4
BADJ{name:SLOW}	B5

Apart from the basic expressions mentioned in (3), an infinite number of variables $VAR\{index:1\}$, $VAR\{index:2\}$, ... with associated basic meanings X_1, X_2, \dots are basic expressions of the example grammar. The names of these basic expressions are represented as x_1, x_2, \dots ²

²The basic expression with name THERE which occurs in (3) does not play a role in the derivation of the example sentence. But it can be used in other sentences, e.g. in *The car passed there*. See section 2.5 for further discussion.

RELATIONS	
det, head, mod, subj	
Type name	Possible values
nametype	A, CAR, PASS, SLOW, LY, THERE, ...
boundnesstype	stem, suffix
numbertype	singular, plural, unspecified
deftype	def, indef
voicetype	active, passive
moodtype	declarative, wh-interrogative, imperative
typetype	main, subordinate
formtype	finite, participle, infinitive, ingform
tensetype	present, past, unspecified
indextype	1,2,3,...
SYNTACTIC CATEGORIES	appropriate AV-pairs
ART	name: nametype
ADV	
BADJ	name: nametype
BADV	name: nametype, boundness: boundnesstype
BN	name: nametype
BV	name: nametype
N	number: numbertype
NP	number: numbertype, def: deftype
S	voice: voicetype, mood: moodtype, type: typetype
V	number: numbertype, form: formtype, tense: tensetype
VAR	index: indextype
LEXICAL CATEGORIES	
ART, ADV, N, V	

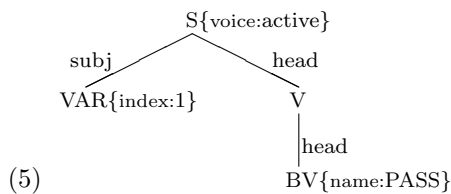
Table 2.1: Relations, categories and A-V-pairs of the example grammar G_E

In the example sentence, there is one expression which is not a basic expression, but which is introduced syncategorematically (i.e. some rule introduces this expression though it is not an argument of the rule).

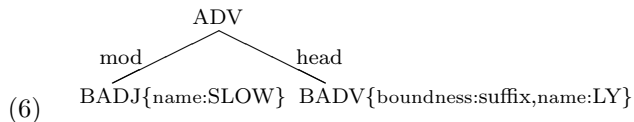
(4) ART{name:A}

The major steps in the example derivation can be described in the following manner. The syntactic component derives a surface tree for this sentence. By recording how the S-tree has been derived, a D-tree is created, which can be mapped onto one or more semantic derivation trees in the semantic component. The morphological component turns the sequence of lexical S-trees occurring in the surface tree into a symbol string.

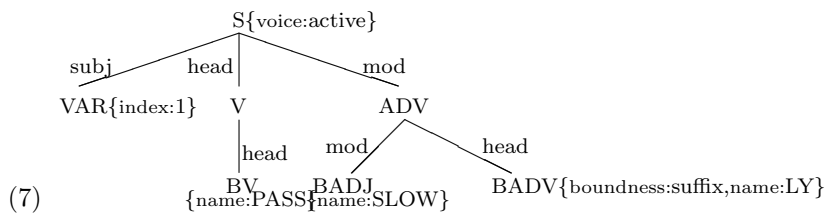
This derivation will now be illustrated in more detail. We start by applying M-rules. A start rule (*RStart1*) is applied to the basic expression BV{name:PASS} and a variable VAR{index:1} to form an active sentence. This yields the following S-tree:



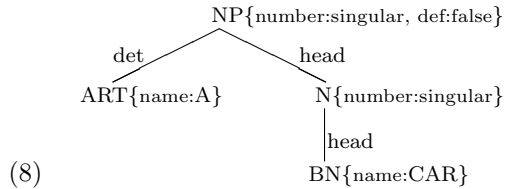
In addition, an adverb is created from the basic expression SLOW (of category BADI) and the suffix LY by the M-rule *Radv*. This yields:



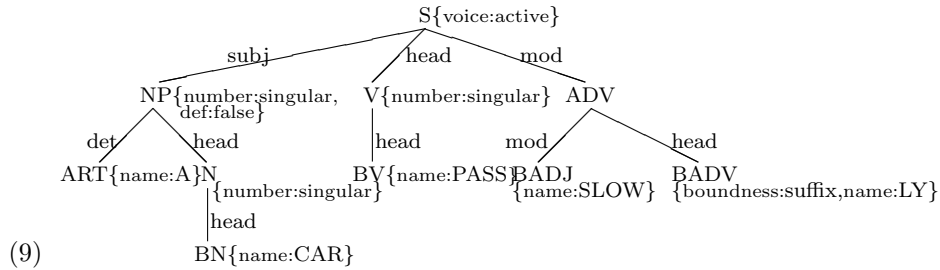
This latter S-tree is combined with the former S-tree by the rule *Rmod* (modification). Applying this rule yields the following S-tree:



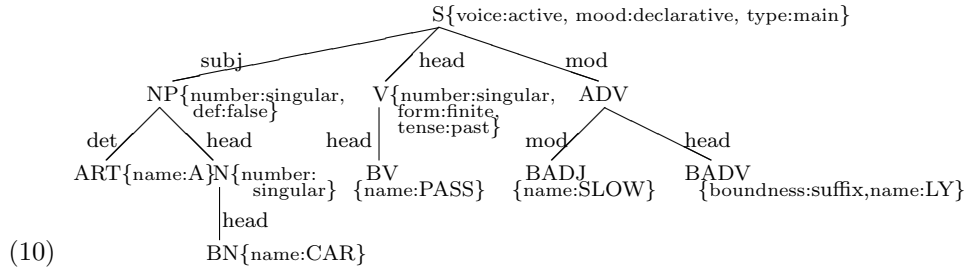
The M-rule R_{NP} forms an S-tree of category NP from the basic expression CAR (of category BN), puts the noun in singular and introduces the article A syntagmatically:



The substitution rule R_{subst} substitutes this NP for the variable with index 1, in S-tree (7) and makes the verb and the subject agree:



The M-Rule R_{past} puts the head verb into past tense, and the rule $R_{maindecl}$ determines that the sentence is a main declarative sentence:



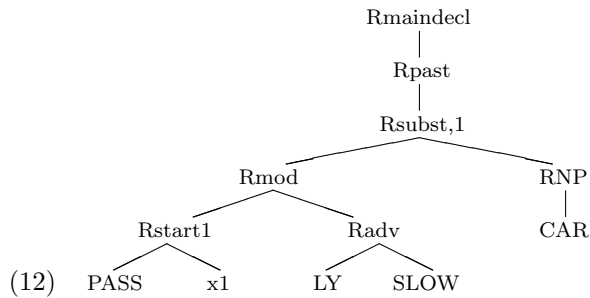
This surface tree can also be represented as in (11):

```

(11)  S  {voice:active, mood:declarative, type:main}
      [ subj / NP  {def:false, number:singular}
        [ det/ART{ name:A },
          head/N{number:singular}
            [ head/BN{name:CAR}]
          ],
        head / V   {form:finite, tense:past, number:singular}
        [ head/BV{name:PASS}]
        mod / ADV  { }
        [
          mod/BADJ{name:SLOW}
          head/BADV{boundness:suffix,name:LY}
        ]
      ]

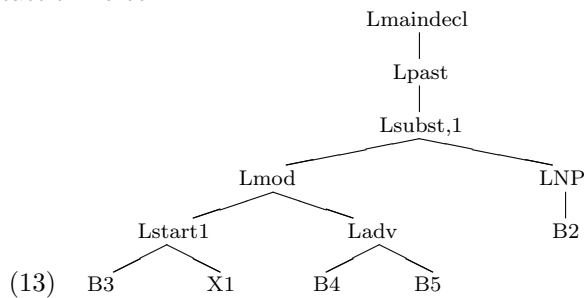
```

The corresponding D-tree is given in (12):



Note that the terminal nodes of D-tree (12) are labeled by unique names of basic expressions.

The semantic component can associate the semantic D-tree (13) to this syntactic D-tree:



The morphological component consists of a lexicon which relates unique names to strings and morphological rules. It turns each lexical S-tree from surface tree (11) into a string, while retaining the relative order of the S-trees. The lexical S-trees of the surface tree (11) are:

- (14) ART{name:A}
 N{number:plural}[head/BN{name:CAR}]
 V{form:finite, number:plural, tense:past}
 [head/BV{name:PASS}]
 ADV{ }
 [mod/BADJ{stem:SLOW},
 head/BADVboundness:suffix,stem:LY}
]

The first lexical S-tree is turned into the string *a*, since the morphological lexicon associates the unique name A to the string *a*. The second one is turned into the string *car*. The third one is changed into the string *passed* by applying morphological rules to form the past tense of the verb PASS which is associated with the string *pass*. The last lexical S-tree has internal structure. First, the parts of this lexical S-tree are turned into the strings *slow* and *ly* respectively by lexicon look-up, and then a rule combines these strings to form the string *slowly*. The sentence *a car passed slowly* has now been derived.

The example illustrated is very simple. In reality it is possible that the M-rules yield more than one S-tree (in the case of paraphrases). Furthermore, the morphological component can also yield several variants (e.g. if there are spelling variants, or if a word can be conjugated in more than one way).

2.4 Controlled M-grammars

The controlled M-grammar framework is an extension of the free M-grammar framework. Many properties of the controlled M-grammar framework are the same as in the free M-grammar framework: the grammar defines a set of surface trees, the grammar is compositional, it contains basic expressions and rules, the objects manipulated are S-trees, the rules to form S-trees are M-rules, the role of derivation trees remains the same, etc.

The properties of the controlled M-grammar formalism which are new in comparison to the free M-grammar formalism can be summarized as follows:

1. A controlled M-grammar consists of a set of *subgrammars*.
2. In a controlled M-grammar a special kind of rules, syntactic transformations, is distinguished. These rules do not contribute to meaning, i.e. they have identity as their meaning operation.
3. The application of rules can be controlled. It is possible to specify by means of a *control expression* in which order the rules and transformations must be applied, whether they are optional, obligatory, iterative, etc.

Subgrammars take a number of S-trees as input, apply a number of rules (including transformations) to them, and yield an S-tree as output. One input S-tree

has a special status. It is called the *head*. By convention, the head is always the first argument in M-rules. Subgrammars are defined by specifying: (1) a name for the subgrammar; (2) a characterization of the *head*; (3) a characterization of the S-tree delivered by the subgrammar, called its *export*; (4) a control expression, which specifies which rules must be applied, and in which order; (5) a characterization of the rest of the input S-trees, apart from the head, called the *import* of the subgrammar. Subgrammars can apply freely in the grammar. Whether they can be applied depends on whether the relevant S-trees match with the *head* and *import* of the subgrammars.

In the current grammars, the characterizations of the head, import and export of the subgrammars is done by specifying the category of the top node of the S-trees, but in principle more detailed characterizations are possible.

This can be illustrated with a simple example. One example subgrammar could have the name *PPformation*, take prepositions (P) as its head and yield prepositional phrases (PPs) as export. Its control expression could specify that several different kinds of rules must be applied in a certain order, e.g. rules to introduce complements of P, rules to introduce modifiers of P, rules to assign case to noun phrases (NPs) inside PPs, etc. (See below for an example control expression). The import for this subgrammar could be e.g. NPs, PPs and adverbial phrases (ADVPs). Given this subgrammar, one can form PPs from a head P by applying rules which add complements and modifiers, or change the structure in some way by M-rules as indicated in the control expression.

*Transformations*³ are normal M-rules in all respects, but they are associated with the identity operation as their meaning operation, i.e. the result of applying a transformation T to S-tree t with associated meaning τ is a new S-tree t' with the same meaning. Transformations always take exactly one S-tree as their input: they cannot combine several S-trees into a new S-tree. This follows from their semantics: the arguments of a rule always have a meaning, at least a basic meaning. The identity operation cannot combine two (or more) meanings into a new meaning, so the identity operation can only be associated to monadic rules. (see [Janssen, 1986]). In this respect, the controlled M-grammar framework is more restricted than non-compositional frameworks which allow rules which are not associated to a meaning operation: in such frameworks there would be no problem having rules which take more than one argument and which are not associated to any meaning, e.g. the generalized transformations as in [Chomsky, 1957] and [Chomsky, 1992]. See chapter 3 for further discussion.

Since the contribution of transformations to meaning is nihil, it is often useful to remove them from a derivation tree when they have applied. Syntactic D-trees from which transformations are removed are called *reduced D-trees*. When no ambiguity can arise, I will also use the simple term *D-tree* instead of *reduced D-tree*.

Control expressions are a restricted variant of regular expressions over names of M-rules. They are used to specify which M-rules are to be applied, and how.

³[Partee, 1973] was the first who proposed to incorporate transformations in a compositional grammar.

It is possible to specify in which order rules must be applied, whether rules are obligatory (they must apply), optional (they can but need not apply) or iterative (they can apply any number of times, including zero), and whether there are alternative rules that can be applied.

(15) is an example control expression:

$$(15) \quad (R_1) \cdot (R_2 \mid R_3) \cdot \{R_4\} \cdot [R_5 \mid R_6]$$

This control expression consists of a number of rules ($R_1 \dots R_6$) plus a specification how they must be applied. The dot (\cdot) stands for ‘followed by’, the symbol \mid separates alternatives, curly brackets indicate iterative rule application (i.e. zero or more times), square brackets enclose optional rules and round brackets enclose obligatory rules. The example control expression must be read as follows: first, apply the rule R_1 obligatorily, then apply either rule R_2 or rule R_3 . Next, apply rule R_4 as many times as one chooses, and finally apply either rule R_5 or rule R_6 or do not apply any rule at all.

Control expressions make the order in which rules must be applied explicit. In the free M-grammar formalism this was left implicit. It followed from the applicability conditions of the individual rules, but this turned out to require a more complex and often somewhat artificial formulation of many rules.

Surface trees are derived in a controlled M-grammar in the following manner.⁴ One starts with basic expressions: these are a special kind of S-trees, basic S-trees. One should find a subgrammar in which the head characterization matches with the basic S-tree. If they match, the rules from the subgrammar can be applied to the input S-tree, in accordance with the control expression specified in the subgrammar. If a rule requires more than one argument, it takes a derived S-tree formed by other applications of subgrammars or a basic S-tree as *import* of the subgrammar. When all rules of a subgrammar have applied, the subgrammar yields a derived S-tree as its *export*. This derived S-tree can then be used by other subgrammars, either as head or as import, to derive larger structures, etc. When a derived S-tree has a category from the set TOPCATS at its topnode, a surface tree has been derived.

2.5 Translation Relation

Up to this point the controlled M-grammar framework has been presented as a framework to describe natural language, and not as part of a machine translation method. Of course, the main emphasis in this book is on monolingual aspects, in particular how syntactically well-motivated analyses can be incorporated in a compositional framework, but it remains important to emphasize that the framework has been developed as part of a method for machine translation. In this section, it will be outlined how the controlled M-grammar framework can be used in the

⁴See [Appelo et al., 1987] and [Rosetta, 1993] for a more formal definition of controlled M-grammars.

adopted method of machine translation, the method of *compositional translation*, to define a translation relation between expressions derived by different grammars.

First, it must be emphasized that we are interested in defining linguistically *possible translations*, i.e. each utterance from language L_1 is associated with a set of utterances from language L_2 which are possible translations of the original utterance if we restrict attention to linguistic issues. General knowledge of the world, or knowledge of the particular situation might be relevant for choosing the correct translation, but this is not taken into account in the method described here.

The basic idea behind the method of compositional translation is the principle of *Compositionality of Translation* (see [Rosetta, 1993]). This principle states that

Two expressions are each other's translation if they are built up from parts which are each other's translation, by means of translation-equivalent rules.

This principle can be realized by tuning compositional grammars of different languages. 'Tuning' of two compositional grammars G_1 and G_2 means that for each basic expression from G_1 there is at least one translation-equivalent basic expression in G_2 , and that for each rule from G_1 there is a translation-equivalent rule in G_2 . Grammars which are attuned to each other in this way are called *isomorphic grammars*, and the approach in which such grammars are used is sometimes called the *isomorphic approach*. Note that transformations are not involved in tuning grammars: transformations need not have a translation-equivalent rule in other grammars. The notion *translation-equivalent* used here entails 'having the same meaning', but it covers not only meaning, but any aspect that might be relevant for translation. For discussion, see [Rosetta, 1993].

When two sentences are built up from parts which are each other's translation by means of translation-equivalent rules from two compositional grammars attuned to one another in the way indicated, they have the same semantic D-tree. The set of semantic D-trees is used as the interlingua. This interlingua is an automatic by-product of the tuned grammars, and it is a valid interlingua only for the languages dealt with.

It is now possible to derive sentences which are translations of each other in parallel, i.e. by recursively applying translation-equivalent rules to the results of their arguments (which are translations of each other), initially translation-equivalent basic expressions.

This will be illustrated by deriving the Dutch sentence *Er kwam langzaam een auto voorbij*, which is a translation of the sentence *a car passed slowly*, which has been illustrated before. The derivation of the Dutch sentence is fully parallel to the derivation of the English sentence, and they have the same semantic D-tree.

First, a second example grammar G_D is defined which describes a (tiny) fragment of Dutch. As before, the grammar has an illustrative purpose only, and no linguistic claims are made with it.

Table 2.2 specifies which syntactic categories and which grammatical relations can be used in the example grammar G_D . For each syntactic category it is spec-

RELATIONS	
det, head, mod, subj, part, top	
Type names	Possible Values
nametype	EEN, AUTO, VOORBLJKOMEN, LANGZAAM, AS, VOORBIJ, ER, ...
boundnesstype	stem, suffix
numbertype	singular, plural, unspecified
deftype	def, indef
boolean	false, true
voicetype	active, passive
moodtype	declarative, wh-interrogative, imperative
typetype	main, subordinate
formtype	finite, passpart, prespart, infinitive
tensetype	present, past, unspecified
indextype	1,2,3,...
Syntactic Categories	& appropriate AV-pairs
ART	stem: nametype
ADV	r: boolean
BADJ	stem: nametype
BADV	stem: nametype, boundness: boundnesstype, r: boolean
BN	stem: nametype
BV	stem: nametype, particle: nametype
N	number: numbertype
PART	stem: nametype
NP	number: numbertype, def: deftype
S	voice: voicetype, mood: moodtype, type: typetype
V	number: numbertype, form: formtype, tense: tensetype
VAR	index: indextype
LEXICAL CATEGORIES	
ART, ADV, N, V, PART	

Table 2.2: Relations, categories and A-V-pairs of example grammar G_D

ified which attributes are appropriate, and what their possible values are. The possible values of an attribute are indicated by a name for a type, which is defined separately.

There are many similarities with the example grammar G_E given earlier, but also some differences. Grammar G_D has different basic expressions, it has a category PART and a grammatical relation *part*. The possible values of the attribute *form* differ slightly, there is an attribute *r* for the categories BADV and ADV, and an attribute *particle* for the category BV.

The grammar contains the M-rules listed in (16). These rules will not actually be formulated here. An informal characterization of their task is given below. The names of the meaning operations associated with these rules are also given:

(16)

Rule Name	Meaning Name
DRstart1	Lstart1
DRadv	Ladv
DRmod	Lmod
DRsubst	Lsubst
DRpast	Lpast
DRmaindecl	Lmaindecl
DRNP	LNP

The basic expressions of the example grammar are given in (17) together with the names of their associated basic meanings:

(17)

Basic Expression	Basic Meaning
BADV{ stem:ER, boundness:stem, r: true }	B1
BN{ stem:AUTO }	B2
BV{ stem:VOORBIJKOMEN, particle:VOORBIJ }	B3
BADV{ stem:AS, boundness:suffix, r:false }	B4
BADJ{ stem:LANGZAAM }	B5

Apart from the basic expressions mentioned in (17) an infinite number of variables $VAR\{index:1\}$, $VAR\{index:2\}$, ... with associated basic meanings X_1, X_2, \dots also belong to the basic expressions of the example grammar. The names of these basic expressions are represented as x_1, x_2, \dots

Apart from basic expressions, grammar G_D has two expressions which are introduced syncategorematically:

(18)

ART{ stem:EEN }
PART{ stem:VOORBIJ }

Note that the example grammars G_D and G_E are isomorphic: for each basic expression from one grammar there is a translation-equivalent basic expression in the other, and for each rule from one grammar there is a translation-equivalent rule in the other.⁵ This is illustrated in (19):

⁵In this example, this is, for simplicity, a 1-1 relationship, but with realistic grammars this relation need not be 1-1.

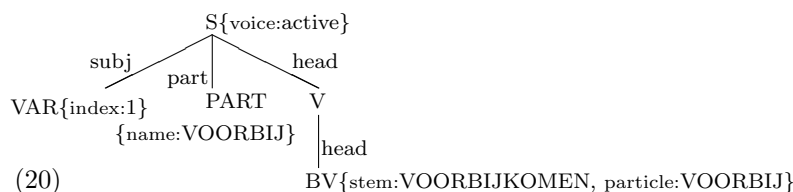
(19)

Rule Name G_D	Meaning Name	Rule name G_E
DRstart1	Lstart1	Rstart1
DRadv	Ladv	Radv
DRmod	Lmod	Rmod
DRsubst	Lsubst	Rsubst
DRpast	Lpast	Rpast
DRmaindecl	Lmaindecl	Rmaindecl
DRNP	LNP	RNP
Basic Expression G_D	Basic Meaning	Basic Expression G_E
ER	B1	THERE
AUTO	B2	CAR
VOORBIJKOMEN	B3	PASS
AS	B4	LY
LANGZAAM	B5	SLOW
x_i	X_i	x_i

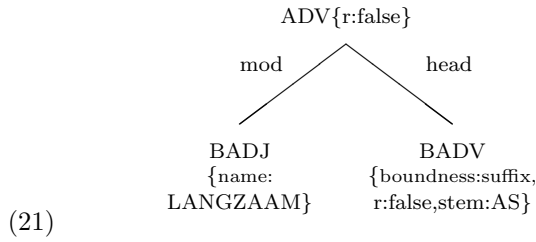
We can now describe how the sentence *er kwam langzaam een auto voorbij* can be derived in a way which is isomorphic to the derivation of the sentence *a car passed slowly*.

As before, the syntactic component derives a surface tree for this sentence. By recording how the S-tree has been derived, a D-tree is created, which can be mapped onto one or more semantic derivation trees in the semantic component. The morphological component turns the sequence of lexical S-trees occurring in the surface tree into a string.

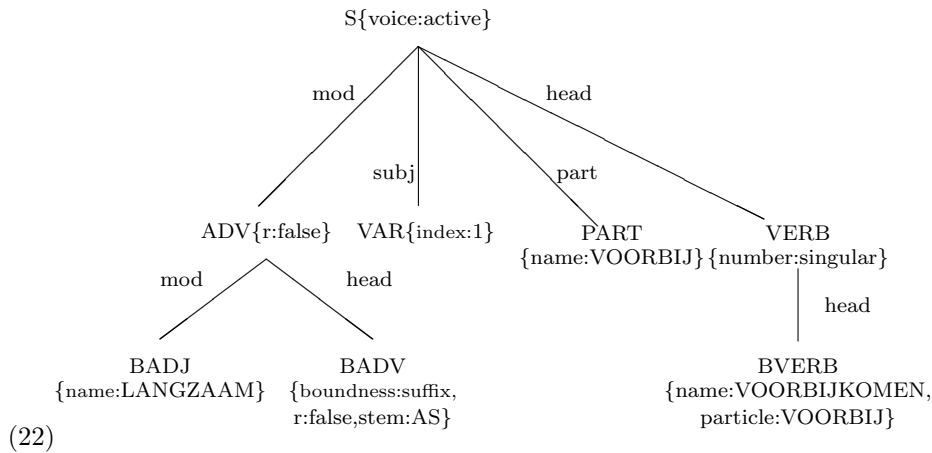
We start by applying M-rules. Start rule (*DRStart1*) is applied to the basic expression $BV\{\text{stem:VOORBIJKOMEN, particle:VOORBIJ}\}$ and a variable $VAR\{\text{index:1}\}$ to form an active sentence. In addition the rule syncategorematically introduces the particle VOORBIJ, triggered by the attribute *particle* of the verb. This yields the following S-tree:



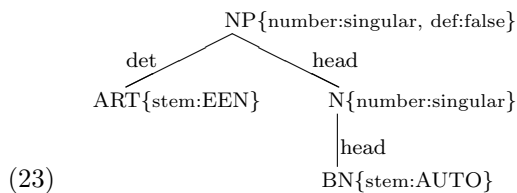
In addition, an adverb is created from the basic expression LANGZAAM (of category BADJ) and the suffix AS by the M-rule *DRadv*. This yields:



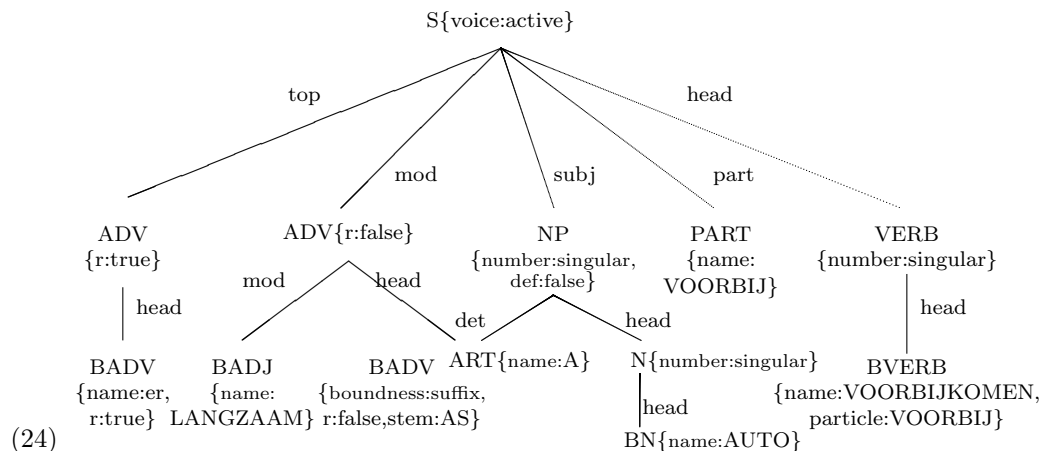
This latter S-tree is combined with the former S-tree by the rule *Rmod* (modification). Applying this rule yields the following S-tree:



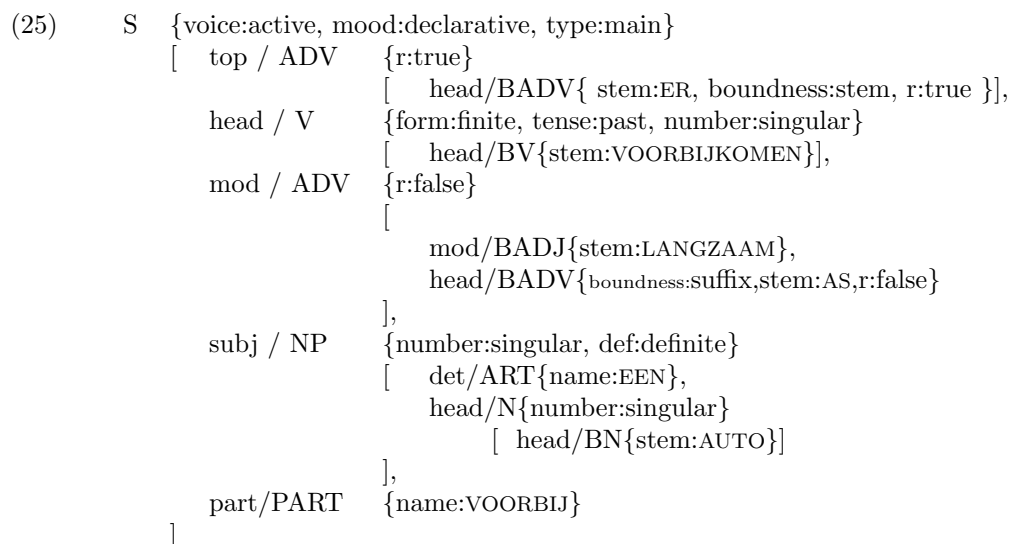
The M-rule *RNP* forms an S-tree of category NP from the basic expression AUTO (of category BN), putting the noun in singular and syncategorematically introducing the article EEN:



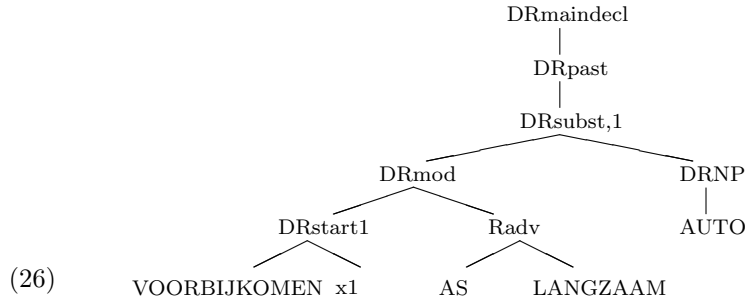
The substitution rule *Rsubst* substitutes this NP for the variable with index 1, it makes the verb and the subject agree, and it introduces the adverb ER syncategorematically in the topic position. The rule introduces this element if the subject NP is indefinite:



The M-Rule *R_{past}* puts the head verb into past tense. The rule *R_{maindecl}* determines that the sentence is a main declarative sentence and puts the finite verb in the ‘second position’. At this point, surface tree (25) has been derived:



The corresponding D-tree is given in (26):



Note that D-tree (26) has the same geometry as (12). They are *isomorphic*: only the names on the nodes differ, but in other respects the trees are identical.

The semantic component can associate the semantic D-tree (13), given earlier when the sentence *a car passed slowly* was derived, to this syntactic D-tree. The sentences *a car passed slowly* and *er kwam langzaam een auto voorbij* share a semantic D-tree. Since they share a semantic D-tree, they are built up from parts which are each other's translation, by means of translation-equivalent rules. As a consequence, these sentences are characterized as translations of each other by the example grammars.

The morphological component turns each lexical S-tree from surface tree (25) into a string, while retaining the relative order of the S-trees. The lexical S-trees of the surface tree (25) are:

- (27)
- ```

ADV {r:true}
 [head/BADV{name:ER, boundness:stem, r:true}]
V{form:finite, number:plural, tense:past}
 [head/BV{stem:VOORBIJKOMEN, particle:VOORBIJ}]
ADV{r:false}
 [mod/BADJ{stem:LANGZAAM},
 head/BADV{boundness:suffix,stem:AS,r:false}
]
ART{stem:EEN}
N{number:singular}[head/BN{stem:AUTO}]
PART{name:VOORBIJ}

```

The first lexical S-tree is turned into the string *er*. The second is changed into the string *kwam* by applying morphological rules to form the past tense of verbs. The third lexical S-tree has internal structure, but the suffix is an abstract expression, i.e. an expression which does not correspond to a string. The rules to form adverbs turn this lexical S-tree into the string *langzaam*. The fourth lexical S-tree is turned into the string *een*, and the fifth into the string *auto*. Finally, the last lexical S-tree is turned into the string *voorbij*. The sentence *er kwam langzaam een auto voorbij* has now been derived.

The example grammars illustrate the use of basic expressions, of expressions which are introduced syncategorematically, and of abstract basic expressions. Ba-

sic expressions such as AUTO, CAR play a role both in semantics and in syntax, and have a corresponding string representation. A basic expression such as AS in  $G_D$  also plays a role in syntax and semantics, but it has no corresponding string representation. It is called an abstract basic expression. Certain expressions are introduced syncategorematically, e.g. A, EEN, VOORBIJ and ER. Expressions which are introduced syncategorematically only play a role in syntax, but not in semantics. They are not represented in the syntactic D-tree, nor in the semantic D-tree, but only in S-trees. The expressions EEN and VOORBIJ are always introduced syncategorematically. But the expression ER is a basic expression, and is introduced as a basic expression in some structures (e.g. in a sentence such as *De auto kwam er voorbij* ‘The car passed there’), but syncategorematically in others (e.g. in the example derived).

Tuning grammars implies that the organization and form of each grammar is influenced by the other grammars that it is tuned to. This might lead to ‘unnatural’ grammars, but it is possible to avoid the most undesirable aspects of tuning grammars by using transformations which are not subject to tuning: phenomena which are purely syntactic peculiarities of a specific language can be dealt with by transformations.

## 2.6 Translation System

In the preceding sections a framework was presented which abstractly and relationally characterizes languages and translation relations which hold between them, but it was not discussed how this framework is to be used in a real language processing system.

In this section, I will show how the framework adopted can form the basis for an actual machine translation system (section 2.6.1), and I will illustrate this system with an example translation (section 2.6.2).

### 2.6.1 Description of the System

If the framework presented is to be used as an actual natural language processing system, it must be possible to use the framework both for analysis and for generation. This requirement is satisfied: controlled M-grammars are *reversible*. A controlled M-grammar  $G$  defines a relation  $R$  between semantic derivation trees and utterances, and two functions can be (automatically) derived from it: one function yields, given an utterance  $s$ , all semantic derivation trees  $D$  for which  $\langle D, s \rangle \in R$  holds; the other function yields, given derivation tree  $D$ , all utterances  $s$  for which  $\langle D, s \rangle \in R$  holds.

An additional constraint, effectiveness, requires that these functions yield their result in finite time, for all input. We will discuss below how this can be achieved. If a grammar satisfies both these requirements, it is *effectively reversible*.

I will describe here informally what form these functions take (for more details see [Landsbergen, 1987], [Appelo et al., 1987], [Rosetta, 1993]), and I will show



how the Dutch example sentence dealt with above can be translated into English by the system.

A controlled M-grammar defines two functions, ANALYSIS and GENERATION. ANALYSIS takes strings of characters as input and yields semantic D-trees. GENERATION takes semantic D-trees as input and yields strings of characters. The functions ANALYSIS and GENERATION are composed out of a number of functions, each of which is derived from the components of a controlled M-grammar.

The two functions derived from the semantic component are called A-TRANSFER (for analysis) and G-TRANSFER (for generation), respectively. A-TRANSFER takes as input syntactic D-trees and yields semantic D-trees, G-TRANSFER takes as input semantic D-trees and yields syntactic D-trees.

The syntactic component yields two pairs of functions. First, there is the pair M-PARSER (for analysis) and M-GENERATOR (for generation). M-PARSER takes as input S-trees and yields syntactic D-trees. M-GENERATOR takes as input syntactic D-trees and yields S-trees. Second, there is the pair S-PARSER (for analysis) and LEAVES (for generation). S-PARSER takes a sequence of lexical S-trees as input, and yields S-trees. LEAVES takes S-trees as input and yields a sequence of lexical S-trees. Note that M-GENERATOR and M-PARSER relate syntactic D-trees to S-trees. The morphological component relates sequences of lexical S-trees to symbol strings. Thus, an additional module is required to relate the S-tree to a sequence of S-trees. In generation, this is no problem: the module LEAVES picks out the lexical S-trees which occur in the S-tree. In analysis, a more complicated module is required. In fact, a separate grammar, called surface grammar, defines a set of surface trees. S-PARSER uses this grammar to turn a sequence of lexical S-trees into a surface tree. The surface grammar is in essence a context-free grammar, which makes it possible to use efficient parsing algorithms. However, the surface grammar defines a superset of the language defined by M-GENERATOR and M-PARSER, so that S-PARSER will yield a set of candidate surface trees for an input sentence. This set always contains the surface trees which are accepted by M-PARSER, but it also contains surface trees which will not be accepted by M-PARSER. The surface grammar has no principled status. Its main function is to insure that a set of candidate S-trees can be found by S-PARSER efficiently.

The morphological component defines two functions, A-MORPH and G-MORPH. A-MORPH takes a string of characters as input and yields a sequence of lexical S-trees. G-MORPH takes a sequence of lexical S-trees as input and yields a string of characters.

Each of these functions is set-valued. If the set contains more than one element in analysis, we speak of *ambiguity*. If the result of a function is the empty set, analysis (or generation) blocks.

ANALYSIS is the ‘composition’ of the functions A-MORPH, S-PARSER, M-PARSER and A-TRANSFER. GENERATION is the ‘composition’ of the functions G-TRANSFER, M-GENERATOR, LEAVES and G-MORPH. The set-valued na-

ture of the functions, however, is taken into account:<sup>6</sup> each element from the resulting set of the preceding function is offered as input to the next function.

As pointed out before, ANALYSIS and GENERATION must be effective functions. This can be achieved by guaranteeing that the grammars satisfy *measure conditions*. These measure conditions must guarantee that the rules in an iterative rule class can be applied only a finite number of times. See [Appelo et al., 1987] and [Rosetta, 1993] for discussion.

Each of the functions is implemented as a module in the system. A graphical representation of a complete system, with 8 modules is given in figure 2.1.

## 2.6.2 Example Translation

In this section I will illustrate how the translation system operates by running the example sentences given earlier through the system. The grammars described above characterized the sentences *A car passed slowly* and *Er kwam langzaam een auto voorbij* as translations of each other. It will be shown that inputting the first sentence (*A car passed slowly*) to the ANALYSIS function derived from  $G_E$ , followed by the GENERATION function derived from grammar  $G_D$  yields the second sentence (*Er kwam langzaam een auto voorbij*) as a result.

The input string is *A car passed slowly*. The first module to apply is A-MORPH. A-MORPH analyzes each word from this string and turns it into a lexical S-tree. If a word can be analyzed in more than one way, A-MORPH will yield two or more lexical S-trees for it. The output of A-MORPH is a set of sequences of lexical S-trees. In the example sentence, each word except *passed* can be analyzed only in one way. The word *passed* can be a past tense form or a participle, so that two lexical S-trees are associated with this string. In this way, A-MORPH introduces an ambiguity. The ambiguity will turn out to be a temporary ambiguity in this example, since it is resolved before we reach semantic D-trees. If the input string contains an unknown word (e.g. a misspelling), A-MORPH will yield the empty set as a result, and the derivation blocks.

The two sequences of lexical S-trees have been represented in (28):

---

<sup>6</sup>This is why the word *composition* was quoted in the preceding sentences.

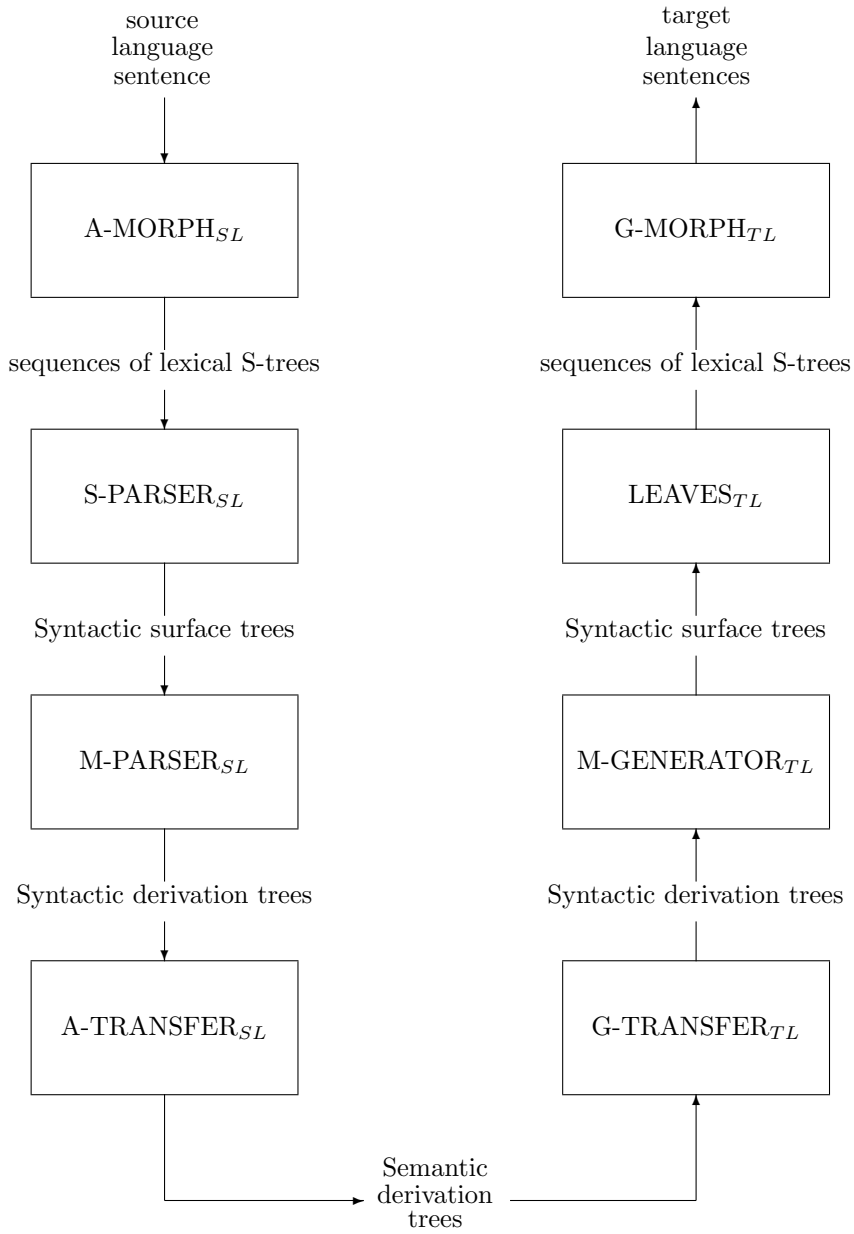


Figure 2.1: The 8 modules of the Rosetta machine translation system

- (28) a Result 1:  
 ART{name:A}  
 N{number:singular}[head/BN{name:CAR}]  
 V{form:finite, number:plural, tense:past}  
 [ head/BV{name:PASS}]  
 ADV{ }  
 [mod/BADJ{stem:SLOW},  
 head/BADV{boundness:suffix, stem:LY}  
 ]
- b Result 2:  
 ART{name:A}  
 N{number:singular}[head/BN{name:CAR}]  
 V{form:participle, number:unspecified, tense:unspecified}  
 [ head/BV{name:PASS}]  
 ADV{ }  
 [mod/BADJ{stem:SLOW},  
 head/BADV{boundness:suffix, stem:LY}  
 ]

The sequences of lexical S-trees are each input to the next module, S-PARSER. S-PARSER attempts to construct a surface tree with the lexical S-trees as its 'leaves'. It can only yield a result for the sequence of lexical S-trees in which the lexical S-tree for *pass* is a past tense form. So, the ambiguity introduced by A-MORPH is immediately resolved in S-PARSER. S-PARSER can resolve ambiguities created by A-MORPH, but it can also create ambiguities of its own. Such ambiguities are called surface structural ambiguities. No ambiguities are created in this example, but this is a quite exceptional situation. As with A-MORPH, it is possible that S-PARSER cannot find a single possible analysis for its input sequences of lexical S-trees. In that case, S-PARSER yields the empty set as a result, and the analysis process blocks. The result of S-PARSER for the example sentence is represented in (11) above.

The result of S-PARSER is input to the module M-PARSER. M-PARSER attempts to break down the surface tree into smaller S-trees, ultimately into basic S-trees, by applying M-rules 'in reverse'. It is recorded in a syntactic D-tree which M-rules apply and which basic expressions occur. The output of M-PARSER is a set of syntactic D-trees, which describe how the surface structures have been analyzed. M-PARSER can resolve ambiguities created earlier, or it can create its own ambiguities. Such ambiguities are called derivational ambiguities. If M-PARSER cannot associate any syntactic D-tree to its input surface trees, it yields the empty set, and the derivation blocks. This is not applicable in this simple example, where M-PARSER yields only the single syntactic D-tree represented in (12) above.

The next module, A-TRANSFER, takes this syntactic D-tree as its input and yields a set of semantic D-trees by mapping the names of basic expressions onto the names of their basic meanings, and by mapping the names of the M-rules

onto the names of their associated meaning operations. A-TRANSFER can create ambiguities, in particular if a word has two or more meanings with the same syntactic and morphological properties. Ambiguities created in A-TRANSFER are called *semantic ambiguities*. A-TRANSFER can also block derivations, e.g. it will not yield a semantic D-tree for the syntactic D-tree of a sentence such as *Hij gaf er de brui aan* ('He quit') under a literal interpretation, since *brui* has no literal meaning.

A-TRANSFER turns the syntactic D-tree (12) into the semantic D-tree (13).

Having dealt with the function ANALYSIS of  $G_E$ , we now turn to the function GENERATION of  $G_D$ .

The semantic D-tree (13) is turned into a syntactic D-tree of  $G_D$  by G-TRANSFER. G-TRANSFER replaces each name of a basic meaning by the name of the corresponding basic expression, and each meaning operation name by a rule name. G-TRANSFER generally creates multiple output: one meaning operation name usually corresponds to several rule names. The result of applying G-TRANSFER to the example is represented in (26) given above.

The next step is M-GENERATOR. M-GENERATOR takes as input a (reduced) syntactic D-tree and turns it into one or more surface trees by applying the rules occurring in the syntactic D-tree and transformations in accordance with the control expressions. If one syntactic D-tree corresponds to several surface trees, we speak of paraphrases. The result of applying M-GENERATOR is represented in (25)

This surface tree is input to LEAVES, which takes the subtrees dominated by a node with a lexical category out of the surface tree. The result is given in (27) above.

Finally, G-MORPH applies to this sequence of lexical S-trees, turning each lexical S-tree into a character string. The result is the sentence *Er kwam langzaam een auto voorbij*.

In this way it is possible to use the controlled M-grammar as the basis for an actual natural language machine translation system.

## 2.7 Some Specific Assumptions

The preceding sections described the M-grammar formalism. This formalism entails only a limited number of linguistic claims (e.g. that S-trees play a role in adequate linguistic descriptions). The formalism itself should not be seen as a theory or a description of language: it is a formalism to express such a theory or description in. The formalism is compatible with a wide variety of linguistic theories and descriptions. (See [Shieber, 1987] for the distinction between these notions).

In this section I will discuss a number of specific assumptions which have been made with respect to the grammars used in the system. These assumptions are not part of the formalism, but are a specific linguistic description expressed in the formalism.

A fairly standard set of syntactic categories has been assumed in the grammars of the Rosetta3 system, though certain adaptations to the framework have been made and certain assumptions have been made explicit.

For the familiar lexical categories *noun*, *verb*, *adjective*, and *adverb* three series of categories are used.

The categories *BN* (*basic noun*), *BV* (*basic verb*), *BADJ* (*basic adjective*) and *BADV* (*basic adverb*) are used for entries of the respective categories as they occur in the lexicon that cannot be analyzed further.

The categories *SUBN* (*subnoun*), *SUBV* (*subverb*), *SUBADJ* (*subadjective*) and *SUBADV* (*subadverb*) are used for the results of compounding and derivation rules. These categories did not appear in the example grammars for expository convenience.

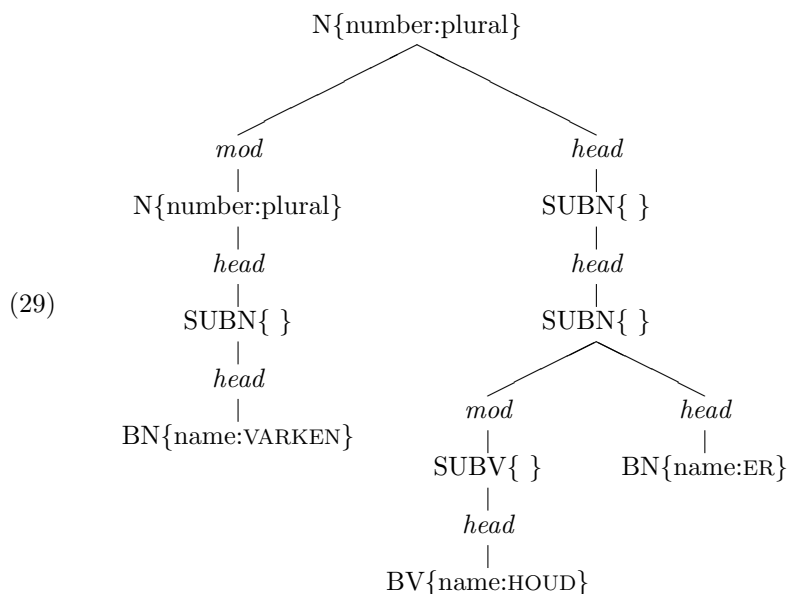
The categories *N* (*noun*), *V* (*verb*), *ADJ* (*adjective*) and *ADV* (*adverb*) are used for (possibly inflected) words.<sup>7</sup>

In an example such as *varkenshouders* (lit. ‘pigs keepers’) *varken* ‘pig’ is a BN and *houd* ‘keep’ is a BV. The word *varkens* is plural N. The part *houder* is a SUBN created by derivation from the BV *houd* and the suffix BN *-er* ‘-er’. The part *varkenshouder* is a SUBN derived by compounding the N *varkens* and the SUBN *houder*. The whole word is a plural N.<sup>8</sup> The structure is represented in (29):

---

<sup>7</sup>With regard to the categorial distinction between adjectives and adverbs in Dutch, we followed traditional grammar. It is necessary anyway to distinguish between words such as *graag*, which can be used only adverbially, and other words which can be used predicatively and adverbially. This distinction is encoded categorially in the Rosetta grammars.

<sup>8</sup>This description might suggest that the word-internal structure subdivides the actual string. In reality, the relation between the word-internal structure and the string is much more complex, e.g. inflectional suffixes are not represented in the structure, see e.g. the lexical S-tree in (29). This will not be discussed any further, since it is not relevant here.



Apart from the major lexical categories a number of minor categories have been assumed, of which *PREP* (*adposition*)<sup>9</sup>, *ART* (*article*) and a series for different kinds of pronouns are the most important.

Concerning phrasal categories, a familiar set consisting of *NP*, *VP*, *ADJP*, *ADVP* and *PP* has been assumed. In addition, a less familiar set of categories *NPP*, *VPP*, *ADJPP*, *ADVPP* and *PPP* has been assumed for ‘small clauses’ (see chapter 4).

Though the formalism does not force one to assume headed structures, it has been assumed that most phrasal categories have an element which bears the grammatical relation *head*, though these need not be present in surface trees. In surface trees, no lexical or phrasal categories can occur which do not dominate anything. Therefore, headless structures must be allowed in surface trees to deal with examples such as *hij zag er gisteren twee* ‘he saw two yesterday’, in which *twee* is a headless NP.

For sentences, the categories *CLAUSE* and *S* (*sentence*) are used. A *CLAUSE* represents a partially derived sentence which is overtly marked for tense and aspect, but to which mood (declarative, yes-no interrogative, etc.) has not yet been assigned.

Utterances have the category *UTT*. An utterance can dominate any phrasal category.

The set *LEXCATS* (lexical categories) contains the categories *N*, *V*, *ADJ*, *ADV*, *PREP*, *ART* and a number of other minor lexical categories.

The set *TOPCATS* consists of the category *UTT*.

<sup>9</sup>This category is used for prepositions, post-positions and circumpositions in Dutch. A separate attribute is used to distinguish among these.

Each category has a set of appropriate attribute-value pairs. I will not discuss these here, but introduce the relevant ones where appropriate in the chapters to follow.

In S-trees, grammatical relations must be specified. In the grammar implemented, however, grammatical relations are often used in a different manner than is usual in most grammatical frameworks. The relations have been used as names for *qualified positions*, i.e. relations are used to indicate named slots. A partial order has been defined on relations which corresponds to left-right order of phrases bearing these relations. In cases where this partial order does not apply, certain positions are distinguished purely by name (hence the term *qualified positions*) where necessary.

This can be illustrated with the following example. Traditional grammar would say that the boldface phrases are *subjects* in each of the following examples from Dutch:

- (30) a las **de man** een krant?  
       ‘Is the man reading a newspaper’  
       b **de man** las een krant  
       ‘The man is reading a newspaper’  
       c Er las **een man** een krant  
       ‘Someone was reading a newspaper’  
       d ‘Er werd **een krant** gelezen  
       ‘A newspaper was read’

In the Dutch grammar which is implemented, however, only the bold face phrase in (30a) bears the grammatical relation *subject*. The other phrases bear different grammatical relations, *shift* (in (30b)), *postsubject* (in (30c)) and *object* (in (30d)). These different grammatical relations correspond to different positions. Thus, a phrase bearing the relation *shift* always precedes a phrase bearing *subject*, and this in turn precedes a phrase which bears the relation *postsubject*, which in its turn precedes any phrase bearing the relation *object*. The reasons for assigning these grammatical relations will not be given here,<sup>10</sup> but it must be pointed out that the position a phrase occupies co-determines which grammatical relation it bears.

A number of grammatical relations (from Dutch) which play a role later will be introduced here. Topicalized phrases and interrogative phrases occupy a sentence-initial position and bear the grammatical relation *shift*. Subordinate conjunctions, and finite verbs when in ‘second position’ occupy a position called *conj*. Immediately to the right of *conj* there is a position called *subject*. It is occupied by subjects or by expletive elements serving as ‘anticipatory subjects’. Indefinite subjects can occupy a position further to the right. If they occur there, the subject position is usually occupied by an expletive element (cf. the presence of *er* in (30c)). Direct

<sup>10</sup>The specific assignment is inspired by recent analyses of the structure of the sentence in Dutch. Compared with theories developed in the Principles and Parameters framework (see [Chomsky and Lasnik, 1991] for a recent description of this theory), one can basically equate *shift* to SpecCP, *subject* to SpecIP, *postsubject* to SpecVP, and *object* to [NP, X].



objects can occur to the left of sentential adverbs in Dutch. An example sentence is given in (31):

- (31) Hij heeft het boek waarschijnlijk niet gelezen  
He has the book probably not read  
'He has probably not read the book'

There is potentially an (in principle) unlimited number of positions to the left of each adverb. These positions are called *preadv* positions.

## 2.8 Concluding Remarks

In this chapter the grammatical formalism adopted has been outlined, and some more specific linguistic assumptions have been made.

One of the virtues of this framework is that it is possible to write grammars from a generative point of view which can also be used for analysis. There are only two conditions which must be obeyed to achieve this effect. First, the grammar rules must be reversible. This does not give problems in practice. In many cases, the rule notations have been chosen to automatically guarantee reversibility.

Second, the grammars must satisfy the measure condition (see section 2.6.1). Our experience is that in practice this condition hardly ever causes problems.

If the grammars satisfy these two conditions, they can be used both for generation and analysis. For this reason, there is no need to discuss both aspects separately in this book. We can, and will, restrict ourselves to grammatical descriptions from a generative point of view, and ignore aspects regarding analysis. Only in exceptional cases, e.g. when S-PARSER is described, will analytical aspects be discussed.



## Chapter 3

# Compositionality and Syntactic Generalizations

### 3.1 Introduction

In the preceding chapter a special kind of compositional grammar, controlled M-grammar, was introduced to describe the form and meaning (or, in the context of machine translation: translational equivalence) of expressions. The basic idea of compositional grammars, e.g. Montagues PTQ ([Thomason, 1974]), is that they relate form and meaning in a very direct way. Such grammars usually have a strong semantic bias, in the following sense: in a compositional grammar, basic expressions are associated with a basic meaning, syntactic rules are associated with a meaning operation, and it is the proper composition of the meaning of an expression which determines how many and which syntactic steps to form an expression are distinguished. This semantic bias is essential given the translation method adopted, in which isomorphic grammars are used to define a translation relation, but it might lead to a grammar which is less adequate from a syntactic point of view.

The introduction of *syntactic transformations* in controlled M-grammars, makes it possible to overcome this disadvantage. Syntactic transformations do not contribute to the meaning and are not involved in the isomorphy relation.

In this chapter (1) I will supply additional evidence that such transformations are essential to achieve an adequate modularization of the syntax, so that optimal descriptions of syntactic phenomena and their associated meanings can be obtained; (2) I will describe in detail for a number of constructions which rules must be meaningful rules and which rules must be transformations, and why; (3) I will show that the addition of transformations makes the relation between form and meaning considerably more indirect than one might expect.

This does not imply that the relation between form and meaning is unclear in the grammar. However, it will become evident that in many cases the rules which

account for the form differences are not directly associated to the related meaning operations. Often, the meaning operation is associated with a rule which changes the expression only at an abstract level. Two situations are typical, though they do not necessarily exhaust the possibilities: (1) the application of a meaningful rule to an abstract structure leaves a trace, which is used to trigger syntactic transformations that actually perform the relevant operations on the form, and (2) syntactic transformations apply blindly, creating a number of alternative structures, after which meaningful rules check whether the right configurations have been created for them to apply. The form changes performed by these meaningful rules tend to be rather small and often have no direct reflex on the surface.

In short, I attempt to contribute to the search for an optimum in the relationship between syntax and semantics, so that complexity and redundancy can be minimized in syntax, in semantics and in their mutual relationship.

The emphasis in this chapter is on an investigation into the question which rules are required to adequately describe certain syntactic constructions, and which of these rules must be meaningful rules and which must be transformations. In chapter 5 and 6 we will concentrate on how certain transformations must be formulated, which formal operations are required to obtain adequate syntactic descriptions, and how the resulting analyses relate to analyses in more theoretically oriented frameworks.

In section 3.2 the problem is introduced in more detail using the syntax of auxiliaries and inversion in English as an illustration. The solution proposed is outlined in general terms, and applied to the specific case at hand. In the sections to follow it is shown that this kind of problem occurs quite frequently, and that the solution proposed is applicable to all of these problems. Phenomena relating to mood in Dutch (section 3.3), variants of sentences which differ only in linear order (section 3.4), phenomena relating to *wh*-movement (section 3.5), and the proper treatment of genericity (section 3.6) are used to argue this.

## 3.2 Auxiliaries and Inversion in English

In this section I will introduce the general problem using the peculiar properties of English auxiliaries as an example (section 3.2.1). A general characterization of the solution to this kind of problem, and an instantiation of this kind of solution, specific to the problem at hand, will be described in section 3.2.2.

### 3.2.1 The Problem

Let us start by considering the following two sentences

- (32) a He bought a book.  
b Did he buy a book?

At first sight, these sentences appear to be excellent examples for an analysis within a compositional approach. The two sentences clearly differ in meaning

(declarative sentence versus yes-no question) and corresponding to this meaning difference there is a difference in form. Sentence (32a) contains a sentence-initial subject followed by an inflected main verb, while sentence (32b) contains a sentence-initial inflected auxiliary verb, followed by the subject, followed in its turn by the main verb in its basic form.<sup>1</sup> It is fairly straightforward how two different rules can be constructed, each associated with its own meaning, and each carrying out different formal changes to an input structure. Let us assume that certain rules have derived a structure which I will represent informally with  $S(\textit{he bought a book})$ . The properties common to (32a,b) have been expressed in this structure, which will be the input for the rules to form declarative and yes-no interrogative sentences. One rule, which will be called RDECL, will take this structure as its input, and it simply resets the value of the attribute *mood* of the S-node dominating the structure from *unspecified* to *decl*. The rule is associated with the meaning of declarative sentences. A second rule, which will be called RYNQ, also takes this structure as input, and it effects the following changes: the attribute *mood* of S is changed from *unspecified* to *ynq*, the auxiliary verb *do* is introduced in a sentence-initial position, the tense features are copied from the main verb onto this auxiliary, and then deleted from the main verb. This rule is associated with the meaning of yes-no questions.

If we extend the grammar with additional rules, however, it soon becomes clear that the simple approach outlined above runs into severe trouble. I will extend the grammar with a number of rules using the compositional approach as indicated, to illustrate in what ways such an analysis is unsatisfactory.

Consider the following sentence:

(33) Which book did he buy?

This sentence differs from the sentences discussed above in several respects. Let us compare it with (32a). First, there is a semantic difference: (32a) has the meaning of a declarative sentence, (33) has the meaning of a wh-question. Second, corresponding to this semantic difference, there is a formal difference: sentence (33) has a sentence-initial NP, with the determiner *which* instead of the article *a*; this NP is followed by an inflected auxiliary verb, which is followed by the subject, and the main verb in its basic form is sentence-final. Again, it is quite straightforward to write a rule to form such sentences, and to associate this rule with its own meaning. If we assume that the input for the rule is  $S(\textit{he bought which book})$ , the rule, which will be called RWHQ, has to perform the following formal changes: the attribute *mood* of the S-node must be reset from *unspecified* to *whq*; the NP *which book* must be preposed to a sentence-initial position, the auxiliary verb *do* must be inserted, the tense features of the main verb must be copied onto this auxiliary verb, and then these tense features must be deleted from the main verb.

Let us add a rule to introduce adverbs. This rule will be called RADV. It takes a sentence and an adverb of a certain kind (e.g. *never*) as input, and inserts the

---

<sup>1</sup>Details relating to punctuation will be ignored.

adverb into the sentence immediately before the main verb. Given  $S(\textit{he bought a book})$  and the adverb *never* as input, we can derive  $S(\textit{he never bought a book})$ . This structure can then be input to the rules introduced earlier to derive the sentences *he never bought a book* and *did he never buy a book* and from a structure such as *he never bought which book* we can derive *which book did he never buy?*.

Let us add a further rule. Consider the sentence in (34):

(34) Never did he buy a book.

As before, this sentence differs from the sentences discussed earlier in form and in meaning, or at least in pragmatic aspects (e.g. topic-comment partitioning). And, again, it is fairly easy to write a rule which can take these phenomena into account. Let us call this rule *RTOP*. One possible input for this rule is  $S(\textit{he never bought a book})$ . The rule must perform the following changes: the adverb *never* must be preposed to a sentence-initial position, the auxiliary verb *do* must be inserted to the right of it, the tense features of the main verb must be copied onto this auxiliary verb, and then the tense feature must be deleted from the main verb.<sup>2</sup>

We continue to add rules. Let us assume a rule to introduce the adverb *not*. This rule (*RNEG*) takes a sentence as input, e.g.  $S(\textit{he bought a book})$ , and it performs the following formal changes: the adverb *not* is introduced, the auxiliary *do* is introduced, the tense features are copied from the main verb onto this auxiliary verb, and the tense features are deleted from the main verb. The meaning of the rule corresponds to logical negation. The rule derives sentence (35) from a structure such as  $S(\textit{he bought a book})$ :

(35) he did not buy a book

Finally, let us introduce a rule to form emphatic affirmative clauses. This rule, called *REMPH*, has the semantics of focusing, and it performs the following changes: it introduces the auxiliary verb *do* to the right of the subject, it copies the tense feature from the main verb onto this auxiliary verb and deletes these features from the main verb. The result, when applied to the structure  $S(\textit{he bought a book})$ , is sentence (36):

(36) he did buy a book

We have now postulated several rules, five of which we want to discuss further: *RYNQ*, *RWHQ*, *RTOP*, *RNEG* and *REMPH*. Each of these rules has its own meaning, and each performs its own formal operations. They could readily be incorporated into a compositional grammar. Nevertheless, the approach adopted here has serious shortcomings. Note that certain operations occur in each of these rules:

---

<sup>2</sup>Of course, I grossly oversimplify the description of the relevant facts here. In particular, I left rather vague what properties cause inversion in these constructions. This is a very complex issue which we will not go into here.

- introduce the auxiliary verb *do*
- copy the tense features from the main verb onto the auxiliary verb *do*
- delete the tense features from the main verb.

These operations must be described separately in each of these rules. It is clear that this is an undesirable state of affairs. In fact, the situation is worse than described, because the operations were only defined for a limited variety of sentences. As soon as we try to extend the analysis to cover an interesting fragment of English, the common operations invariably become considerably more complicated: different operations must apply if the input already contains auxiliary verbs, the introduction of the auxiliary verb is optional if the main verb is *have*, and then the main verb *have* starts to occupy the position where otherwise the auxiliary verb is introduced, etc., etc. The important point here is that all these complications are the same for all the rules introduced. By formulating these rules separately in each rule given above, this fact is described as being accidental; a considerable redundancy is introduced into the system; the rule system is unnecessarily large, and maintaining and updating the system of rules is made considerably more complex. In short, it is clear that a number of linguistic generalizations have been missed.

### 3.2.2 A Solution

It is obvious how the problems described in the previous section can be solved: the operations common to all these rules should be factored out. Each individual rule can then be considerably simplified, the grammar as a whole becomes smaller, the common properties are explicitly identified and isolated, and maintainability and updating become easier. A future change in an operation that has been factored out will immediately have effects in each of the constructions mentioned, clearly a desirable result.

One way to factor out these operations is by writing these operations as separate rules of grammar. This is the method which will be illustrated here, and has actually been used in the Rosetta system. This method has several advantages. It turns out that most of these operations resemble syntactic rules to a high degree. The method adopted accounts for this fact immediately: it supplies one with a format and a notation to describe the operations, and the modes of interaction with existing rules are also immediately determined. In addition, it makes it possible to keep the relevant operations local, as will be illustrated in section 3.5. An alternative method could consist of writing a set of functions which are called by rules, as suggested by [Partee, 1977], [Partee, 1979] and [Partee et al., 1990, 318-9]. Though I did not compare the methods in detail, the latter method requires additional clarification of the operations allowed, a format and a notation for them and a specification of the possible modes of interaction with rules, and it appears that certain operations cannot be kept local under this approach (see section 3.5). To my knowledge, this latter method has never been applied in a real large-scale

system, while the method adopted here has proven its usefulness in the Rosetta system.

There are, however, examples where it is advantageous to call functions from within rules, and where an analysis in terms of separate rules does not really achieve the desired result. These examples will be described in section 3.7. However, the functions called can perform only a limited set of well-defined operations which are independently required in the grammar.

Applying the solution, in which the common properties are factored out in separate rules, to the analysis of auxiliaries in English, we observe that the rule in which the common properties of the preceding rules are described performs operations on auxiliary verbs. It is not at all obvious that a meaning can be associated with this rule. The rule performs a number of the operations needed to form yes-no questions, wh-questions, topicalized constructions, negated sentences and emphatic sentences. Its semantics should be the ‘intersection’ of the meanings of these constructions, but this intersection is most probably empty. In addition, it is highly unlikely that translation-equivalent rules can be found in other languages. This rule must therefore be a *transformation*, i.e. a rule which is associated with the identity operation as its meaning. The rule will be called RAUX.

When we consider the rule RAUX in a realistic system, it proves desirable to split it up into two different rules. Note that RAUX introduces an auxiliary verb in all examples, but in certain cases it should do so to the right of the subject (in the constructions for which RNEG, REMPH have been proposed), and to the left of the subject in others (in the constructions for which RYNQ, RWHQ, RTOP have been proposed). When another auxiliary verb is present in the structure, no additional auxiliary verb should be introduced, but in the case of RYNQ, RWHQ and RTOP the auxiliary present has to change its position and be placed before the subject. By isolating and factoring out this operation of inverting the subject and the auxiliary verb, the complicated rule RAUX can now be broken down into two simpler rules: one which introduces the auxiliary verb *do* if no other auxiliary is present, and one which inverts the subject and the auxiliary in certain configurations. The rule of inversion is necessary anyway, to deal with cases of inversion involving auxiliaries other than the auxiliary verb *do*. This rule can now also be simplified, since it need no longer exclude the auxiliary verb *do*: it can simply state that in certain configurations any auxiliary verb must invert.

In fact, the rules and the whole grammar can be simplified still further if the auxiliary verb is introduced into *all* structures which do not already contain an auxiliary verb, and by postulating a rule which copies the tense features from the auxiliary *do* onto the main verb and deletes *do* if, at the end of the derivation, *do* and the main verb are adjacent. This simplifies the rule of auxiliary introduction, simplifies the rules assigning values to tense attributes (these rules can now refer in all cases to the auxiliary verb), and simplifies the rule of subject-verb agreement (which also need only refer to the auxiliary verb). The only problem it creates is that this solution will not work for REMPH (e.g. in *he did buy a book* the auxiliary verb *do* and the main verb are adjacent, but *do* cannot be deleted). This, however, can be solved in a very simple manner, e.g. by marking *do* in this construction as



[+stressed], and formulating the rule of auxiliary deletion in such a way that it does not apply to [+stressed] elements, or by postulating an abstract element EMPH which occupies the position between *do* and the main verb, and is deleted later in the derivation. This abstract element could be a basic expression and correspond to Dutch *wel* which must occur in translations of this English construction: *he did buy a book – hij kocht wél een boek*. The analysis of the English auxiliary system presented here is, of course, based to a large extent on the analysis given originally by [Chomsky, 1957] and in several studies since.

To summarize, the resulting analysis which has actually been implemented, can be described as follows. There is a transformation which introduces the auxiliary *do* to the right of the subject in all finite sentential structures, unless there is already an auxiliary or modal verb. The meaningful rules RNEG, REMPH, RYNQ, RWHQ and RTOP still apply, but they have been considerably simplified since the operations common to them have been factored out. There is an inversion transformation which inverts the auxiliary or modal verb and the subject in certain configurations created by the meaningful rules RTOP, RWHQ and RYNQ. And there is a transformation which copies the tense features of the auxiliary *do* onto the main verb and deletes this auxiliary if it is adjacent to a verb.

In the analysis, the relation between the meaning and the form of the sentence has now become very indirect. Though the sentences (37a,b,c) differ formally only with regard to the presence of an auxiliary and to the relative positions of the subject and the auxiliary, neither of these formal differences is accounted for by a rule which takes care of the semantic differences.

- (37) a he bought a book  
 b he did buy a book  
 c did he buy a book?

This situation is typical of most phenomena in natural language. The kind of approach described here, in which common properties of rules are maximally factored out and written as separate transformations, turns out to be useful for many constructions. This will be illustrated with a number of examples in the sections to follow. It will be shown that common properties have been factored out of semantically motivated rules, and that this leads to operations which operate in a purely formal, syntactic way, making the relation between meaning and form considerably more indirect.

### 3.3 Mood in Dutch

An account similar to the one given in the preceding section can also be given for mood<sup>3</sup> in Dutch. Consider the following Dutch sentences:

- (38) a Hij koopt een boek lit. ‘he buys a book’  
 b Koopt hij een boek? lit. ‘buys he a book’

---

<sup>3</sup>I use the term ‘mood’ here to distinguish different types of sentences, e.g. declaratives, interrogatives, imperatives, etc.

Here, as before, we have two different sentences: the sentences differ in form, and in meaning, and it would be straightforward to formulate rules to form each of the sentences and associate each with its own meaning. The formal differences between the two sentences relate to the position of the finite verb, and to the position of the subject. However, when one considers more relevant facts from Dutch, it appears that the simplest description for the facts concerning the placement of finite verbs and subjects in main clauses involves two rules: (1) a rule putting a finite verb into the position of a subordinate conjunction (provided none is present), and (2) a rule preposing the subject in a topic position. The first rule applies in a variety of constructions: in main clauses, in subordinate clauses, and even in certain kinds of adverbial clauses. It appears impossible to associate this rule with any semantic aspect. For the second rule it also appears to be implausible that semantics can be attached to it: one might think that this rule plays a role in topic-comment relationships, but in fact the rule also applies to subjects which are completely meaningless grammatical elements such as expletives, the Dutch variants of weather-*it* and extraposition-*it*, and idiom chunks, for which topic-comment relations appear irrelevant. And even for semantically non-empty phrases, the particular rule does not seem to change topic-comment relations at all: the position of the subject relative to other arguments and adjuncts in the clause does not change. In addition, it is not clear whether this rule can be associated with rules from other languages: there appears to be nothing that could correspond to this rule in English or Spanish, which also makes it difficult (or at least unnatural) to associate a meaning to this rule. Thus again we find a situation where a certain meaning difference is correlated to a formal difference, but where none of the rules taking care of the formal aspects is itself directly responsible for the meaning difference.

In the actual Dutch grammar in the Rosetta system, the relevant facts are accounted for in the following manner. There is a transformation which optionally puts phrases (among them: subjects) into a topic position. There is a set of meaningful rules, called *mood rules*, which assign mood to structures (i.e. declarative, yes-no question, imperative, etc.) and which read off the structure whether this is possible. Thus, the mood rule to form declarative main clauses requires the presence of a [-wh]-phrase as a topic. The mood rule to form wh-interrogatives requires the presence of a [+wh]-phrase as a topic, and the mood rule to form yes-no questions disallows any topic at all. Where appropriate, the mood rules introduce subordinate conjunctions. After the mood rules, transformations apply to put the finite verb in its appropriate position. These transformations apply blindly, and put the finite verb in the position of the subordinate conjunction if none is present (see chapter 5 for a more detailed treatment of these rules). So here again, the relation between the form and the meaning is very indirect. The whole grammar is considerably simplified by the approach adopted, and linguistically established generalizations are adequately expressed.

### 3.4 Order variants

Many sentences have variants in which the phrases have a slightly different relative order, sometimes accompanied by small additional formal differences (e.g. presence of a preposition). Some examples are given in (39-43).

- (39) a He gave the boy a book  
 b He gave a book to the boy
- (40) a Hij heeft waarschijnlijk die jongen gezien  
 He has probably that boy seen  
 ‘He has probably seen that boy’  
 b Hij heeft die jongen waarschijnlijk gezien  
 He has that boy probably seen  
 ‘He has probably seen that boy’
- (41) a Hij heeft naar het programma gekeken  
 He has at the program looked  
 ‘He has looked at the program’  
 b Hij heeft gekeken naar het programma  
 He has looked at the program  
 ‘He has looked at the program’
- (42) a Hij heeft aan de jongen dat boek gegeven  
 He has to the boy that book given  
 ‘He has given the boy that book’  
 b Hij heeft dat boek aan de jongen gegeven  
 He has that book to the boy given  
 ‘He has given that book to the boy’
- (43) a He bought a book for the boy  
 b He bought the boy a book

In (39a) the indirect object precedes the direct object, but in (39b) it is the other way around (and the indirect object is accompanied by the preposition *to*). In (40a) the direct object follows the sentential adverb, but in (40b) it precedes this adverb. In (41a) the prepositional complement precedes the verb, but in (41b) it follows. In (42a) the indirect-object phrase headed by the preposition *aan* precedes the direct object, but in (42b) it follows. And in (43a) the direct object precedes the beneficiary phrase (headed by the preposition *for*), but in (43b) it follows.

Semantic differences are frequently (but not always) associated with these different variants. These differences are sometimes rather subtle, but in other cases quite obvious. Thus, the sentences often differ with regard to topic-comment relations, and often there are scope differences. Of course, there have to be syntactic rules to describe the form differences between these sentences. Since there are

meaning differences between these variants as well, the issue arises whether such rules must be meaningful rules, and what the corresponding rules are in the other languages dealt with. If the rules to create these order variants are to be associated directly with meaning operations, a number of generalizations will be missed. Each of the relevant rules will have to be tripled, because in certain cases application of the rules does not correlate with any meaning difference and in other cases the meaning difference relates to scopal properties, while in yet other cases only topic-comment relations are changed. It is very difficult, if not impossible, to describe the relevant rules in such a way that the right meaning differences are correctly associated. In addition, the effects of such rules are global, i.e. moving a phrase into another position might change its scope relative to some other phrase which does not figure in the rule or its specific application at all. Finally, it is very difficult to set up rules in other languages which can correspond to these rules: they should have the same semantic effects, and apply in the same manner in the corresponding syntactic configurations.

This suggests that the rules responsible for the phenomena indicated should be transformations, and the semantic effects should be dealt with by rules which can take the global structure of the sentence into account. Phrases are therefore not introduced directly into a structure. Instead, *syntactic variables*, i.e. small syntactic trees consisting of exactly one node with a special attribute *index*, are initially introduced into the structure. Transformations to yield different orders then apply. And, finally, the semantic aspects, e.g. scope, are taken care of by *substitution rules*. These are rules which substitute a phrase for a variable with a specific index occurring in a syntactic tree. These substitution rules take global aspects of the sentence structure into account. For example, if a certain scope difference is expressed by linear order, they take linear order into account. This approach to scope, with variables and substitution rules, has been adopted directly from [Montague, 1974b]. The only difference with [Montague, 1974b] is that there are no syntactic transformations which can apply to syntactic variables in that analysis.

A similar approach can be adopted to account for the syntactic realization of topic-comment relations. Meaningful rules which mark a certain variable as being *topic* or *comment* might be assumed; several transformations can move variables into other positions, and the substitution rules which substitute the actual phrases for these variables should take syntactic conditions on topic-comment realization into account, e.g. be applicable only if the variable marked as comment occurs in a specific position. If this were combined with a pragmatic theory of which topic-comment relations in a sentence are appropriate in a specific context, we would have the beginnings of a full-fledged account of topic-comment relations.

In these examples we again see that certain semantic and pragmatic differences correspond to certain formal differences, in particular, differences in word order. But the rules which are responsible for the formal differences are not themselves directly responsible for these semantic or pragmatic differences.

### 3.5 Wh-movement

Consider the following sentences:

- (44) a He can see something  
 b He can see what?  
 c What can he see?

These three sentences differ semantically and pragmatically: (44a) is a simple declarative sentence, sentence (44b) is an echo-question, and sentence (44c) is a wh-question. Correspondingly, they also differ in certain formal respects. First, (44a) contains the word *something* as a direct object, and the other two examples contain the word *what* instead. Second, in (44b) the phrase *what* follows the verb, but in (44c) this phrase occupies a sentence-initial position. At first sight, an ideal situation for the compositional approach, but again, things will turn out to be slightly more complex.

Sentences (44a) and the other two sentences differ because they contain different basic expressions: *something* in (44a), *what* in the other two sentences. It remains to account for the difference between (44b) and (44c). The only difference between these two sentences is the position of the phrase *what*. So, it appears natural to postulate a rule, RMWH, which preposes the phrase *what* from the normal position of direct objects into a clause-initial position, and which accounts for the meaning of (44c). In addition, a rule is needed to assign a meaning to (44b).

The first problem for this simple approach, however, is that not only the fact that the phrase *what* is preposed to a clause-initial position is relevant, but also *how far* it is preposed. This is illustrated in (45):

- (45) a He will know what he can see  
 b What will he know he can see?

In both sentences of (45) the phrase *what* has been preposed to a clause-initial position, but they clearly have a different meaning. The only difference between the two sentences is the distance over which the phrase *what* has been preposed. So this aspect must be taken into account as well. Second, it has been argued by Chomsky ([Chomsky, 1977b]) that the rule preposing *what* in such sentences plays a role in many constructions, e.g. in relativization, in topicalization, clefting, pseudo-clefting, tough-constructions, comparative deletion, subdeletion in comparatives, equative constructions, complements to *too*, etc. In all these constructions an operation of preposing is required which obeys exactly the same restrictions.<sup>4</sup> So it is desirable to analyze all of these constructions as involving the rule RMWH as a suboperation. But now it has become quite difficult to assign a unique meaning to this rule. It cannot concern questioning, and in fact we should find a common semantic factor in all of these constructions and show that it has to be associated to this rule (and not to other operations). A third problem is created by languages which can form questions such as (44c) without applying a preposing rule. In such

---

<sup>4</sup>Where there are differences, these can be accounted for by independent factors.

languages the phrase corresponding to *what* can create wh-questions without being preposed. French is an example. The sentence *il peut savoir quoi?* can have the meaning of a normal question such as (44c). This would require that in French there is a rule corresponding to English RMWH which performs no formal change.

These arguments suggest that the rule RMWH should be a transformation, and that the difference between (44b) and (44c) should be dealt with in some other manner. In fact, as before, a substitution rule, in combination with the mood rule RWH proposed earlier, can achieve exactly the desired effect. The transformation RMWH preposes [+wh] variables to a clause-initial position. If there is more than one clause-initial position, RWH will, in principle, yield two results. A substitution rule substitutes a wh-phrase for the [+wh] variable which is clause-initial. And the mood rule RWH accounts for the fact that the clause which the interrogative phrase happens to introduce is turned into a wh-question. Thus, the combination of a substitution rule and the mood rule RWH will automatically control the application of the preposing rule RMWH, and account for the relevant semantic aspects, so that rule RMWH need not be associated with a meaning and can apply freely. This is a typical example of the situation mentioned in the introduction, in which transformations apply freely, after which their results are input to meaningful rules, which perform additional form changes and account for the semantic aspects. The approach described here is compatible with the semantic analysis of wh-complements as outlined by [Groenendijk and Stokhof, 1982], and therefore it makes it possible to integrate a well-motivated syntactic analysis of wh-movement with a well-motivated semantic analysis of wh-questions.

If we can also show that in all the other constructions mentioned above other rules than the rule RMWH are required and that the relevant semantics must be associated to these rules rather than to RMWH, we can fully integrate Chomsky's ([Chomsky, 1977b]) insight that the RMWH plays an important role in all these constructions. I think that this can indeed be done, and I have actually implemented such analyses for relativization and topicalization.

These facts show again that the relation between form and meaning is rather indirect. The sentences (44b,c) differ in meaning and in form, but the only formal difference which is directly visible (the position of *what*) is accounted for by a rule (RMWH) which itself is not directly responsible for the meaning difference.

The construction dealt with here also provides us with an argument in favor of the method in which operations which are factored out are written as separate rules, and against the method in which functions are called from within rules. The movement operation illustrated here can apply, in principle, over indefinite distances. By formulating the movement operation as a separate rule, the rule can be applied in a successive cyclic manner, so that each application is local. In the alternative approach, in which functions are called from within a rule, the function will be called in the rule RWH, and the movement must be able to apply over indefinite distances. Of course, all other things equal, local rule-applications are to be preferred to global applications, and in addition, there is empirical evidence that the rule involved here must apply in a successive cyclic manner (see chapter 5).

### 3.6 Generic Sentences

Consider the following sentences from Dutch:

- (46) a Hij zegt dat honden blaffen  
       ‘he says that dogs bark’  
       b Hij zegt dat er honden blaffen  
       He says that there dogs bark  
       ‘He says that dogs are barking’  
       c Blaffen honden?  
       Bark dogs?  
       ‘Do dogs bark?’  
       d Blaffen er honden?  
       Bark there dogs?  
       ‘Are dogs barking?’

The subordinate clause in (46a) must be interpreted as describing a generic<sup>5</sup> property of dogs, while the subordinate clause in (46b) must be interpreted as a statement pertaining to a specific situation in which some dogs happen to be barking. Similarly, (46c) is a question pertaining to a generic property of dogs, but (46d) asks whether — in the particular situation at hand — it happens to be the case that dogs are barking. Corresponding to these semantic differences, there are formal differences: the sentences (46b,d) contain the word *er* which does not occur in the other sentences. Again, we find a situation where a semantic difference appears to correlate fairly directly with a formal difference at first sight. It appears straightforward to write rules which account for these facts. However, as before, I will show that the situation is actually far more complex than one might think.

If one writes a rule to deal with sentences (46b,d) directly, then this would be a rule, say *RNONGEN* which introduces *er* and relates it to some NP, so that the NP is interpreted non-generically. And an additional rule would be required, say *RGEN*, which would account for the generic interpretation of this NP. An initial problem with this approach is that it is unclear how the use of *er* in this construction can be related to other uses of *er* in a principled manner. The pronoun *er* is also used in impersonal passives (47a), in passives of verbs with sentential complements (47b), in sentences where the subject must be interpreted non-specifically (47c), and in existential sentences (47d):

- (47) a Er werd gedanst  
       There was danced  
       ‘There was dancing going on’  
       b Er wordt beweerd dat hij ziek is  
       There is said that he ill is  
       ‘He is said to be ill’

---

<sup>5</sup>[Zwarts, 1988] contains a discussion of genericity in the framework assumed here.

- c Er kocht iemand een boek  
There bought someone a book  
'Someone was buying a book'
- d Er zijn geen boeken  
'There are no books'

A second problem is that one can show that absence or presence of *er* is not a crucial factor in all cases: in (48) the NP can be interpreted non-generically, but no *er* is present:

- (48) a Altijd lopen daar mannen  
Always walk there men  
'There are always men walking there'
- b Altijd spelen in het stadion voetballers  
Always play in the stadium soccer players  
'In the stadium soccer players always play'

The crucial factor here is a difference in position. And in fact, this holds for the earlier examples as well: the occurrences of the phrase *honden* occupy different positions in (46a,c) and (46b,d) as well. This is shown in (49):

- (49) a Hij zegt dat honden altijd blaffen  
'He says that dogs always bark'
- b hij zegt dat er altijd honden blaffen  
'He says that there always dogs bark'  
'He says that there are always dogs barking'
- c Blaffen honden altijd?  
Bark dogs always?  
'Do dogs always bark?'
- d Blaffen er altijd honden?  
'Are there always dogs barking?'

In (49a,c) the adverb *altijd* follows the phrase *honden*, while in (49b,d) it precedes this phrase. This order is obligatory:

- (50) a \*Hij zegt dat altijd honden blaffen  
He says that always dogs bark
- b \*Hij zegt dat er honden altijd blaffen  
He says that there dogs always bark
- c \*Blaffen altijd honden?  
Bark always dogs?
- d \*Blaffen er honden altijd?  
Bark there dogs always?

This suggests that we should consider introduction of *er* to be a purely syntactic transformation, and that the rule accounting for the generic or non-generic interpretation of NPs must be formulated in terms of the position the NP occupies.



This would make it possible to relate the use of *er* in these constructions to other uses of *er*, and it is now also possible to account for the examples where no *er* occurs.

There are, however, reasons for assuming that the generic v. non-generic interpretation should not be associated to the different positions of the phrase either, though this is more controversial. First, an account which is formulated purely in terms of positions would be insufficient in two ways: on the one hand there are many NPs that can never be interpreted generically, whatever position they occupy, and on the other hand there are NPs that can be interpreted generically or non-generically irrespective of the position they occupy (e.g. bare plurals in direct-object position: *he bought flowers* v. *he likes flowers*). So, there must be additional rules to specify which NPs can and which cannot be interpreted generically. Second, generic NPs cannot co-occur in just any type of sentence. There are strong restrictions on the aspectual and temporal properties of generic sentences. E.g. a generic NP can rarely occur in a sentence with the main verb in perfect tense (unless certain adverbs, e.g. *altijd* ‘always’ are present): *??Honden hebben gebiaft* ‘dogs have barked’. So there must be conditions on the combination with tense and aspect forms of verbs as well.

In Spanish, generic and non-generic NPs differ in form. The NP *honden* in its generic interpretation is translated into *los perros*, and this NP in its non-generic reading is translated into *perros*. This clearly shows that rules forming NPs must take genericity into account: there must be rules to generate the NP *perros* as an NP which cannot be interpreted generically, and the NP *los perros* as one which can. In English the distinction between the generic and the non-generic interpretation is formally encoded here by different tenses of the verb:

- (51) a He says that dogs bark  
 b He says that dogs are barking  
 c Do dogs bark?  
 d Are dogs barking?

In order to account for all these facts, I assume that the semantic difference between generic and non-generic NPs must be ascribed only to rules forming NPs. Bare plural noun phrases in Dutch and English are formed by two different sets of rules, one set yielding the generic NP, and the other yielding the non-generic NP. The restrictions on the occurrence of generic and non-generic NPs are accounted for by formulating appropriate conditions on the substitution rules that introduce such NPs into a clausal structure. The rules forming NPs account for the right form of the NP, and substitution rules check that generic and non-generic NPs are inserted into a structure in the right positions and only if the sentence has specific temporal and aspectual properties. In Spanish, *perros* and *los perros* are formed by different rules, and substitution rules check that they occur in the right environment.

This approach to these phenomena again leads to an indirect relation between form and meaning. Though initially it appeared necessary to ascribe the difference between generic and non-generic NPs to the presence of *er*, we soon concluded that

this was not the relevant factor, though this was the only directly visible formal difference between the initial examples. We shifted the semantics to rules dealing with the positions of NPs in a clause, but again we soon concluded that this would not be sufficient. And finally we decided to ascribe the semantic difference between a generic NP and non-generic NP to rules forming these NPs, leaving the syntactic and semantic conditions which hold with respect to the distribution of such NPs to substitution rules and their corresponding meaning operations. Such substitution rules are meaningful rules, but their semantics do not relate to genericity or non-genericity of NPs. The resulting situation implies that the semantic difference between (46a,c) and (46b,d) is not to be found in the absence v. presence of *er* (the only directly visible difference between these sentences), nor in the position that the NPs occupy (these positions are different, but this cannot be observed directly), but in a difference between the NPs occurring in these sentences (a difference that is only visible in abstract intermediate representations of the sentence). This is a clear example where the relation between form and meaning is very indirect: the rules that account directly for the visible form differences are not themselves associated with the meaning differences. Nevertheless, it appears to me to be the simplest analysis possible and the relation between form and meaning is clear in the grammar.

### 3.7 A Different Kind of Modularization

In the preceding sections we saw that in many cases it was useful to identify and isolate ('to factor out') common operations from several complex rules, and put these common operations in separate transformations. In this way each of the rules is simplified, and the relevant operations common to each of these rules need to be stated only once, which increases their maintainability and updatability considerably, and which expresses in a principled way that they are actually the same operations.

In this section, I will introduce a different kind of phenomenon, which resembles calling separate functions from within rules. These functions, however, are not arbitrary operations, but applications of normal rules of grammar, as we will see below. It turns out that the rules which are normally associated with a meaning must apply as a meaningless operation when called from within a rule. All cases involve *syncategorematic introduction* of elements.

The kind of modularization described here has not been implemented, but it is quite clear that this is a defect of the current grammars, and that modifying the grammars in the way described below would constitute a considerable improvement.

Let us start with a simple observation. There must be rules in the grammar to form NPs, and such rules must be able to create the proper syntactic structure to account for the syntactic behavior of a pronoun such as *it* when used in a clause. Similarly, there must be rules to form prepositional phrases, e.g. prepositional phrases expressing locations (e.g. *aan het water, by the water*), to express

directions (*e.g. door de tunnel, through the tunnel*), etc.

On the other hand, there are also many rules which introduce elements *syncategorematically*, i.e. they introduce lexical elements and their associated syntactic structure though these lexical elements or rules forming their associated syntactic structures are not represented in the derivation tree at all. Many examples can be given. The rule which forms sentential structures on the basis of verbs such as *rain, snow* etc. must introduce an NP headed by the lexical element *it*. The rules forming passives must introduce the preposition *door* in Dutch and *by* in English, combine them with one of the arguments of the verb, and construct a PP-node on top of them. The rules dealing with the syntactic realization of complements must sometimes introduce prepositions such as *aan, to* or *voor, for* and build the relevant prepositional structure on top of it, etc., etc. There is no practical problem in doing this. M-rules, the rules operating in controlled M-grammars, have the capacity to do this. But something is fundamentally wrong if we actually do it this way.

Note that the syntactic structures built by rules which introduce elements syncategorematically are identical to one another: whatever preposition is introduced syncategorematically in English, *by* or *to* or *for*, etc., the preposition always precedes its complement, and the PP node on top is created in the same way in all examples. If, however, the syncategorematically introduced structure is created separately in different rules, this correspondence is not expressed in any way. What is more, the form of the syntactic structures associated with elements that have been introduced syncategorematically are not only identical to one another, but they are also identical to the syntactic structures generated by the normal rules to form syntactic structures on top of such lexical items: the prepositions which are introduced syncategorematically precede their complements in English just as the prepositions which are not introduced syncategorematically. This fact must be expressed somehow. This is necessary in order to capture the generalization observed and express it in the grammar, and to increase the maintainability and updatability of the grammar. If each rule which introduces a preposition syncategorematically describes individually what additional syntactic structure is to be introduced, then a simple change in the rules forming prepositional phrases will require changes in all of these rules.

Before we outline the solution to these problems, it is first necessary to introduce the concept of a *model*. M-rules are rules which relate tuples of S-trees to S-trees. So, M-rules must contain descriptions of the S-trees they relate. Such descriptions are expressed by *S-tree models*. An S-tree model describes a class of S-trees. Put informally, we can say that an S-tree model is an S-tree containing variables. When an M-rule is applied, the input S-trees are matched with the S-tree models which are part of the M-rule. The variables which occur in the S-tree models are instantiated by this matching process. A separate S-tree model describes the class of output S-trees. When the variables are all instantiated, this S-tree model describes the form of the output S-tree exactly.

The relevant generalizations can be captured by extending the concept of model from S-tree models to *D-tree models*.

The first generalization to capture is the fact that structure which is introduced syncategorematically has the same form as structure which is introduced categorically. This is achieved by creating the former structure by applying the normal rules of grammar in the normal manner. The only difference with categorically introduced structure is that these rules apply *within a rule*. A D-tree model can be used to describe the class of relevant D-trees. After instantiation of the variables of the D-tree model, it is uniquely specified which rules must be applied and in which order they must be applied to derive the relevant structure. Thus, the rule which introduces the structure uses existing syntactic rules as specified in a D-tree model.

Second, each rule which introduces this structure syncategorematically contains a reference to this D-tree model. In this way, it is not necessary that each rule which introduces structure syncategorematically specifies the relevant D-tree model separately. The reference is possibly parameterized, to instantiate the variables occurring in the D-tree model. An example of this might be a D-tree model which specifies how a PP is to be constructed which has a variable for the lexical prepositional head of the construction. See below for more detailed examples.

Application of a rule which introduces some element syncategorematically must now be interpreted as follows: for all normal elements occurring in a rule, the normal interpretation holds, but as soon as a reference to a D-tree model is found, this D-tree model must be located and the variables it contains must be instantiated. The result is a D-tree which is converted into an S-tree by applying the rules of the D-tree in the normal manner. The resulting S-tree is inserted in the S-tree at the point where the reference to the D-tree model was found.

This method complicates the rule notation and its interpretation, but it solves the problems mentioned above associated with syncategorematic introduction. Small changes in individual rules playing a role in the formation of syntactic structures which can arise as a consequence of syncategorematic introduction are taken into account immediately: when the normal rule changes, the structures resulting from syncategorematic introduction change as well, because the same rule is used. Larger changes, which affect the rule interaction (e.g. rules are removed, new rules are added, or the application order of rules is changed) require only one change in the relevant D-tree model for all cases.

The method will be illustrated with a few examples. Let us start with a simple example where the D-tree model contains no variables. There are many rules in the Rosetta grammars which must syncategorematically introduce an NP headed by the pronoun *it*. Examples are the start rules for impersonal verbs such *rain*, *hail*, the start rules for impersonal adjectives (e.g. *it is cold*), several rules which account for the complement structure of verbs and rules which fill the subject position with expletive pronouns if it is empty. Each of these rules has a reference, to a D-tree model (let's call it *ItNP*) to form such an NP. In this way it is guaranteed that all NPs headed by *it* which are introduced syncategorematically are formed in the same way.

The D-tree model *ItNP* itself describes a class of D-trees to form such an NP. Since *ItNP* contains no variables, it describes exactly one such D-tree. The D-tree

model contains the names of the normal rules which are used to form such NPs. In this way, it is guaranteed that all NPs headed by *it* are formed in the same way, independent of whether they are introduced syncategorematically or not. The D-tree model to form such an NP might be a very simple model (of course, the exact form depends on the grammar used), consisting of a top node labeled with a rule name (e.g. *RNPformation3*) which dominates a node labeled with the name of the basic expression *it*:

$$(52) \quad \begin{array}{c} \text{RNPformation3} \\ | \\ \text{IT} \end{array}$$

The rule *RNPformation3* puts an NP-node on top of the S-tree for the basic expression *it*, it assigns a relation *head*, and it accounts for the setting of all kinds of attributes of the NP-node, e.g. to indicate that the NP is singular, definite, specific, headed by the pronoun *it* (cf. the attribute NFORM in HPSG, [Pollard and Sag, 1987, 62]), 3rd person, neuter gender and certain other properties.

A second example, which deals with the syncategorematic introduction of prepositional phrases, is slightly more complicated. All rules which must introduce a prepositional phrase syncategorematically (e.g. the rule which introduces *by*-phrases in passives) contain a reference to a D-tree model *MakePP* which contains two variables. The first is a variable for the name of the preposition which will serve as the head of the PP. The second is a variable for an S-tree which is to serve as the complement of the preposition.

The D-tree model *MakePP* might take the following form:

$$(53) \quad \begin{array}{c} \text{RPP} \\ | \\ \text{TassignCase} \\ | \\ \text{RstartPP} \\ \wedge \\ \text{T1} \quad \text{T2} \end{array}$$

In this D-tree model, the variable *T1* will be instantiated by a PREP before the resulting extended D-tree is evaluated, and the variable *T2* is instantiated by the S-tree which serves as the complement of the preposition.

When variable *T1* is instantiated by the preposition *by* and variable *T2* by the S-tree *NP[head/BN{name:JOHN}]*, the D-tree can be evaluated by applying the rules it contains. The rule *RstartPP* combines the preposition *by* with the instantiation of *T2* as its complement, in accordance with the lexical specifications of the preposition *by*. Since the rule *RstartPP* is the rule which is normally used in the grammar to combine a preposition with its complements, the preposition will precede its complement (in English). Thus, it is now adequately expressed that prepositions in English precede their complements, irrespective of how they have been created (by syncategorematic or by categorematic introduction). The

transformation *TassignCase* assigns *Case* to the complement of the preposition, so that the complements of syncategorematically introduced prepositions receive case in the same way and by the same rule as complements of prepositions which are introduced categorically.

The rule RPP puts the top node labeled PP on the structure and assigns specific values to attributes of PP, e.g. whether the PP is a wh-PP, which preposition heads the PP, etc.

This second example is more complex than the first example. The concept of D-tree model, in which variables can occur, plays a crucial role here.

The method can be seen as a special instance of the method where common properties in rules are factored out and put into special functions. But it does not suffer from the defects of the method mentioned earlier: the functions are not just arbitrary functions but have only limited power: they can only derive structures by the normal rules of grammar.

The method described here is actually a special (simpler) case of the method described by Schenk ([Schenk, 1989]) to describe idiomatic expressions in the controlled M-grammar formalism.

Note that the method requires that rules which are normally associated with a meaning operation apply here without this associated meaning operation. This is true for all these rules irrespective of the number of arguments they take. As a consequence, we find examples of dyadic rules here which are not associated to any meaning. See chapter 2.

It is quite clear that this method is to be preferred to other methods in which the syncategematic elements, and their associated structures, are introduced by separate statements in different rules which are unrelated to the normal rules which create such structures.

### 3.8 Concluding Remarks

In this chapter I have shown that the controlled M-Grammar formalism makes it possible to find a proper balance between purely syntactic requirements and the requirements of a compositional grammar in which form and meaning must be associated in a direct way. The controlled M-Grammar formalism is a compositional formalism and has for this reason a strong semantic bias. Though such a strong semantic bias could easily lead to syntactically inadequate analyses, there is compensation by allowing syntactic transformations, rules which have the identity operation as their meaning. In addition, we found that it was necessary to allow that rules which are normally associated with a meaning are applied without this associated meaning to form structures which are introduced syncategorematically.

In this way it is possible to express syntactic generalizations, to simplify individual rules of the grammar, to reduce the overall size of the grammar, and to avoid redundancy, in short, to approach a more modular organization of the grammar, so that it becomes easier to maintain and update.

The combination of meaningful rules and syntactic transformations, each of

which may be relatively simple, increases the amount of rule interaction and makes the relation between form and meaning rather indirect in the following sense: many rules which account for form differences should be transformations in a syntactically adequate description, even though the form differences appear to correlate with meaning differences. Other rules, which often apply to more abstract representations of the sentence and therefore do not have directly visible formal effects in the surface, account for the meaning differences.

The facts discussed show that the relation between form and meaning is quite subtle and that there is no direct correspondence between meaning aspects and surface properties of utterances. Note that this does not have negative consequences for the relation between form and meaning in the grammar. In the grammar this relation is clear, which is a main virtue of compositional grammars, including the controlled M-grammar framework.





## Chapter 4

# Predicate-Argument Relations

### 4.1 Introduction

In this chapter, it will be described how the relation between predicates, mainly verbs, and their arguments is accounted for in the Rosetta grammars. A number of characteristics of this relation will be given, and it will become clear that many of these follow from the compositional character of the grammars.

First, I will make some general remarks about the relation between syntax and semantics with regard to predicate-argument relations (section 4.2). Next, I will discuss the treatment of arguments and formulate three of its characteristics (section 4.3). Arguments need not have an overt realization in surface trees. Such arguments are called *covert arguments*. Section 4.4 describes how different kinds of covert arguments are dealt with. In section 4.5 a distinction between *bound adjuncts* and *adverbials* (or *free adjuncts*) is introduced and explained. In section 4.6 the status of various kinds of *small clauses* in the Rosetta grammars is discussed. Section 4.7 describes the necessity of expressing systematic relations between different, but related lexical items, especially when they differ in the number of arguments they take. In section 4.8 some conclusions and topics for further research are summarized.

### 4.2 Predicate-Argument Relations and Compositionality

The grammars used in the Rosetta machine translation system are compositional. It is natural that this compositional approach, in which there must be an intimate relation between basic expressions and basic meanings, and between meaningful syntactic rules and meaning operations, has determined the character of the treat-

ment of the relation between verbs and their arguments in the Rosetta syntax to a high degree.

The strong bond between syntax and semantics with respect to argument structure also ties in nicely with certain syntactic theories. Thus, [Chomsky, 1981, 36] introduces the  *$\theta$ -Criterion*, which requires a very strict relation between a lexical specification of the arguments which a verb (or a word of a different category) takes in terms of thematic roles and the realization of arguments at syntactic levels of representation.

Though the  $\theta$ -Criterion played an important role in theories in the Principles and Parameters (P&P) framework since the appearance of [Chomsky, 1981], recently [Chomsky, 1991, 45-46] claims that the  $\theta$ -criterion is not required as an independent principle. The fact that it holds is supposed to be a direct consequence of the fact that each element at an interface level must receive an interpretation in the components that the level interfaces. In the controlled M-grammar framework, this is guaranteed by the compositional nature of the grammars, in a very strict way: syntactic derivation trees and semantic derivation trees are isomorphic, so any element appearing in one of them must have a unique corresponding element in the other.

In my view, the elimination of the  $\theta$ -Criterion as an independent principle in the P&P framework is desirable anyway, because it is unclear how this criterion can be formulated elegantly (if at all) directly on syntactic structures which resemble surface trees: one has to take into account that all kinds of ‘meaningless’ material can intervene between predicate and argument, e.g. meaningless prepositions and virtually any material in the case of idiomatic expressions; and one must take into account that the predicates themselves may not be easily identifiable: in the case of idioms they can take a wide variety of forms, be discontinuous, etc. None of these problems hold for syntactic derivation trees in the controlled M-grammar formalism: here the relation between predicates and arguments is always expressed as a relation of sisterhood between simple elements in a D-tree.

The *Projection Principle* has played an even greater role than the  $\theta$ -criterion in theories in the P&P framework since 1981. This principle states that the  $\theta$ -Criterion does not only hold at the interface level to semantics (LF in the P&P framework), but also at the levels D-Structure and S-Structure. Within the controlled M-grammar framework no such levels exist, so the Projection Principle is simply irrelevant. And recently, [Chomsky, 1991] initiated a program to get rid of these levels of representation. If this program succeeds, the Projection Principle becomes meaningless.

One might wonder whether the  $\theta$ -Criterion (or some variant of it) holds for surface trees in the controlled M-grammar framework: the answer is no. I think that it is impossible to state the  $\theta$ -Criterion elegantly in terms of surface trees, for the reasons mentioned above. In addition, it is almost impossible to maintain the  $\theta$ -Criterion if no abstract elements are allowed. Since surface trees cannot contain lexical S-trees which correspond to the null string, the use of abstract elements in surface trees is highly limited.

From this discussion it is clear that the compositional nature of the controlled

M-grammar has immediate and far-reaching consequences for the treatment of predicate-argument relations: the very design of the grammar immediately derives a variant of the  $\theta$ -criterion in a very strict way. In the sections to follow we will see more consequences of the compositional character of the grammars for the treatment of predicate-argument relations.

### 4.3 Arguments

I assume that the parts that make up a sentence<sup>1</sup> can be partitioned into a number of classes. For predicate-argument relations the following classes are relevant: (1) the predicate; (2) the arguments; (3) the adjuncts.

In this section we will deal mainly with arguments. Section 4.5 will discuss adjuncts in more detail.

It is not always easy to distinguish arguments from non-arguments, but I will assume that this is clear for a core set of cases. A number of criteria can be derived from the characterization of the notion argument in this chapter. A number of global criteria have been used in practice. These criteria relate to obligatoriness, meaning, selection restrictions and form restrictions. [Pollard and Sag, 1987, 135-137] mention criteria such as order-dependence of content, constancy of semantic contribution, iterability, relative ordering and the possibility of internal gaps as possible criteria to distinguish arguments from free adjuncts.

*Arguments* have an intimate connection with and are fully dependent on the predicate word: both their presence and their form is determined by the predicate word to a high degree. As far as translation is concerned, the syntactic realization of the translation of the argument (i.e. what grammatical relation it bears, whether other words should be present, etc.) is fully dependent on the properties of the translation of the predicate word. Typical examples of *arguments* are subjects, direct objects, indirect objects and prepositional objects.

*Adjuncts* have a much looser syntactic connection to the predicate word with which they co-occur. They are always optional, generally can appear with all predicate words (though sometimes there are semantic restrictions), their form is not determined by the predicate, and in principle they can be iterated (subject to semantic constraints). Typical examples of *adjuncts* are temporal and locative expressions.<sup>2</sup>

As will become clear in this chapter, there may be arguments which have no overt realization in surface trees. Such arguments are called *covert arguments*. They will be discussed in section 4.4.

If a phrase is an argument, it is an argument of some expression (often, but not necessarily, a single word). This expression will be called a *predicate*. As pointed out in chapter 2, there is a special class of rules, called *start rules*, which combine a predicate with its arguments. Such a rule combines a predicate that takes  $n$

<sup>1</sup>These 'parts' are not necessarily visible in surface trees or in the surface string.

<sup>2</sup>Though locative expressions can sometimes be arguments, e.g. to a verb such as *wonen* 'to live'.

arguments (an  $n$ -ary predicate) with  $n$  arguments. The arguments of a predicate are syntactic variables in the current grammar. To illustrate, the instantiation of such a start rule to combine the predicate *give* with its arguments to form the sentence *he gave the book to John* could be:<sup>3</sup>

$$(54) \quad R(\text{give}, x_1, x_2, x_3)$$

where substitution rules will substitute *he* for  $x_1$ , *the book* for  $x_2$  and *John* for  $x_3$  at a later point in the derivation. Given the compositional nature of the Rosetta grammars, there must be a meaning operation associated to a start rule, which applies to the meaning of the arguments of the rule. It follows that the arguments of the rule must have a meaning. This immediately leads to the first characteristics of arguments in Rosetta:

**Semantic-Arguments Characteristic** An *argument* must have a meaning.

This property need not be stipulated independently, but is a direct consequence of the compositional nature of the grammars. It is nevertheless important to emphasize it, because the notion *argument* is often understood in a purely syntactic way.

I will therefore make an explicit distinction between *semantic arguments* (i.e. arguments which have meaning) and *syntactic arguments* (i.e. phrases which occur in positions where semantic arguments can occur). The notion of *semantic argument* is the one that plays an important role in the grammars. The concept of *syntactic argument* actually plays no role in the grammars, but it is convenient to have a term for this concept.

I will sometimes use the word *argument* on its own, without any modification. In that case I invariably mean *semantic argument*.

To illustrate the difference between semantic and syntactic arguments, the following examples can be used: the boldface words or word groups in the following examples are **NOT** considered to be *semantic arguments* in Rosetta:

---

<sup>3</sup>The start rule given combines all arguments with a predicate at one time. One might also combine a predicate with its arguments by a number of rule applications, adding one argument at a time. In the actual grammar, the start rules combine the predicate with all its arguments at one time. The main reason for doing this is to be found in simplifying the tuning of grammars of different languages. This method is feasible because the arguments that a start rule combines with their predicate are variables.

- (55) a **It** is raining  
 b **There** is a man in the room  
 c **Er** kwam een man aan ('There arrived a man')  
 d Ik schaam **me** er niet voor ('I am not ashamed of it')  
 e Ik scheer **me** iedere dag ('I shave every day')  
 f **It** is said that he is ill  
 g Hij betreurt **het** dat hij ziek is ('He regrets it that he is ill')  
 h Hij rekent **erop** dat hij mag komen ('He counts on it that he can come')  
 i He is looking **at the girl**  
 j She was looked at **by him**  
 k She gave a bottle of wine **to her friend**  
 l He kicked **the bucket** (idiomatic reading)  
 m He considered her **interesting**

Weather-*it* (55a), existential *there* (55b), expletive *er* (55c), inherently reflexive pronouns (both so-called *necessarily* (55d) and *accidentally* (55e) inherent reflexives), extraposition-*it* (55f,g), or variants of it (*er* in the Dutch example (55h)), prepositional phrases which contain a prepositional object (55i) or which express the *by*-phrase (55j) or the indirect object (55k), idiom chunks (55l) and parts of small clauses (55m) are not arguments in the Rosetta grammars. In the examples (55i,j,k) only the NP contained in the PP is an argument, but the PP as a whole is not. In example (55m) the combination *her interesting* is an argument, but the part *interesting* on its own is not.

Start rules combine predicates with their arguments to form sentences, but they do so under abstraction from many generally applicable syntactic rules such as extraposition, heavy NP-shift, topicalization, focus movement, wh-movement, clitic rules, relativization rules, etc; so-called NP-movement, hence passivization, subject-to-subject raising, subject-to-object raising, (object) NP-movement in the case of ergative verbs and passive verbs, cases of external passivization (*he was taken advantage of*), etc.; verb movements such as Verb Second, Verb Raising, inversion, etc. This makes it possible to realize predicate-argument relations in local configurations in S-trees.

The second important property with respect to arguments in Rosetta, which might be called the *Fixed-Arity Condition* can be formulated as follows:

**Fixed-Arity Condition** Every predicate has a fixed arity.

This property is important for the treatment of phrases which are optionally present. It is not possible to say in the system proposed here that e.g. a verb such as *to eat* takes one or two arguments (cf. *he is eating* v. *he is eating a sandwich*), for that would violate the *Fixed-Arity Condition*. In such a case there are two options: (1) either the relevant predicate is analyzed ambiguously, e.g. as *eat<sub>1</sub>* and *eat<sub>2</sub>*, and each of these words individually satisfies the *Fixed-Arity Condition*, or (2) the relevant word is not analyzed ambiguously. In the latter case its arity is taken to be the maximum of the alternatives (in the case of *to eat*: 2). It is specified that one of the arguments of this word can be realized as

a special phrase, called *EMPTY*. This special phrase is deleted somewhere in the derivation, so that no overt argument appears in the surface tree. See section 4.4 for further discussion.

The Fixed-Arity Condition is necessary to ensure that the relation between syntactic derivation trees and semantic derivation trees can be defined completely in terms of relations between nodes of these trees, and finds its basis in the logical foundations of the compositional approach. It also is useful in that it provides one with a strict guideline to distinguish different uses or meanings of predicates, so that it can be used as a heuristic for this purpose. However, it also has a number of drawbacks for certain phenomena. These will be discussed in section 4.5, along with a remedy to remove these drawbacks.

A third important property of the concept of argument in Rosetta relates to translation and can be described as follows:

**Arity-Preservation Condition** The arity of a predicate and its translation must be the same.

This is a very important property, that determines the structuring of the lexicon to a high degree. The property is a direct consequence of the isomorphic approach. It is important to emphasize this property because it does not hold in other frameworks, e.g. in the EUROTRA framework as described by [Arnold et al., 1986, 300]. In this framework two identical co-indexed arguments can be rendered as one argument in a different language.

In the remainder of this section I will discuss a number of additional properties of arguments. In section 4.3.1 a convention for ordering arguments in D-trees will be introduced. In section 4.3.2 the distinction between internal and external arguments will be described and in section 4.3.3 the attribute-value pairs to characterize properties of predicates are given.

### 4.3.1 Argument-Ordering Convention

Start rules combine an expression with its arguments in an S-tree. In this S-tree constituent structure and grammatical relations (such as subject, object, predicate, etc.) are expressed.

Consider the following example sentences:

- (56) a he gave the book to John  
       b he gave John the book

A grammar must express in a principled manner that sentences (56a) and (56b) are synonymous. One necessary condition to achieve this is to be able to associate the corresponding arguments in these sentences, i.e. it must somehow be expressed that the phrase *the book* plays the same role in both sentences. One way of doing this is by labeling these arguments with the same label, e.g. a semantic role such as *theme*. In Rosetta, this has been achieved by adopting a convention to put the

arguments in a specific order in the start rule. The correct instantiations of the start rules for the two sentences are given in (57):<sup>4</sup>

- (57) a R(gave, he, the book, John)  
 b R(gave, he, the book, John)

In order to guarantee that this is being done systematically, a special convention has been adopted, called the *Argument-Ordering Convention*. It orders the arguments in D-trees on the basis of the grammatical relations they bear. This convention is a guideline which has been adopted, but which can be overruled if required for the purposes of translation. Such cases will not be dealt with here. See [Appelo, 1993] for discussion.

**Argument-Ordering Convention** subject  $\prec$  direct object  $\prec$  sentential complement<sup>5</sup>  $\prec$  prepositional object, predicates, locatives, directionals  $\prec$  indirect object, *aan*-PP (*a*-PP), *voor*-PP (*for*-PP, *para*-PP)

Notice that a distinction is made between direct objects and sentential complements. This is necessary, because direct objects and sentential complements can co-occur with certain verbs (e.g. *dwingen* ‘to force’ in Dutch). I will illustrate the argument-ordering convention with a number of examples:<sup>6</sup>

- (58) 

|                                                                       |
|-----------------------------------------------------------------------|
| He <sub>1</sub> gave Peter <sub>3</sub> a book <sub>2</sub>           |
| He <sub>1</sub> gave a book <sub>2</sub> to Peter <sub>3</sub>        |
| It <sub>1</sub> irritated him <sub>2</sub>                            |
| I <sub>1</sub> forced him <sub>2</sub> [to do this] <sub>3</sub>      |
| I <sub>1</sub> promised him <sub>3</sub> [to do this] <sub>2</sub>    |
| I <sub>1</sub> ordered him <sub>3</sub> [to do this] <sub>2</sub>     |
| I <sub>1</sub> accused her <sub>2</sub> of a crime <sub>3</sub>       |
| I <sub>1</sub> painted [the door] <sub>2</sub> green <sub>3</sub>     |
| I <sub>1</sub> sent him <sub>2</sub> [to Amsterdam] <sub>3</sub>      |
| This <sub>1</sub> costs us <sub>3</sub> [three guilders] <sub>2</sub> |
| omdat mij <sub>2</sub> [dat mooi] <sub>1</sub> leek                   |
| Dat <sub>1</sub> viel hem <sub>2</sub> op                             |

Notice that the order of the argument need not correspond to the left-right order in S-trees. Thus indirect objects obligatorily precede direct objects and the NP in *aan*-PPs can follow direct objects in Dutch, but in both cases they are the third argument for a verb such as *geven* ‘to give’.

<sup>4</sup>I abstract here from the fact that start rules actually combine variables with a predicate.

<sup>5</sup>When a sentential complement is introduced by a preposition, it counts as a prepositional object.

<sup>6</sup>The latter two examples are Dutch, because English has no equivalents. The gloss of the penultimate is *because me that beautiful seemed* (‘because that seemed beautiful to me’); the gloss of the last example is *that fell him up* (‘that struck him’).

### 4.3.2 External and Internal Arguments

Start rules combine a predicate with its arguments in an S-tree. The predicate can specify whether its first argument will bear the grammatical relation subject or not. The argument that is made the subject of a structure when it is combined with a predicate and other arguments is called the *external argument*. All other arguments are *internal arguments*. Recall that a predicate and its arguments are combined at a point in the derivation where generally applicable syntactic rules are abstracted away (see above). After the application of start rules which turn arguments into external or internal arguments, pattern rules check additional properties of internal arguments and syncategorematically add the grammatical words required. This can be illustrated by the following examples:

- The verb *see* takes two arguments. The first argument is made the subject by start rules, and the second one is made a VP-internal argument. Pattern rules determine later that the internal argument bears the grammatical relation *object*. In a sentence such as *John saw the book*, *John* is the external argument, because it is made the subject of the sentence by the start rules.
- The Dutch verb *komen* ‘to come’ takes one argument. This argument is not an external argument. It is made the VP-internal argument when it is combined with the verb, and pattern rules determine later that it bears the grammatical relation *object*. Such verbs, which take no external argument, but do take (at least one) internal argument of the category NP, are called *ergative* or *unaccusative* verbs. Though the argument of the verb *komen* starts out as a direct object, it very often does not end up as a direct object, but as a subject in the surface tree. There are special rules, which are also used for passive structures, to turn these direct objects into subjects.
- The verb *seem* can be used in many ways. We will consider here its usage as illustrated in the sentence *it seems to me that he is ill*. The verb *seem* takes two arguments. Neither of these arguments is made the subject by start rules. Both arguments are VP-internal arguments. Pattern rules determine later in the derivation that the first argument (*that he is ill*) functions as the complement sentence, and that the second argument (*me*) functions as a prepositional object governed by the preposition *to*. There is no external argument. The phrase *it* is not an argument, but a semantically empty syntactic filler (an expletive pronoun), which must occur in English whenever no other subject occurs.

For each individual expression it is stipulated whether it takes an external argument, or not. The choice between these two options is an empirical matter, i.e. one should choose that option which yields the simplest description of the relevant facts. There is a lot of literature arguing for the distinction in general, or arguing in favor of or against the status of a class of expressions as taking external arguments (e.g. [Perlmutter, 1978], [Burzio, 1981], Den Besten ([den Besten, 1982]), [Hoekstra, 1984], Den Besten ([den Besten, 1985]),



[Belletti and Rizzi, 1988], [Cinque, 1990] and many others). The terminology (*external argument*) has been taken from [Williams, 1981]. In the next section I will describe how it is indicated in Rosetta whether an expression takes an external argument, or not, and I will give some more illustrations of the distinction.

No categorial restrictions or other syntactic restrictions have been imposed on the external argument in Rosetta, though the formalism does not prevent one from doing so.<sup>7</sup> Maybe, an external argument cannot be just of any category, but if so, this is not determined by individual expressions, but holds for the whole class of external arguments. Many languages (e.g. English, Dutch) have a restriction on the possible syntactic category of the subject in surface trees: the surface subject must be an NP. For this reason, special rules are present to guarantee that external arguments of categories other than NP do not remain subjects (e.g. extraposition rules). It may appear that there are restrictions on the possible syntactic categories of subjects (e.g. a sentential subject for the verb *eat* is ill-formed: *\*that he bought the book is eating a sandwich*), but these can always be better accounted for by semantic-type restrictions on the subject. A verb such as *eat* requires animate subjects, and this does not only exclude all sentential subjects, but also NPs with an inanimate meaning (e.g. *\*his buying of the book is eating a sandwich*).

[Pollard and Sag, 1987, 129-131] give three possible arguments in favor of the necessity to specify form restrictions on the subject. The first argument relates to the option of choosing expletive *it*, expletive *there* or a ‘normal’ NP as a subject. As far as I can see, the choice between the expletive *there* and the expletive *it* can be made on the basis of a general rule. Therefore, it is not necessary and in fact even undesirable to specify this information with each relevant lexical item. This reduces the problem to a choice between an expletive or a non-expletive: but this choice correlates perfectly with the argument v. non-argument status of the subject, and can therefore also be derived by general rule.

The second argument relates to the fact that certain predicates in Icelandic select subjects with a specific case. We did not consider this type of fact when designing the grammars, since Icelandic was not one of the languages we dealt with. But it might very well be the case that this constitutes a real example of verbs imposing form restrictions (relating to case) on subjects. Of course, nothing in the formalism prevents one from specifying form restrictions on the subject, but no evidence to do so was found in most other languages.

The third argument relates to a number of predicates (e.g. *make*, *result*, *be incoherent*) in English which apparently do not allow a finite clause, but do allow an NP as their subject. Pollard and Sag themselves note that their conclusions are somewhat controversial. The analysis they suggest is incompatible with the analysis of sentential subjects as given by [Koster, 1978b] (see also [Emonds, 1976, 127,130]). Koster argues that sentences are never subjects, so that it is impossible in this analysis to capture the relevant facts by imposing form restrictions on the

---

<sup>7</sup>Semantic restrictions, e.g. semantic-type restrictions, can be imposed by a predicate on each of its arguments. Such restrictions are imposed in the semantic component, but this has not been implemented yet. See [Grimshaw, 1979], who motivates a similar distinction between syntactic subcategorization and semantic selection.

subject. This suggests that the relevant facts should be accounted for in some other way, perhaps a semantic account, a possibility not fully excluded by Pollard and Sag, though, as they note, such an analysis still has to be developed.

For the non-external arguments (the *internal arguments*) strong syntactic restrictions can be specified, on their categorial status, which grammatical relation they bear, whether they must be accompanied by specific grammatical words (e.g. specific prepositions, specific subordinate conjunctions or specific infinitival markers such as e.g. *to* in English), which case they must bear, etc., i.e. typical strict subcategorization information.

### 4.3.3 Attribute-Value Pairs to Specify Arguments

In this section I will describe in more detail how the distinctions introduced in the preceding section are represented in the Rosetta grammars. A special attribute, called *thetavp* is used to indicate the number of arguments an expression takes and whether there is an *external argument* or not, and to indicate in which order the other arguments must be put in the structure. The maximum number of arguments allowed is 3.

The second attribute, the *pattern* attribute, specifies the following information for each *internal argument*:

- what its syntactic category must be
- what grammatical relation it must bear
- whether additional words should be present (e.g. specific prepositions)

There is a unique pattern name for each possible realization of the internal arguments, and the actual realization is performed by pattern transformations. Start rules make an S-tree using the attribute *thetavp* in which the external argument (if there is one) is made the subject of the structure, and the internal arguments (if any) are put inside VP with no further information except that they are arguments. Pattern transformations take this structure as input and change it in accordance with the value of the *pattern* attribute.

I will now describe the relevant attributes, and their possible values in more detail. Attention is limited to monolingual issues:

**thetavp** The possible values of this attribute have the general format *vp###*, where each # is replaced by a different value from the set {1,2,3} or by 0.

It is assumed that there are three positions to realize an argument: the subject position (for the external argument), and two positions inside VP (for the internal arguments). Each # represents such a position. The first # is for the subject position, the second # is for the leftmost position inside VP, the third # is for the rightmost position in the VP.

It is furthermore assumed that arguments (maximally three) come in a certain order. Let us represent arguments by the letters *x,y* and *z*.

A 0 replacing # means that *no* argument is realized in the corresponding position. A 1 replacing # means that the *first* argument is realized in the corresponding position. A 2 replacing # means that the *second* argument is realized in the corresponding position. A 3 replacing # means that the *third* argument is realized in the corresponding position.

I will illustrate this with some examples:

- Suppose a verb with *thetavp* = *vp120* is combined with two arguments, *x* and *y*, in that order. Then the first argument (*x*) is realized as a subject (external argument), because the digit 1 is the first digit in *vp120*. The second argument (*y*) is realized as the leftmost and only argument inside VP. This is a typical value for transitive verbs (e.g. *see*).
- Suppose a verb with *thetavp* = *vp010* is combined with one argument, *x*, then this argument is realized as the leftmost and only argument inside VP, because the digit 1 is the second digit in *vp010*. There are no other arguments, so there is no external argument. This is a typical value for monadic *ergative verbs* (e.g. Dutch *komen* ‘to come’).
- Suppose a verb with *thetavp* = *vp012* is combined with two arguments, *x* and *y*. Then, both arguments are realized inside VP. This is a typical value for dyadic ergative verbs (e.g. Dutch *opvallen* ‘to strike’) and for verbs such as *seem* in English.

The *Arity-Preservation Condition* only requires that the number of arguments is preserved under translation. It does not say anything about the realization of the arguments. Hence, words with *thetavp* = *vp010* can be synonymous with and translated into words with *thetavp* = *vp100* (e.g. ergative verbs into intransitive verbs), words with *thetavp* = *vp120* can be translated into words with *thetavp* = *vp012*, etc.

From a monolingual point of view all values for *thetavp* can have the digits in ascending order (except for 0), i.e. from a purely monolingual point of view there is no need for values such as *vp210*, *vp132*, etc. However, this is necessary to adequately deal with translation. In certain cases, the conventions for ordering the arguments (see below) yield different results in different languages. In that case, in one of the languages the value for *thetavp* must be changed into a value with descending digits. I will not discuss this further, but see [Appelo, 1993].

The possible values for this attribute are:

**vp000** Verbs not taking any arguments, e.g. *rain*, *snow*, *hail*.

**vp010** Ergative verbs taking one argument, e.g. Dutch *vallen* ‘to fall’, *aankomen* ‘to arrive’, *sterven* ‘to die’;<sup>8</sup> subject-raising verbs taking

<sup>8</sup>I use Dutch examples here (and below, to illustrate the value *vp012*), since there is strong evidence for considering these verbs ergative in Dutch. For English, there is far less strong evidence for the existence of ergative verbs.

one argument, e.g. Dutch *schijnen* ‘seem’, *blijken* ‘turn out’; copular verbs such as *become*, *get*; impersonal transitive verbs such as *regenen* ‘to rain’, *hagelen* ‘to hail’ as in *het regende verwijten lit.* ‘it rained reproaches’; *het hagelde kogels lit.* ‘it hailed bullets’.<sup>9</sup>

**vp100** Intransitive verbs, e.g. *to dance*, *to work*, *to sleep*.<sup>10</sup>

**vp120** Dyadic verbs, including transitive verbs, e.g. *to see*, *to build*, *to kill*, *to try*, *to believe*.

**vp012** Dyadic ergative verbs, e.g. *opvallen* ‘to strike’, *meevallen* ‘to turn out better than expected’, *ontgaan* ‘to fail to notice’; dyadic raising verbs (e.g. *lijken* ‘to seem’).

**vp123** Verbs which take three arguments, e.g. ditransitive verbs (*to give*, *to sell*), *to accuse*.

**patterns** This attribute takes a set of verb patterns as its value. The number of possible verb patterns is very large (approx. 125 possible values) so I will not list all of them here.

A *verb pattern* is an atomic value that specifies how the internal arguments of a verb should be realized syntactically, i.e. what their syntactic category should be (*NP*, *PP*, *ADJP*, etc.), what grammatical relation (*direct object*, *indirect object*, etc.) they bear, whether accompanying words should be present (e.g. certain prepositions, or certain expletive pronouns (e.g. *het*, *it*), etc.

The verb patterns of Rosetta only specify how the *internal arguments* are to be realized syntactically. Verb patterns that specify the syntactic realization of more than one argument are also atomic.

This is one of the reasons why there are so many verb patterns. Another reason why there are so many verb patterns can be found in the fact that *optional arguments* do not exist in Rosetta (cf. *Fixed-Arity Condition*). Instead, it is assumed that all arguments are obligatory, but that certain arguments can be realized by the special category *EMPTY*, which has no overt reflex in surface trees. Whether an argument can be realized by the special category *EMPTY* or not is indicated by means of separate verb patterns.

If a verb has more than one possibility to realize its arguments syntactically, then this is specified by a separate verb pattern for each different realization. This is the reason why the verb-pattern attribute takes a *set* of verb patterns as its value.

<sup>9</sup>The distinction between impersonal transitives and ergatives is made by an additional attribute, not discussed here.

<sup>10</sup>Certain normally intransitive verbs can be used transitively when combined with so-called ‘internal objects’, e.g. *to dance the tango*. Such examples must be dealt with in Rosetta as transitive verbs. Such verbs may be related by a rule to their intransitive counterparts, see section 4.7.

## 4.4 Covert Arguments

It has already been pointed out that certain arguments need not have an overt realization in surface trees. In this section the kinds of covert arguments which occur in the grammars will be dealt with systematically.

First, covert direct and indirect objects, and the optional presence of *by*-phrases will be discussed. Such cases are called *implicit arguments*. Next, a special kind of covert argument which occurs with certain verbs which take sentential complements is dealt with. Third, the treatment of various kinds of covert subjects in finite clauses is presented (i.e. non-overt subjects in Spanish finite clauses, and covert subjects in imperative constructions). Finally, I will briefly discuss the treatment of covert subjects of infinitival clauses, and how these are interpreted.

First, verbs that take a direct object optionally such as *eat*, *drink*, *etc.* (pseudo-transitive verbs) are assumed to be dyadic predicates. The first argument is realized as a subject in simple, active sentences. The second argument can be realized by a special element, called EMPTY. Start rules combine a pseudo-transitive verb with two variables; after the application of pattern transformations, the first variable functions as the subject and the second one as the direct object. There are special rules that substitute the special element EMPTY for variables. These rules substitute EMPTY for the object variable, and immediately delete the EMPTY element. As a consequence, there will be no surface realization of the second argument. It is assumed that the element EMPTY has the meaning ‘someone, something’, so that absence of the direct object with pseudo-transitives has the meaning of existential quantification over the second argument.<sup>11</sup> Such arguments, which are dealt with by the special element EMPTY, will be called *implicit arguments*.

The same method is used to account for optional indirect objects (cf. *I gave a book* v. *I gave him a book*), and for optional prepositional objects (cf. *I looked* v. *I looked at her*). In the latter case not only the argument, but also the preposition must be absent.

In passive structures the *by*-phrase is optional. Absence of this *by*-phrase is accounted for by assuming that the relevant argument is the special element EMPTY, which (along with the preposition *by*) is deleted upon substitution.

The variables for covert arguments discussed above are present in syntactic representations until some point in the derivation. There is evidence that this is necessary. Their presence makes it possible to account for the fact that verbs can be combined with certain adverbs which require the presence of an agent and impose selectional restrictions upon these (so-called agent (or better: controller)-

<sup>11</sup>[Zubizarreta, 1985, 250] suggests that the second argument of a verb such as *eat* is a constant with the meaning of ‘food’. By assuming that it is a constant, she accounts for the fact that *eat* can be used intransitively, and by the assumption that the constant has the meaning ‘food’ she accounts for the fact that a sentence such as *the baby is eating* implies that the baby is eating some sort of food, but not e.g. a marble. In the account sketched here, the second argument of *eat* is not a constant, but a variable bound by a covert existential quantifier. So, the verb can be used intransitively. With respect to the semantics, I assume that the type restrictions that *eat* imposes on its second argument (i.e. *edible entities*) account for the implication mentioned.

oriented adverbs, such as *intentionally*, *enthusiastically*). Typical minimal pairs to illustrate this can be constructed by opposing passive constructions with middle constructions: in passive constructions the argument expressed in the *by*-phrase can license the presence of an adverb such as *intentionally*, even when it is not overtly represented (as in 59), but in middle constructions where no comparable argument can be present such adverbs are not allowed:<sup>12</sup>

- (59) a The ship was sunk (by them) intentionally  
 b \*The ship sank intentionally (\*by them)

The covert arguments can also control covert subjects of infinitives, as in (60):

- (60) a He signaled to lower the flag  
 b He suggested not to do this.  
 c It was decided to leave at six  
 d The ship was sunk in order to collect the insurance money

In (60a,b) the covert indirect object is interpreted as the subject of the infinitival clause. In ((60c) the covert *by*-phrase argument controls the subject of the embedded infinitival complement, and in (60d) the covert *by*-phrase argument controls the subject of the adverbial infinitival purpose clause.

In certain examples, covert arguments clearly do not have the meaning of ‘someone, something’, but are non-overt ways to express the meaning of definite pronouns. The English example *I know* does not mean that there is something that I know, but rather that I know something very specific, probably mentioned before in the discourse. The translation into other languages also differs from the translation of the covert arguments discussed above. The examples discussed above are translated into covert arguments in most cases in Dutch and Spanish, or, in certain configurations, into indefinite pronouns. The covert argument in *I know*, however, must be translated into the definite neuter pronoun (*het* in Dutch; *lo* in Spanish). It seems best then, to treat the covert argument in *I know* by an abstract element different from EMPTY. In Rosetta this element is called Pro-Sent, because in effect it is a pronominal element that occurs *instead of* (*Pro*) *sentential structures* (*Sent*). This element happens to be empty in English, but not in Spanish. In Dutch, this element can be overt (*het*, or *dat*), or it can be covert when in the topic position (cf. *weet ik* lit. *know I*, ‘I know’).

In Spanish, pronominal subjects of finite verbs tend to be non-overt. For these elements (they can be arguments or non-arguments, e.g. expletives) it has been assumed that pronominal subjects are present during almost all the derivation, and are deleted late in the derivation. The consequence is that all rules will act ‘as if’ there were a subject throughout the derivation, which yields exactly the desired results, as is well-known.

Imperatives usually do not occur with subjects (though they can under specific circumstances). Again, there is considerable evidence that all rules should act as

<sup>12</sup>See [Chomsky, 1986, 117ff] for discussion of these and related facts in a different framework, and with a partially different analysis.

if there were a (second-person) subject throughout. It has been assumed that imperatives are always accompanied by a subject. Under appropriate conditions this subject is deleted at a late point in the derivation. A similar analysis is used for first-person plural imperatives such as *comamos* ‘let’s eat’ in Spanish.

Many infinitival constructions appear without an overt subject, even if the verb heading the infinitive requires the presence of a subject in all finite clauses and in certain infinitival clauses and requires an additional argument to be saturated. Furthermore, infinitival constructions can take an overt subject under appropriate conditions. It is well-known that many rules act as if there were a subject in infinitival constructions even when it is not overtly present. For these reasons<sup>13</sup> it has been assumed that infinitival constructions always have a subject. When the subject is not overt at the surface, one of the following situations can hold:

- if the infinitival clause is a complement of a so-called *subject-raising* verb, the subject of the infinitive is moved to the subject position of the embedding verb. This position is always unoccupied with such verbs, since they do not take an argument in subject position. An example is given in (61a), where the subject of the finite clause originates as the subject of the embedded infinitival clause, as in (61b):

- (61) a [The man seems [ to be ill]]  
       b [ seem- [ the man to be ill]]

- if the infinitival construction is a complement of a *control* verb (or adjective), the subject of the infinitive is a syntactic variable which must be co-indexed with some argument (the controller) of the embedding verb. The embedding predicate specifies which argument this is.<sup>14</sup> Control transformations check whether the correct relation between the argument of the embedding predicate and the subject of the infinitive holds, and if so, they delete the syntactic variable which occupies the subject position of the infinitival clause. This is illustrated in (62). In (62a) the embedded infinitival complement contains a subject variable which is co-indexed with the subject of the finite superordinate clause. A control transformation deletes this subject variable under appropriate conditions. The result is represented in (62b). After the application of substitution rules to the two variables  $x_1$  and  $x_2$ , a sentence such as (62c) can be derived by substituting the appropriate NPs for the variables. It is now correctly encoded that the subject of the infinitival complement and the subject of the finite verb are co-indexed, and syntactic reflexes of this (e.g. if the embedded clause contains reflexive pronouns) will be accounted for correctly.

<sup>13</sup>Additional reasons relate to translation: it must be possible to translate certain finite clauses into infinitival clauses or vice versa, e.g. Dutch *Ik denk ziek te zijn* into *I think that I am ill* and not into the more literal but ill-formed *\*I think to be ill*. This translation problem is considerably simplified if both infinitival and finite constructions are clauses which contain a subject.

<sup>14</sup>This is done by specifying the grammatical relation of the controller. It is well-known that such a specification is insufficient to deal with all relevant facts, so improvement is possible here. See [Růžicka, 1983] for discussion and an alternative.

- (62) a  $[x_1$  promised  $x_2$  [ $x_1$  to do this]]  
 b  $[x_1$  promised  $x_2$  [to do this]]  
 c [John promised Mary [to do this]]

This approach, in which the control relation involves a relation between syntactic variables, immediately accounts for the fact that non-overt subjects of infinitives in these constructions must be arguments and cannot be meaningless elements such as expletives or idiom chunks (see [Schenk, 1992]).

- if the relation between controller and antecedent is less strict, e.g. if there is more than one controller (split antecedents) or if the controller is the abstract element *EMPTY*, the subject of the infinitive is an abstract indexed pronoun, which is deleted by special control transformations. An example of split antecedents is given in (63):

- (63) a John agreed with Pete to leave together  
 b  $[x_1$  agreed with  $x_2$  [ $PRO_{1,2}$  to leave together]]

In (63) the subject of the infinitive is interpreted as the set containing both John and Pete. This phenomenon of split antecedents has only partially been implemented, since it was assumed that it should be accounted for by more general rules of discourse grammar which interpret pronouns.

An example where the controller is an implicit argument is given in (64): the implicit argument is represented by a variable for the special element *EMPTY* and it is co-indexed with the abstract pronoun *PRO*, as indicated in (64b). A control transformation deletes *PRO*, and the normal *EMPTY* substitution rules substitute *EMPTY* for *EMPTYVAR* and deletes *EMPTY* and *by*, resulting in (64a).

- (64) a It was decided not to do this anymore  
 b It was decided by *EMPTYVAR*<sub>1</sub> [*PRO*<sub>1</sub> not to do this anymore]

- if there is no controller at all, or if the controller is too far away in the structure, the subject is an abstract pronoun. It can either be an abstract pronoun with the approximate meaning of *one* in English, or an abstract personal pronoun. Its interpretation is left to rules which interpret pronouns in discourses. An example of this situation is given in (65). The infinitival relative clause contains an abstract pronoun. After it is interpreted by discourse rules, it is deleted by special control rules.

- (65) a The books to read are on the table  
 b The books [ *PRO* to read] are on the table



## 4.5 Bound Adjuncts and Adverbials

In this section I will discuss the status of adjuncts in more detail. In particular, I will propose a distinction between two kinds of adjuncts, *free* adjuncts and *bound* adjuncts.

As pointed out before, *adverbials* (or *free adjuncts*) have a much looser syntactic connection to the predicate word with which they co-occur. They are always optional, can generally appear with all predicate words (though sometimes there are semantic restrictions), their form is not determined by the predicate, and in principle they can be iterated (subject to semantic constraints). Typical examples of *adverbials* are temporal and locative expressions.

The concept of argument in Rosetta is a very strict one, because of the conditions that it is subject to (i.e. fixed arity, arity preservation). There are phrases where a treatment as arguments would lead to very complex lexical specifications, and to syntactically complex and semantically ill-motivated derivations. This occurs in particular if a whole class of complements is always optional. Let us take, as an example, resultative phrases with transitive verbs. Resultative phrases are almost always, if not always, optional. Treating them as an argument of a verb would require that the arity of the verb be increased (for transitive verbs from 2 to 3), and that the patterns be adapted accordingly. Note that the number of patterns will in many cases be doubled, since both the possible presence and the possible absence of the resultative phrase must be described. So, an initial problem is that the description of subcategorization information becomes more complex. Second, it is not indicated by general rule that resultative phrases are always optional, so that a generalization is missed.<sup>15</sup>

Third, the derivation of a simple sentence which contains a subject, a direct object and a verb which allows a resultative phrase is complicated considerably if the resultative phrase is considered an argument. Though the resultative phrase does not appear in the surface tree, start rules combine the relevant verb with three variables, and one of the variables is replaced by the abstract node EMPTY which is deleted immediately.

Finally, such a derivation makes no sense from a semantic point of view. It seems unnatural to say that e.g. the meaning of *he painted the door* is something like ‘there is a state x such that he painted the door until it was in state x’. The correct semantics for this sentence should not make any reference to a resultant state at all.

These considerations suggest that it is much better not to treat resultative phrases as arguments. On the other hand, it is also not really possible to treat them as adverbials (free adjuncts), since the main verb occurring in a sentence determines whether they can appear in a sentence or not.

Therefore, a new, intermediate type of relation to predicates has been introduced in the grammar, called *bound adjuncts*. Bound adjuncts have some prop-

---

<sup>15</sup>If there are verbs which require the presence of a resultative phrase, then the resultative phrase can and must be dealt with as an argument of this verb. The problems mentioned will not arise for such a case.

erties of adverbials and other properties of arguments. Bound adjuncts look like arguments in that the predicate word determines whether they can appear or not. However, bound adjuncts are always optional and though they cannot be translated completely on their own, their translation depends only on the possible adjuncts of the translation of the predicate word. In addition the translation of the predicate word need not take bound adjuncts at all in certain cases, so there is no variant of the arity-preservation condition for adjuncts.

The introduction of *bound adjuncts* makes it possible to circumvent the problems described above in the following manner: (1) the description of predicates in the lexicon can be simplified if the potential of a predicate to take bound adjuncts is represented in a separate attribute, so that there is no multiplication of verb patterns. (2) Since every bound adjunct is optional, by definition, there is no need to specify this with individual verbs; it can be expressed by a general rule, as it should. (3) Derivations in which adjuncts do not appear in the surface tree are maximally simple, since the bound adjunct does not appear at any point in the derivation, and (4) the derivations of sentences without bound adjuncts make sense from a semantic point of view: there are no bound adjuncts in the derivation if they do not appear at the surface, so they play no role in the semantics of the expression either.

The proper division between arguments, adjuncts and adverbials is determined by considerations of fact, elegance and simplicity: the choice should be made that yields an overall simpler grammar which correctly describes the relevant facts. In the grammars developed a number of constructions have been classified as bound adjuncts, as will be illustrated below.

Each verb has an attribute *adjuncts* that indicates which bound adjuncts it can take. The attribute is set-valued. Possible elements of this set and a description of their meaning are given below:<sup>16</sup>

**ResAP** for resultative adjectival phrases, e.g. *hij brak het kapot* ‘he broke it into pieces’, *hij braadde het vlees bruin* ‘he baked the meat brown’, *hij knipte het haar kort* ‘he cut the hair short’, *hij kookte de aardappels gaar* ‘he cooked the potatoes done’, *hij verfde de deur groen* ‘he painted the door green’. This value is intended for resultative APs with verbs that take a direct object as their argument. This direct object is interpreted as the subject of this AP. It might also be used for a different class of resultatives which have an overt subject and can be combined with intransitive verbs, e.g. *hij heeft zijn schoenen kapot gelopen* ‘he ran his shoes into pieces’, *hij heeft zijn lippen kapot gebeten* ‘he bit his lips into pieces’, *hij heeft de tube leeg geknepen* ‘he squeezed the tube empty’, etc. However, it appears that this construction can be used with any real intransitive (i.e. not ergative) verb which expresses an activity. If so, there is no need to specify the possibility of taking this kind of resultative constructions with each individual verb.

<sup>16</sup>The bound adjuncts *subjcomit*, *objcomit*, *resnp* and *topicadjunct* have been implemented only partially. The rest have been implemented in full.

This value *cannot* be used for transitive verbs that can leave out their direct object and replace it by a resultative small clause, as e.g. in *hij heeft de winkel leeg gekocht* ‘he bought the shop empty’, *hij heeft zijn bord leeg gegeten* ‘he ate his plate empty’, because then the *Fixed-Arity Condition* would be violated. See section 4.7 for a suggestion of how to deal with these cases.

**ResPP** for resultative PPs (with transitive verbs), e.g. *hij zaagde het hout in stukken* ‘he sawed the wood into pieces’, *hij brak het hout in stukken* ‘he broke the wood into pieces’, and for resultative PP-headed small clauses (with intransitive verbs), e.g. *hij deelde de taart in tweeën* ‘he sliced the pie into two parts’, *hij beet het brood in stukken* ‘he bit the bread into pieces’, *hij knoopte de touwen aan elkaar* ‘he knotted the ropes together’.

**SubjComit** for subject-oriented co-arguments, e.g. **met iemand / iets** *afwisselen* ‘to alternate with someone / something’, *iets met iemand bespreken* ‘to discuss something with someone’, **met iemand** *corresponderen* ‘to correspond with someone’, etc. It involves expressions that require a semantically collective subject argument, unless there is a comitative phrase in the sentence related to the subject argument.

**ObjComit** for object-oriented co-arguments, e.g. *iemand met iemand verbinden* ‘to connect someone with someone’, *iets met iets combineren* ‘to combine something with something’, *iets met iets vergelijken* ‘to compare something with something’, etc. It involves expressions that require a semantically collective object argument, unless there is a comitative phrase in the sentence related to the object argument.

**BenefactNP** for beneficiary NPs. These are marginal in Dutch, but fully possible in English and Spanish: *?hij heeft ons een boot gemaakt*, *he made us a boat*, *he baked her a cake*.<sup>17</sup>

**BenefactPP** for beneficiary PPs. These PPs are headed by *voor* in Dutch, by *for* in English and by *para* in Spanish: *Hij heeft voor ons een boot gemaakt*, *he made a boat for us*, *compró un bocadillo para él* ‘he bought a sandwich for him’.

**Diradjunct** for optional directional PPs and ADVPs. There are examples of directional PPs and ADVPs with verbs that take a direct object (transitive or ergative), e.g. *he brought her to the factory*, *she returned to her birthplace*, *he came home*. There are also examples of directional small clauses (with intransitive verbs), e.g. *hij duwde de kar naar boven* ‘he pushed the car up’, *hij pompte het water omhoog* ‘he pumped the water up’, *hij heeft de kinderen naar huis gereden* ‘he drove the children home’, *zij hebben de atleet naar de finish gelopen* ‘they ran the athlete to the finish’.

<sup>17</sup>Henk van Riemsdijk pointed out that ethical datives might be treated in this way as well. I agree that this might be a possibility, but ethical datives have not been dealt with in the Rosetta grammars at all.

**ResNP** This is intended for ‘resultative NPs’. I only know of idiomatic cases such as *hij rende zich een ongeluk* ‘he ran himself an accident’, *hij heeft zich een beroerte gelopen* ‘he ran himself a stroke’. They usually involve the reflexive pronoun *zich* and an additional NP. It is not fully clear what grammatical relations these NPs bear (indirect object + direct object, or direct object + predicative NP).

**TopicAdjunct** for topic phrases, e.g. *he wrote (a book) about linguistics*.<sup>18</sup>

As pointed out above, verbs with certain adjuncts can be translated into verbs without adjuncts. Of course, there must be some way of expressing what is expressed by the adjunct in the source language, but this need not be an adjunct in the target language. E.g. it is possible in Rosetta to translate resultative adjective phrases into subordinate sentences introduced by *until*, e.g. *hij boende de kamer schoon* into *he polished the room until it was clean*. This can in principle be done for the adjuncts *resAP*, *resPP* and *diradjunct* (but only with intransitive verbs). Verbs with *benefactNP* or *benefactPP* can only be translated into verbs with *benefactNP* or *benefactPP*. Verbs with *subjComit* or *objComit* can only be translated into verbs with *subjcomit* or *objComit*.

## 4.6 Small Clauses

*Small clauses* are constructions which consist solely of a subject and a predicate with its internal arguments (see [Stowell, 1981]). They differ from real clauses (which I will call *full clauses*) in that properties such as mood, tense and aspect are not expressed. Small clauses are quite pervasive in the current grammar.<sup>19</sup> All full clauses are built on the basis of partially derived small clauses, by adding formal means to express tense, aspect, mood and additional properties which distinguish full clauses from small clauses. In addition, small clauses can be constructed in subgrammars which are partially isomorphic<sup>20</sup> to the subgrammars which form full clauses. Small clauses and full clauses can be formed around each of the categories noun, verb, adjective, adverb and preposition. The syntactic categories of small clauses are NPP, VPP, AdjPP, etc. and they differ categorially from ‘simple phrases’ such as NP, VP, AdjP, etc.

Small clauses in Rosetta can have overt subjects or non-overt subjects. This deviates from most other analyses in which small clauses play a role. Usually, only small clauses with overt subjects are supposed to exist. (see e.g. [Stowell, 1991]).

<sup>18</sup>Henk van Riemsdijk suggested that expressions with the preposition *van* (e.g. *de man waarvan hij zei dat Piet hem haatte* ‘the man of whom he said that Pete hated him’) might be dealt with in this manner as well. I agree, but such expressions are not dealt with in a systematic manner in the current grammars.

<sup>19</sup>See [Odiijk, 1989] for a description of the pervasive role of small clauses in the Rosetta grammars.

<sup>20</sup>See [Appelo et al., 1987], [Odiijk, 1989], [Rosetta, 1993] and [Appelo, 1993] for the notion ‘partial isomorphy’.

Small clauses with non-overt subjects behave as predicates with one unbound variable. These small clauses have a syntactic variable as their subject plus a marking that this subject must be interpreted by a set of special rules. Their treatment is more or less analogous to the way controlled infinitival complements are dealt with, though the conditions on appropriate antecedents differ considerably. Several examples of small clauses without overt subjects will be discussed in this section.

Special rules check whether there is an appropriate antecedent for the syntactic variable in the subject position, and, if so, the abstract variable is deleted. The following approach has been adopted, which is a variant of a proposal by [Bresnan, 1982, 322]:<sup>21</sup> the antecedent of a small clause is the c-commanding direct object if there is one, otherwise it is the c-commanding subject, if there is one, otherwise there is no antecedent (and the structure will be ill-formed). This formulation generalizes to all uses of small clauses known to me. Two special cases are worth mentioning.

First, pseudo-transitive verbs which are combined with a small clause with a non-overt subject (e.g. *he painted (the door) combined with green*) cannot leave out their object if the small clause is present. This follows from the grammar in the following manner. It is well-known that predicates require that the argument which they interpret as their subject cannot be implicit. This property, which is sometimes called ‘Bach’s Generalization’ ([Bresnan, 1982, 373]) and which they share with a number of other phenomena, e.g. certain kinds of control and binding of reflexives (see [Williams, 1981], [Koster, 1984]), has been captured by formulating the rules for the interpretation of non-overt subjects of small clauses in such a way that the abstract element EMPTY cannot serve as an appropriate antecedent. However, this abstract element is a direct object in a sentence such as *\*he painted green*, and therefore must be the antecedent of the subject of *green*. A contradiction is derived, so the relevant sentence cannot be generated, which is a desirable consequence.

Second, it is interesting how the interaction with passivization is taken care of, and how this accounts systematically for one part of ‘Visser’s generalization’, as [Bresnan, 1982] calls it. As will be argued in chapter 5.2, the formation of passive structures consists of a number of rules. One of these rules has as one of its tasks to turn the subject into a *by*-phrase. A second rule preposes NPs, among them direct objects, into the subject position in passive and other structures. An essential feature of the rules for the interpretation of subjects of small clauses is that they are ordered between these rules to form passive structures. This results in the following consequences. A verb which takes a direct object can be passivized even when there is a small clause: the subject of the small clause can take the direct object as its antecedent, in accordance with the rule given above. A verb which does not take a direct object, but only a subject and a small clause as one of its arguments, cannot be passivized: the rule which creates the *by*-phrase removes

---

<sup>21</sup>Bresnan’s proposal is defined in terms of LFG’s *f*-structure objects, and mentions a grammatical relation *OBJ2* which is not assumed here. Furthermore, Bresnan’s proposal is intended as the unmarked case for all ‘functional control’ relations, while we restrict attention to the interpretation of subjects of small clauses.

the subject, so that there is no longer an antecedent for the subject of the small clause. As a consequence, verbs with so-called ‘subject complements’ (e.g. *smell*, *sound*, but also verbs such as *cost*, *weigh*, *last*, *etc.*, which are analyzed as taking ‘subject complements’ in the current grammar) cannot passivize. So one part of Visser’s generalization is derived by the design of the grammar.

Small clauses can be used in larger structures, e.g. as an adjunct or as an argument to a predicate. Though small clauses are formed as constituents in the syntax, there is no need for them to remain constituents throughout the derivation. In surface trees, small clauses need not occur, and in the current grammar they usually do not occur there. If a small clause with an overt subject is a complement to a verb, this overt subject is usually removed from the small clause and made a direct object. The remaining small clause is turned into a simple phrase (NP, AdjP, etc). If a small clause does not have an overt subject, the abstract subject is deleted under identity with an antecedent, and the remaining subjectless small clause is turned into a simple phrase. Note that the formalism does not force one to such an analysis: one might adopt an analysis where small clauses occur in surface trees as well, and one might also adopt a mixed analysis, in which certain small clauses occur in surface trees but others don’t. In the current grammar, however, there happen to be no cases of small clauses in surface trees, though there might well be constructions not dealt with now where this might be required, e.g. absolutive *with* constructions.

I will illustrate how different kinds of small clauses can be used in larger structures. Restricting attention to small clauses headed by adjectives, we can distinguish the following cases:<sup>22</sup>

- (66) a *hij* is *gek* ‘he is mad’  
 b *hij* werd *gek* ‘he got mad’  
 c *hij* kneep *de tube leeg* ‘he squeezed the tube empty’  
 d *hij* kocht *de winkel leeg* ‘he bought the shop empty’  
 e *hij* vond *haar aardig* ‘he found her nice’
- (67) a *hij* schoor *de man glad* ‘he shaved the man smooth’  
 b *hij* at het vlees *naakt* ‘he ate the meat naked’  
 c *hij* at het vlees *rauw* ‘he ate the meat raw’  
 d *dat* klinkt *mooi* ‘that sounds beautiful’

First, the examples can be divided in small clauses with an overt subject (as in (66)) and small clauses without an overt subject (as in (67)). Note that the subject of the small clause need not (and usually does not) remain the subject of the small clause in surface trees.

In (66a,b) small clauses headed by adjectives are used in combination with copular verbs. The subjects of the small clauses have become the subjects of the copular verbs in these examples. Copular verbs do not have a special status in

<sup>22</sup>Sentence (66d) has an additional reading which can paraphrased as ‘he bought the shop when it was empty’. This interpretation must be ignored here, and in the rest of this chapter.

Rosetta: they are verbs which happen to take an internal argument which is a small clause with an overt subject and no external argument. Only the verb *to be* has a special status in the grammar. This verb is not a basic expression, but a verb which is introduced syncategorematically. This is motivated by the desire to be able to paraphrase small clauses by full clauses and vice versa, so that it also becomes possible to translate between small and full clauses. For this reason (66a) and (66b) are derived in different ways, though their surface trees are similar.

In (66c,d,e) small clauses with overt subjects are combined with a normally intransitive verb (66c), with a normally transitive verb (66d) and with a verb which can take small clauses as complements (66e). I will discuss these examples and the differences between them in more detail below.

In (67a) a transitive verb is combined with a small clause with a non-overt subject. This example will be discussed in more detail below, where it is compared with the examples (66c,d,e).

In (67b,c) small clauses are used as ‘secondary predicates’ to the subject and the direct object respectively. It is assumed that the small clause which takes the subject of the clause as the antecedent for its subject occupies a position other than the small clause which takes the object of the clause as its antecedent: the first one occupies a position outside VP, the second one inside VP. Though the existence of ‘secondary predicates’ have always been taken into account in the analysis, they have not yet been incorporated in the current grammars.

Examples such as (67d)<sup>23</sup> have been dealt with in the following manner. It is assumed that the verbs take two arguments. One argument is an external argument, and the other is an internal argument which must be a small clause with a non-overt subject headed by an adjective. The subject of the small clause takes the subject of the verb as its antecedent (See also [Hoekstra, 1984, 120] for some discussion of a specific way of using the verb *smaken*). This analysis might be controversial in at least three respects. These can be summarized as follows:

- Do these verbs take adjectives or adverbs as complements?
- Do these verbs take an external argument?
- Is the subject an argument of the verb?

With respect to the first issue, this is hardly a significant question in Dutch, since no relevant formal distinctions can be made between predicative adjectives and (the relevant class of) adverbs in Dutch. In English, however, these verbs usually occur with adjectives, not with adverbs (e.g. *it sounds good*, *it smells nice*).

With respect to the second issue, it was decided to deal with these verbs in Rosetta as verbs which take an external argument, but the issue has not really been investigated. A quick first glance at the facts, however, does not clearly indicate that these verbs are ergative. The verbs are conjugated in Dutch with the verb *hebben* (not with *zijn*). The verbs cannot be passivized (*\*Er werd lekker gemaakt door het ijsje* lit. *It was tasted nice by the ice-cream*). This would follow from an

<sup>23</sup>Other examples are Dutch *smaken*, *ruiken*, *aanvoelen(?)* and English *taste*, *smell*, *feel*.

ergative analysis, but it also follows from the fact that the verb takes a ‘subject complement’. The verbs cannot appear without their subject as a complement to the causative verb *laten* ‘to let’. In general, one can say that this is possible for verbs taking an external argument (cf. *Hij liet er dansen* lit. *He let dance there* v. *\*Hij liet er vallen* lit. *he let fall there*), though additional conditions may blur this simple picture. From this fact, however, no conclusions can be drawn. Even if the verbs are analyzed as taking an external argument, these constructions will be ill-formed for independent reasons: if the subject is absent, the small clause will not have an antecedent, and the sentence is ill-formed. The verbs cannot appear as participles in nominal expressions (*\*de lekker gesmaakte ijsjes* lit. *the nice tasted ice-creams*, *\*de mooi geklonken geluiden* lit. *the beautiful sounded sounds*). It is not fully clear what this shows. If such constructions are possible for ergative verbs and impossible for non-ergative verbs, the verbs under consideration must be considered non-ergative. If these constructions can only be formed from ergative verbs which are conjugated with *zijn* ‘to be’, then this test is neutral with respect to the ergativity of these verbs. The facts relating to extraction in *wat voor* constructions do not point to the ergativity of these verbs either: sentences such as *\*Wat klonk er gisteren voor geluid mooi?* lit. *What sounded there yesterday for sound beautiful?*, *Wat smaakte er gisteren voor gerecht lekker?* lit. *What tasted there yesterday for dish nice?* are completely out, though *Wat voor geluid klonk er gisteren mooi?* lit. *What for sound sounded there yesterday beautiful?* ‘What kind of sound sounded beautiful yesterday?’ and *Wat voor gerecht smaakte er gisteren lekker?* lit. *What for dish tasted there yesterday nice?* ‘What dish tasted nice yesterday?’ are perfect. This would be consistent with a non-ergative analysis; it is perhaps consistent with an ergative analysis, since *wat voor* split always yields somewhat deviant results if the relevant NP is followed by AP. The verb *smaken* ‘to taste’ can optionally take an indirect object: *De maaltijd smaakte ons goed* ‘The dish tasted good to us’. If one believes that indirect objects cannot appear if there is no direct object, then this would point to an ergative analysis. Under such an analysis, one would also expect that the indirect object can precede the subject (as in *?Heeft jou de maaltijd goed gesmaakt?* lit. *Has you the meal good tasted?* ‘Did you like the meal?’ and *?...dat ons deze maaltijd goed gesmaakt heeft* lit. *...that us this meal good tasted has* ‘...that we liked this meal’), but judgments on these sentences vary. Reflexives and reciprocals can hardly occur with these verbs. Only the verb *smaken* ‘to taste’ can take an additional argument, and sentences such as *??hij smaakt zichzelf goed* ‘he tastes good to himself’, *??zij smaken elkaar goed* ‘they tasted good to each other’ appear not only semantically odd, but also syntactically deviant. This suggests that an ergative analysis might be more appropriate. In short, a first glance at the facts does not yield a clear conclusion, and the issue requires further investigation. If it turns out better to analyze these verbs as ergative verbs, nothing prevents one from changing the relevant specifications of these verbs in the dictionary. The rest of the grammar will function normally, and probably better if the arguments for treating these verbs as ergatives are sound.

Concerning the final issue, it has been decided in the Rosetta grammar to



adopt an analysis with two (or for *smaaken* even three) arguments for the following reason. Sentence (68a) is deviant:

- (68) a ??Dat geluid smaakt goed ‘That sound tastes good’  
 b Dat geluid klinkt/is goed ‘That sound sounds/is good’  
 c Dat gerecht smaakt goed ‘That dish tastes good’

The reason for this deviance cannot be found in the fact that the NP *dat geluid* and the adjective *goed* are incompatible in such constructions (see (68b)). It cannot be caused by the combination of the verb *smaaken* ‘to taste’ and the adjective *goed* ‘good’, see (68c). It must then be a consequence of the combination of the NP *dat geluid* ‘that sound’ and the verb *smaaken* ‘to taste’, a plausible conclusion on intuitive semantic grounds as well. These facts can be accounted for by assuming that the verb *smaaken* imposes type restrictions on the NP which are incompatible with the NP *dat geluid*. Type restrictions, however, cannot be imposed by a predicate on free adjuncts, so it must be concluded that *dat geluid* is an argument of the verb *smaaken*.

Of course, the relevant NP is also subject to type restrictions imposed by the adjective. Examples such as ??*de aardappels smaken geschift* ‘the potatoes taste curdled’ v. *de melk smaakt geschift* ‘the milk tastes curdled’, or ??*het water smaakt onrijp* ‘the water tastes unripe’ v. *de appel smaakt onrijp* ‘the apple tastes unripe’ might be used to illustrate this. But this does not necessarily imply that the NP is an argument of this adjective. The adjective can impose type restrictions on the (abstract) subject of the small clause it heads. Since the subject of the verb serves as an antecedent for this subject, these type restrictions will hold for the subject of the verb as well.<sup>24</sup>

It remains to discuss the examples (66c,d,e) and (67a). For convenience, I will repeat the relevant sentences in (69).

- (69) a hij kocht de winkel leeg ‘he bought the shop empty’  
 b hij kneep de tube leeg ‘he squeezed the tube empty’  
 c hij vond haar aardig ‘he found her nice’  
 d hij schoor de man glad ‘he shaved the man smooth’

All these examples can be characterized as a sequence  $NP_1 V NP_2 AP$ . Though all these examples have surface trees which look very much alike, there are considerable differences between these constructions and they are derived in different ways. I will first outline how each of these examples is analyzed, and then sketch

<sup>24</sup>Henk van Riemsdijk suggested considering the adjective in these constructions as bound adjuncts. This is certainly a possibility, but whether it is feasible depends on one’s analysis of the sentences in which the ADJP is absent. Such sentences are possible (cf. *Dat smaakt!* ‘That tastes!’), but one cannot really say that the ADJP has just been left out, as with other adjuncts. Sentences in which the ADJP do not occur require a special intonation (cf. the exclamation mark), and are interpreted as if the verb is still modified (by an ADJP with a somewhat indefinite meaning). I have assumed that such sentences are to be analyzed as containing an abstract ADJP with exclamatory force, which induces the special intonation. Under such an analysis it is assumed that the ADJP is obligatory, therefore it cannot be an adjunct.

some of the considerations that lead to these analyses. Example (69a) is analyzed as containing a verb which takes a resultative small clause with an overt subject as its obligatory internal argument. Example (69b) is analyzed as containing an intransitive verb which is accompanied by a resultative small clause with an overt subject as a bound adjunct. Example (69c) is analyzed as containing a verb which takes two arguments, a subject, and a non-resultative small clause with an overt subject. Example (69d) is analyzed as containing a transitive verb which takes a subject and a direct object as its arguments and which is accompanied by a resultative small clause with a non-overt subject as a bound adjunct.

The differences between these analyses can be represented as follows: first, the analyses differ with respect to the status of NP<sub>2</sub>: is it an argument of the verb (in (69d)), or not (in (69a,b,c)). Second, is the combination NP<sub>2</sub> AP an argument of the verb (in (69a,c)) or not (in (69b,d)). This is represented in the table in (70):

|                                | (69a) | (69b) | (69c) | (69d) |
|--------------------------------|-------|-------|-------|-------|
| (70) NP <sub>2</sub> arg of V? | no    | no    | no    | yes   |
| NP <sub>2</sub> AP arg of V?   | yes   | no    | yes   | no    |

These criteria do not distinguish (69a) and (69c), which differ in the resultative v. non-resultative nature of the small clause argument.

The analyses are supported by several facts. First, let's consider omissibility of NP<sub>2</sub> and AP. In (69a) neither NP<sub>2</sub>, nor AP, nor their combination can be omitted (cf. *hij kocht \*(de winkel) \*(leeg)*). It might look as if the AP on its own can be omitted in the example given, but this is accidental. The resulting sentence happens to be grammatical because the verb is transitive and the semantic type of the NP happens to be compatible with the possible types of the object of this verb. But the resulting sentence cannot be seen as a variant of the original sentence in which the AP is left out. Semantic considerations corroborate this view (see below). In general, AP cannot be left out in this construction (cf. *hij at zijn bord ??(leeg), hij at zijn buikje ??(vol/rond)*). Both can be omitted together in (69b), but neither of them can be omitted on its own (cf. *hij kneep, \*hij kneep de tube, \*hij kneep leeg*). NP and AP cannot be omitted in (69c), neither on their own, nor together (cf. *\*hij vond, \*hij vond haar,<sup>25</sup> \*hij vond goed*). AP can, but NP cannot, be omitted in (69d) (cf. *hij schoor hem, \*hij schoor glad*).

These facts follow immediately from the analysis proposed, in combination with independent properties of the relevant verbs and of predicative APs. In (69a) NP and AP together form a small clause which is an obligatory argument of the verb, so that NP and AP cannot be omitted together. The NP cannot be omitted on its own, since a small clause must contain a subject. The AP cannot be left out on its own, because a small clause must contain a predicate.

In (69b) NP and AP can be omitted together because together they are a bound adjunct, and adjuncts are always optional. NP and AP cannot be omitted on their own for the same reasons as given for (69a).

<sup>25</sup>This sentence is well-formed under an irrelevant reading, which must be ignored here.

In (69c) NP and AP cannot be omitted together, because together they form an obligatory argument of the verb. They cannot be omitted on their own for the same reasons as given for (69a).

In (69d) NP cannot be omitted, because it is an obligatory argument of the transitive verb. The AP can be omitted because it is the surface reflex of a small clause adjunct, and adjuncts are always optional.

A second argument relates to entailments. The issue is whether certain entailments can be formulated at all. An initial example is the implication  $NP\ V\ NP\ AP \Rightarrow NP\ V\ NP$ . This implication cannot be formulated for (69a,b,c), but it can be formulated for (69d). (cf. *\*Hij kneep de tube, \*Hij vond haar, hij schoor de man glad*  $\Rightarrow$  *hij schoor de man*).

It may appear that the entailment can be formulated for (69a) as well, but this is just accidental. If other examples of the same type are taken into account, one observes that in general the entailment cannot be formulated (e.g. *hij at zijn buikje \*(rond), hij at zijn bord \*(leeg)*).

A number of conditions determine whether the entailment holds or not. However, if the entailment can be formulated at all for examples such as (69a), it never holds, though in (69d) and similar examples it may hold, cf. *Hij kocht de winkel leeg*  $\Rightarrow$  *hij kocht de winkel* is false; *hij schoor de man glad*  $\Rightarrow$  *Hij schoor de man* is true.

A second example is the entailment  $NP\ V\ NP\ AP \Rightarrow NP\ V$ . This implication can be formulated for (69b), but it cannot be formulated for (69a,c,d) (cf. *\*hij kocht, hij kneep de tube leeg*  $\Rightarrow$  *hij kneep, \*hij vond, \*hij schoor*).

Generally, such entailments can be formulated if the omitted phrases in the right hand part of the implication are an adjunct. In the analysis sketched, the relevant facts follow immediately: in the cases where the implication can be formulated, the omitted phrases are bound adjuncts; in the cases where the implication cannot be formulated, the omitted phrases are not adjuncts.

A third argument which provides evidence for the analysis proposed concerns paraphrases. How can the relevant sentences be paraphrased, in particular can they be paraphrased by expressions of the form  $NP_1\ V\ NP_2\ \text{zodat (so that)}\ NP_2\ AP\ V'$ , where  $V'$  is an appropriate copula (e.g. *worden, raken* 'become'), and where the second occurrence of  $NP_2$  is perhaps better replaced by an appropriate pronoun; or can they be paraphrased by expressions of the form  $NP_1\ V\ \text{dat (that)}\ NP_2\ AP\ V'$  where  $V'$  is an appropriate form of the copula *zijn*; or can they be paraphrased by an expression of the form  $NP_1\ V\ \text{op zo'n manier, dat (in such a way, that)}\ NP_2\ AP\ V'$ , where  $V'$  is an appropriate copula (e.g. *worden, raken* 'become')?

An expression of the first form can be a paraphrase for (69d), cf. *hij schoor de man zodat die glad werd*, but not for the other examples. An expression of the second form can be a paraphrase for (69c), cf. *hij vond dat zij aardig was*, but not for the other examples. An expression of the third form can be a paraphrase for (69b), cf. *hij kneep op zo'n manier, dat de tube leeg raakte*.

These facts follow from the analysis in the following way. A necessary condition for the formulation of an expression of the first form is that  $NP_2$  is an argument

of V. This is the case, according to the analysis proposed, only in (69d), but not in the other examples. A necessary condition for the formulation of the second expression is that the combination  $NP_2 AP$  is an argument of V, which is only the case in (69c). And a necessary condition for the formulation of the third expression is that the combination  $NP_2 AP$  is an adjunct, which is only the case in (69b).

A fourth argument in favor of the analysis proposed concerns selectional or type restrictions: can V impose type restrictions upon  $NP_2$ ? If so, then  $NP_2$  is an argument of V, if not, it is not. The verb can impose type restrictions upon  $NP_2$  in (69d), cf. *??hij schoor de tafel glad* ‘he shaved the table smooth’, but not in the other examples. Other examples of this kind also show that the verb can impose selectional restrictions upon the direct object, e.g. *??hij brak het water (kapot)* ‘he broke the water into pieces’, and with verbs such as *zetten, leggen, doen, stoppen* (all meaning ‘to put’). These latter verbs take a direct object and a prepositional or adjectival small clause with a covert subject. The combination of the semantic type of the direct object and the nature of the locative predicate determine which of these verbs is appropriate. The following examples show that the semantic type of the direct object is crucial:

- (71) a Hij ??zette / legde / stopte / deed de zakdoek in de tank  
       ‘He put the handkerchief into the tank’  
       b Hij ??zette / ??legde / stopte / deed benzine in de tank  
       ‘He put gasoline into the tank’  
       c Hij ??zette / ??legde / stopte / deed gas in de tank  
       ‘He put gas into the tank’  
       d Hij zette / legde / stopte / deed het boek in de tank  
       ‘He put the book into the tank’

This can be accounted for by means of selectional restrictions if the object NP is an argument of the verb.

Similar examples can be constructed with verbs such as *benoemen (tot), aanstellen (als)* ‘appoint’, which take a prepositional predicate phrase as one of their complements. The direct object of these verbs must be at least animate and is usually human (cf. *??hij benoemde de eerlijkheid tot de hoogste deugd* ‘He appointed honesty as the highest virtue’).

In all other examples, the verb does not impose selectional restrictions on  $NP_2$ . In (69c),  $NP_2$  can even be ‘weather’-*it*, as in *hij vond het koud* ‘he considered it cold’.

There are further differences between these constructions. One set of differences concerns the issue whether the AP can be replaced by the wh-word *hoe* (which, being a wh-word, will be preposed). This is possible for (69c,d), but impossible for (69a,b):<sup>26</sup>

<sup>26</sup>Example (72a) is well-formed under an irrelevant reading, which will be ignored, as before.

- (72) a \*Hoe kocht hij de winkel? (answer: leeg)  
       ‘How did he buy the shop?’ (answer: empty)  
       b \*Hoe kneep hij de tube? (answer: leeg)  
       ‘How did he squeeze the tube?’ (answer: empty)  
       c Hoe vond hij haar? (answer: aardig)  
       ‘How did he find her?’ (answer: nice)  
       d Hoe verfde hij de deur? (answer: groen)  
       ‘How did he paint the door?’ (answer: green)

I cannot explain why these facts are the way they are, but they clearly indicate that the examples (a,b) must be analyzed in a way which differs from the analysis for the examples (c,d).

There has been some debate regarding how these constructions should be analyzed. Basically, two points of view can be distinguished.<sup>27</sup> One view, adopted e.g. by [Hoekstra, 1984], [Hoekstra, 1988] and Van Gestel ([van Gestel, 1989]) claims that all these constructions involve small clauses with overt subjects. The second point of view, adopted e.g. by [Neeleman and Weerman, 1992] claims that small clauses play no role in the derivation of these sentences at all.

The position adopted here regarding these constructions is intermediate, in two respects. First, I claim that small clauses do play a role in these constructions, but I do not claim that all small clauses are present at all syntactic levels of representation (some are, others aren't). Second, I claim that the NP complement of verbs such as *verven* (and similarly *zetten*, *leggen*, *stoppen*, *doen* (all meaning ‘to put’), and *staan*, *liggen*, *zitten* (all meaning ‘be (located)’)) is normally<sup>28</sup> not the subject of the small clause but an argument of the verb. The analysis of these constructions differs from the analysis of examples such as (69a,b,c) where the NP is the subject of the small clause. [Hoekstra, 1984] and [Hoekstra, 1988] analyze all these constructions in the same manner, with the NP as the subject of the small clause. The evidence for this analysis, however, is rather weak for examples with the verbs mentioned above. The only real argument Hoekstra presents concerns the interpretation of the NP complements in such constructions. According to Hoekstra, these NP complements are obligatorily interpreted as *affected* complements, but not as *effected* complements (cf. *He painted the house* (affected or effected reading) v. *He painted the house green* (affected complement reading only). Though this may be correct, it is unclear what the point shows: it can count as an argument only if it can be shown additionally that this difference in meaning can only be accounted for by adopting the small clause structure proposed by Hoekstra, and not by adopting some other structure. But Hoekstra does not show this, so no conclusions can be drawn on the basis of this fact.

<sup>27</sup>I will restrict attention here to analyses of these constructions dealing with Dutch facts. Of course, there is a much wider literature dealing with these or similar constructions in other languages.

<sup>28</sup>In many cases, an analysis like the one proposed by Hoekstra is possible as well, e.g. Hoekstra gives the example *I have painted my fingers black and blue*, which by all criteria mentioned above must be analyzed as taking a small clause argument.

## 4.7 Systematic Relations between Predicates

Our experiences with the system for predicate-argument relations as discussed in the preceding sections in grammars for a machine translation system, and especially our experiences with filling a formal lexicon with this information yielded as one result that it is necessary to capture regular relationships between lexical items by means of rules. This has been done to an insufficient degree in the current grammar, and requires amelioration. This section describes a number of examples where such regular relationships hold, and a method is suggested to express them.

One possible type of rules which might serve this purpose are lexical rules. Lexical rules operate in the lexicon and take as input a lexical entry and yield as output a new lexical entry.

In a grammar in which lexical rules play a role, a lexical item  $L$  belongs to the set of lexical items of a language, if  $L$  belongs to the set of initial lexical entries (which are defined by enumeration), or if there is a lexical rule by which  $L$  can be related to a different lexical item  $L'$  which belongs to the set of lexical items of the language. If the lexical rules only define a finite set of related lexical entries, the lexical rules can be used ‘off-line’, e.g. the lexical rules can be used to expand a set of initial lexical entries to the full set of lexical entries. This full set is actually used during processing. If the lexical rules are used ‘on-line’, i.e. they are applied during processing, this is to a large extent equivalent to formulating these rules as rules operating in syntax. In this case, if recursion is allowed, the lexical rules must obey a measure condition in order to guarantee termination in analysis (the output of the application of a lexical rule must be ‘smaller’, according to some measure, than the input).

A difference between formulating these rules as lexical rules and formulating them as syntactic rules in the grammatical framework adopted here, is that syntactic rules should be associated to a meaning and to a translation-equivalent operation in the grammars of other languages. Wherever possible, it would be most natural within the framework adopted here to deal with such phenomena by means of syntactic rules which have translation-equivalent rules in the grammars of other languages.<sup>29</sup> But this is probably not possible for all examples, since it is often not the case that translation-equivalent rules are applicable to translation-equivalent basic expressions. In that case, either lexical rules which apply in the lexicon are necessary, or syntactic rules which apply to larger structures than lexical items are required. An example of the latter case might be relevant to middle formation. Both English and Dutch have a rule of middle formation. These rules can be associated to the same meaning and be translation-equivalent. But the rules do not apply to translation-equivalent lexical items (e.g. the English rule does apply to *bribe*, but the Dutch rule does not apply to *omkopen*). In this

<sup>29</sup>In the controlled M-grammar framework both ‘syntactic’ rules and ‘lexical’ rules can apply in the syntactic component. If both types of rules apply in syntax, the distinction between the two types of rules is made by imposing different restrictions on their applicability. ‘Lexical’ rules can apply only to S-trees corresponding to a word (lexical S-trees) and yield S-trees which correspond to a word (lexical S-trees), while ‘syntactic’ rules can apply to larger structures, and yield larger structures.

case one might consider setting up two translation-equivalent rules which apply to larger structures and which form normal middles whenever possible, but turn to different constructions if a middle cannot be formed. To give an example, for Dutch *omkopen* the rule might yield a construction such as *zich laten omkopen*.

A problem with lexical rules applying in the lexicon concerns their status in multilingual systems. If we derive some lexical entry L2 from some lexical entry L1 by lexical rule R1, then how should L2 be translated into another language? This is taken care of automatically if the lexical rule applies in syntax, but not when it applies in the lexicon. For lexical rules which apply in the lexicon one should not only specify for each lexical entry which lexical rules are applicable to it, but also what the translation equivalents of the derived entries are in the lexicons of other languages. This can be achieved by specifying a unique name for the derived entries. Multilingual dictionaries can then be formulated in terms of relations between unique names for lexical entries from different languages. A consequence is that lexical rules applying in the lexicon can derive only a finite number of non-basic lexical entries.

Passivization is often mentioned as a possible lexical rule. The idea is that such a rule would change the subcategorization properties and additional properties of a lexical item, and yield a new lexical item which is the passive of this lexical item. However, it is dubious whether passivization can be considered a lexical rule, for a number of reasons. For discussion, see chapter 5.2.

Examples where lexical rules might be useful in dealing with variations in argument structure are listed below. Most examples are from Dutch, though some apply to English and Spanish as well, and examples from other languages can readily be found.

1. The regular alternation in Dutch in ‘verbs of manner of movement’ such as *lopen* ‘walk’, *schaatsen* ‘skate’, *rijden* ‘drive’, etc., which can be used either as an activity verb (*hij heeft geschaatst* ‘he has skated’), or as a directional verb (*hij is naar Franeker geschaatst* ‘he has skated to Franeker’), or as a transitive verb (*hij schaatste de elfstedentocht* ‘he skated the Elfstedentocht’) or as a verb with a directional small clause (*hij heeft de winnaar naar de finish geschaatst* ‘he skated the winner to the finish’). The forms of all these verbs do not differ, but all kinds of syntactic and semantic properties (number of arguments, event type, choice of auxiliary verb, passivizability, etc) are different for these different uses. It is not possible to describe these differences within one lexical entry.
2. A class of verbs which express activities having to do with body care (*kammen* ‘to comb’, *wassen* ‘to wash’, *scheren* ‘to shave’, *afdrogen* ‘to dry’) can be used transitively, but also reflexively with the reflexive pronoun *zich*. This differs from the use in which the verb has a reflexive pronoun as its direct object (in such a case the reflexive pronoun *zichzelf* should be used), and it is more a special kind of intransitivization.
3. A class of transitive verbs can also be used with a resultative small clause

instead of the direct object, e.g. *hij kocht een appel* ‘he bought an apple’ v. *hij kocht de winkel leeg* ‘he bought the shop empty’ (i.e. he bought so much that the shop got empty). These uses cannot be considered as two uses of the same lexical entry, as discussed above.

4. The opposition transitive v. middle (*hij verkocht het boek* ‘he sold the book’ v. *de boeken verkopen goed* ‘the books sell well’), and perhaps also intransitive v. middle (*hij fietst* ‘he bicycles’ v. *deze weg/fiets/band fietst lekker* ‘this road/bicycle/tyre bicycles nicely’).
5. The opposition result-event meaning with nominalizations. Perhaps it is possible to simply indicate for certain words that they are a nominalization of a certain verb. The noun *destruction* can be used to illustrate this. The noun can be used as the result of a ‘destroy’ event, but also to describe the event itself. In this latter usage it takes arguments in a way similar to the related verb *destroy*. If it is indicated that *destruction* is a nominalization of the verb *destroy*, then its complementation properties can perhaps be derived by rule. However, this clearly requires a more thorough investigation.
6. Oppositions such as *hij kneep* ‘he squeezed’ v. *hij kneep de tandpasta uit de tube* ‘he squeezed the tooth paste out of the tube’ en *hij danste* ‘he danced’ v. *hij danste zijn longen leeg* ‘he danced his lungs empty’, if these cannot be accounted for by general rule.
7. Certain intransitive-transitive alternations, e.g. *hij liep* ‘he ran’ v. *hij liep de marathon* ‘he ran the marathon’; *hij danste* ‘he danced’ v. *hij danste de tango* ‘he danced the tango’.

Lexical rules might also be useful to express systematic relations between lexical entries which do not directly relate to predicate-argument relations, e.g. for the following:

1. the opposition mass-count, to indicate a quantity, e.g. *beer tastes good* v. *please give me two beers*)
2. the opposition mass-count, to indicate a kind, e.g. *wine tastes good* v. *several nice wines*.
3. the opposition count-mass, e.g. certain words for animals can be used as a mass noun if it denotes food (e.g. *Er liep \*(een) kip* ‘A chicken was running’ v. *hij at kip* ‘he was eating chicken’).

The list given above is just a short list of potential examples where lexical rules might be useful. Many other examples can easily be found. Most, however, require additional investigation to see whether an approach in which lexical rules are used is really fruitful.



## 4.8 Concluding Remarks

In this chapter I discussed how predicate-argument relations are dealt with in the Rosetta grammars. I showed that the compositional nature of controlled M-grammar has immediate and far-reaching consequences for the treatment of predicate-argument relations: the very design of the grammar immediately derives a variant of the  $\theta$ -criterion in a very strict way; it implies that only the notion ‘semantic argument’ is relevant, and it has immediate consequences for the treatment of optional arguments.

To avoid the undesirable consequences of this approach, it has been proposed that a class of phrases should not be dealt with as arguments, but as *bound adjuncts*, which are intermediate between arguments and adverbials.

The treatment of small clauses, which play an important role throughout the grammar, has been illustrated.

Finally, it has been observed that many systematic and regular relationships between lexical entries or between different uses of one and the same entry are not expressed adequately, and it has been suggested that rules should be introduced to express these. One can partially do this by means of local rules which apply in the isomorphic grammars, or by syntactic rules which operate on larger structures, or, to capture monolingual generalizations only, by lexical rules which operate in the lexicon. Such lexical rules, however, can only apply to lexical entries which specify a unique name for an output lexical entry, in order to be able to express the correct translation relations.



## Chapter 5

# Some Constructions

### 5.1 Introduction

In this chapter I will illustrate how a number of constructions have been incorporated in the Rosetta grammars. I will discuss the treatment of passivization (section 5.2), Verb Second in Dutch (section 5.3), unbounded dependencies (section 5.4) and crossing dependencies in Dutch (section 5.5). For all these constructions a particular existing analysis has been adopted as a starting point, and an attempt has been made to incorporate these analyses in the framework adopted here. This required adaptation of the original analyses in certain cases, as will be illustrated, though in general the existing analyses could be incorporated fairly directly.

In the analysis of passives I attempted to incorporate some of the insights developed in the *Move  $\alpha$*  program with respect to passivization into a framework in which there are language-specific and construction-specific rules. The major tenet of the *Move  $\alpha$*  program is that it claims there are no language-specific or construction-specific rules. Though this approach is incompatible with the framework adopted here, it is clear that much can be gained by making rules less construction-specific.

The result is an analysis of passives in which a conglomerate of rules, only one of which is particular to passives, derives passive constructions. Most of the rules postulated are also used in the derivation of completely different constructions, and are for this reason not construction-specific. The operations which are common to the derivation of several constructions have been identified and isolated (factored out) in separate rules, so that undesirable duplication has been avoided.

The analysis of Verb Second in Dutch incorporates the analysis developed originally by Den Besten ([den Besten, 1983]) in a fairly direct manner. It is argued that this analysis, which requires a movement rule, is superior to alternative analyses developed in the GPSG and HPSG frameworks. The analysis of Verb Second as developed by Den Besten states that the application of the relevant rule is conditioned by syntactic structure: no other factors are relevant and there is no direct

link with semantic properties. Though the position of the finite verb in a sentence may appear to have a direct influence on the interpretation of a sentence (cf. *hij loopt* ‘he is walking’ vs. *Loopt hij?* ‘Is he walking’). Also see chapter 3), the rule accounting for the position of the verb has been formulated as a transformation which is sensitive to syntactic structure only. It plays a role in the derivation of many different kinds of sentences (declaratives, wh-questions, yes-no questions, conditional subordinate clauses, etc.) and is therefore not construction-specific. The fact that Den Besten’s analysis can be incorporated directly provides an illustration of one of the main virtues of the controlled M-grammar framework: the transformations allowed in this framework make it possible to factor out common properties in separate rules.

For unbounded dependencies, an analysis is adopted in which phrases are actually moved from one position in the tree to another position in a successive cyclic manner, as originally proposed in [Chomsky, 1973]. It is argued that such a movement analysis is superior to analyses using feature-transportation mechanisms to describe unbounded dependencies. The evidence is derived from the behavior of idiom chunks and from facts relating to inversion in Spanish, as originally pointed out by [Torrego, 1984].

The movement rule involved in describing these unbounded dependencies is also a transformation which is sensitive to syntactic structure only. It is not associated to a specific meaning aspect. This makes it possible to use this rule in all constructions where unbounded dependencies play a role, so this rule is also not construction-specific.

Finally, I will discuss the treatment of crossing dependencies in Dutch. I adopted the analysis by [Evers, 1975] as a starting point, and incorporated it, with a few minor modifications, into the grammar. I will show that the controlled M-grammar formalism makes it possible to incorporate such a well-motivated syntactic analysis fairly directly, and to express the relevant syntactic generalizations concerning the distribution of verbs in Dutch with a minimum of machinery.

## 5.2 Passivization

There are basically two major types of analyses of the passive construction in modern formal theories of grammar. They are called the *syntactic* and the *lexical* analysis of passive, respectively.<sup>1</sup> If we restrict attention to verbs, then it can be said that both theories assume that an essential part of passivization is an operation on the verb. The lexical analysis states that passivization is to be characterized as an operation which changes the subcategorization frame of a verb: the prototypical case of a transitive verb takes a direct-object NP in active structures, but it doesn’t take one in passive structures. In addition, the subject is turned into

---

<sup>1</sup>A third type of analysis, which might be called the *metasyntactic* analysis of passive, is employed in GPSG (see [Gazdar et al., 1985]). Passive structures are formed by syntactic rules which have been formed by metarules from other syntactic rules. As far as we can see, most of the drawbacks of the lexical analysis apply to the metasyntactic analysis as well.

a *by*-phrase. The syntactic analysis states that passivization does not change the subcategorization frame: a passivized transitive verb takes a direct object, just as its active counterpart does; passivization only suppresses the subject, which can (optionally) be realized in a *by*-phrase. One essential difference between the two classes of theories can be summarized as follows: under the lexical analyses it is assumed that a passivized verb does not take a direct object at syntactic levels of representation, though the syntactic analyses claim that passivized verbs do take a direct object at (at least some) syntactic levels of representation. Another essential difference between the two approaches can be described as follows: in the syntactic approach it is assumed that many operations involved in the formation of passive structures are also involved in the formation of completely different structures, so that these operations should not be part of the operation to form passives. Examples of these operations will be given below. In the lexical analyses, however, some of these operations are part of the lexical rule to form passivized verbs.

Representatives of lexical analyses of passivization are Lexical Functional Grammar ([Kaplan and Bresnan, 1982]), Head Phrase Structure Grammar ([Pollard and Sag, 1987]), and certain others. Representatives of syntactic analyses of passivization are the Principles and Parameters Framework<sup>2</sup> (which originated with the Government Binding Theory), see [Chomsky, 1981], [Chomsky and Lasnik, 1991], and Relational Grammar ([Perlmutter, 1977]).

It must be clarified what a ‘syntactic’ analysis vs. a ‘lexical’ analysis of passive means in the context of the controlled M-grammar framework. In ‘lexical’ analyses of passivization it is usually assumed that the rule to form passives applies in the lexicon. Such an analysis cannot be incorporated directly in the controlled M-grammar framework which does not have rules in the lexicon. It might be possible to use rules of the kind described in section 4.7, but these do not yet exist. As noted before, in the controlled M-grammar framework both ‘syntactic’ rules and ‘lexical’ rules must apply in the syntactic component. The distinction between the two types of rules is made by imposing different restrictions on their applicability. ‘Lexical’ rules can apply only to S-trees corresponding to a word (lexical S-trees) and yield S-trees which correspond to a word, while ‘syntactic’ rules can apply to larger structures, and yield larger structures.

In Rosetta, the syntactic analysis of passivization of Government Binding (GB) theory has been taken as a basis for the analysis. The main reason for this is that the lexical analysis combines a number of operations which in my view are independent and should be described as independent operations. This will become clear below, where I will describe and justify the analysis adopted in Rosetta. Many of the arguments for certain assumptions have been derived directly from GB-like analyses of passive constructions.

The analysis of passive in the Government Binding Theory is an excellent example to illustrate what the *Move  $\alpha$*  program has led to. I have attempted to incorporate some of the positive results of this approach into the controlled M-

---

<sup>2</sup>A lexical analysis for passives is not excluded in this framework, but most researchers working in this paradigm assume a syntactic analysis.

grammar framework, especially with respect to the fact that most rules involved in forming passive structures are not specific to the passive construction, but have independent motivation and are used to form completely different constructions as well.

I will first outline globally how passive structures are derived in the grammar, and then explain in more detail why this analysis has been adopted.

The first step to form a sentence consists of combining a verb with a number of variables into a propositional structure consisting of a predicate and optionally a subject, by start rules (see chapter 2). The position of the variables and the grammatical relations they bear in this structure correspond to their positions and relations when used in an active sentence. Thus, for instance, the two-place verb *kiss* can be combined with two variables,  $x_1$  and  $x_2$ , where  $x_1$  is the subject and  $x_2$  is the direct object in the propositional structure headed by *kiss*.

Voice rules apply to determine the voice of the structure. The rule for active voice leaves the structure unchanged, and only changes the value of the voice attribute of the proposition node from *unspecified* to *active*.

The rule for passive voice (called *Rpassive*) changes the value of this voice attribute into *passive*, but a number of other changes are also performed. The rule *Rpassive*, however, is not really comparable to the traditional *passive transformation* (e.g. as in [Chomsky, 1957]). The traditional passive transformation has been decomposed into a number of separate rules and transformations, which together form passive structures. This approach is inspired directly by the treatment of passives in [Chomsky, 1981].

In order to illustrate this, I will summarize the differences between a typical passive sentence and the corresponding active sentence, and indicate where and how the differences among them are accounted for:

**Voice** The attribute **voice** has the value *active* in active sentences, the value *passive* in passive sentences. As pointed out above, the value of this attribute is set in the voice rules called *Ractive*, *Rpassive*.

**Subject** The first argument of a verb in active sentences is expressed as a *subject*. In passive sentences it is expressed as a complement to the preposition *door* in Dutch, *by* in English or *por* in Spanish. This change is performed in a separate rule.

**Verb Form** In active sentences the form of the main verb is not fixed. It can be a finite verb, an infinitive, a present participle, a past participle (if an auxiliary verb *hebben* or *zijn* (in Dutch), *have* (in English) or *haber* (in Spanish) is present); in passive structures the main verb is a past participle. The form of the verb is set in the rule *Rpassive*.

**Auxiliaries** In active sentences no auxiliaries to express voice are present. In passive sentences, however, auxiliary verbs (Dutch: *worden* or *zijn*; English: *be*; Spanish: *ser*) are present. These auxiliary verbs are introduced by rules turning the propositional structure into a clause.

**NP-movement** One of the NP-arguments inside VP in an active sentence is usually realized as a subject in the corresponding passive sentence. I will call this NP the *moving NP*. It can be a *direct object*, (**The girl** *was kissed*), an *indirect object* (**The man** *was given a book*), a *prepositional object* (**The girl** *was looked at*), the subject of an embedded clause (**The man** *was believed to be ill*) or the subject of an embedded small clause (**He** *was considered a fool*). The operation to turn the moving NP into a subject will be called *NP-movement*. NP-movement is performed by a special transformation.

As is clear, the traditional effects of the passive transformation are performed by the interaction of voice rules, rules to form *by*-phrases (or their equivalents in other languages), rules to form clauses and transformations to turn NPs into subjects. The reasons for doing this in this way are:

**by-phrase** This phrase is not formed in the rule *Rpassive*, because it is used in other constructions as well. Dutch has a construction which can occur only with the verb *laten* ‘to let’ in which the Dutch equivalent of a *by*-phrase (formed with the preposition *door*) occurs in a structure which is active in all other respects: *hij liet het huis door hen schoonmaken* lit. *he let the house clean by them* ‘he had them clean the house’. As in passives, this phrase is optional (cf. *Hij liet het huis schoonmaken*). *By*-phrases also occur in nominal constructions, and in so-called ‘modal passives’, a construction in which *te* + infinitival V has the meaning ‘can/must be V-ed’ (e.g. *de deur ons schoon te maken huizen* lit. *the by us clean to make houses* ‘the houses which are to be cleaned by us’), though we did not deal with these systematically in the Rosetta system. It is clear then, that *by*-phrases are not specific to passive structures, and that their formation should not be carried out in the rule *Rpassive*.<sup>3</sup>

**NP-movement** This operation is not performed inside the rule *Rpassive* because it is not essential to the passive voice at all:

- There are passive sentences with no VP-internal NP at all (impersonal passives, for instance *Er werd gedanst* (lit. *There was danced*) in Dutch, and passives of verbs that take a sentential complement, for instance *It*

---

<sup>3</sup>In many analyses, the *by*-phrase is not analyzed as an argument, but as an adjunct. Such an analysis is possible in the framework assumed here as well. However, for translation purposes we want to derive active and passive constructions isomorphically. The adjunct analysis can then only be maintained if the adjunct is introduced syncategorematically. In addition, in each analysis it must be guaranteed that the NP in the *by*-phrase is linked to the external (implicit) argument of the verb. The approach adopted takes care of this immediately, and does not require special linking mechanisms.

*was said that he was dishonest*). If NP-movement were part of the rule *Rpassive*, a separate rule for such passives would have to be stated.

- In personal passive sentences the moving NP sometimes remains inside VP, and is not moved into the subject position. Some relevant examples are: *De mannen werd **het boek** gegeven* (lit. *The men was the book given*), *Er werd door hem **een boek** gelezen* (*There was by him a book read*), for which it can be argued convincingly that the boldface NPs are in object position, not in subject position (cf. Den Besten([den Besten, 1981])).

- A distinction is made between *ergative* and *non-ergative* verbs (see chapter 4)

Ergative verbs are verbs that take an NP as argument but do not realize any of its arguments as an *external argument*. For these verbs there must be a rule to turn an NP into a subject as well, though this involves only active structures. It is natural to use the same rule both in passives and in ergative structures.

- NP-movement must also be performed for NPs that are not yet present in the structure at the moment that the rule *Rpassive* applies. As pointed out above, the first rules to form sentences combine verbs with a number of variables, and it is the case that for a sentence such as *She was considered smart by us* the structure contains a variable for the subject (*we*), and a variable for the ‘small clause’ *she smart*, but nothing that corresponds directly to the NP *she* (which is the NP which should become the subject) at the moment that the rule *Rpassive* must apply.
- Control transformations dealing with *obligatory control* (in the sense of [Williams, 1980]) can be simplified if the object is still an object at the moment that they apply. This makes it easier to account for the fact that a derived subject can still function as a controller, though an NP in the *by*-phrase cannot. Recall here that there are no traces. Some relevant sentences to illustrate the phenomenon are: *She was forced by him to defend herself* vs. \**She was promised by him to defend himself*. The structure of these sentences, at the point in the derivation where control transformations are applicable, can be represented informally as in (73a):

- (73) a [ was forced x2 [by x1] [ x2 to defend herself]]  
 b [ was promised x2 [by x1] [x1 to defend himself]]

In (73a) lexical properties of *force* state that the direct-object variable (x2) must be the controller of the subject of the infinitive. Since there is a direct-object variable, and since it has the same index as the subject variable of the infinitival complement, the control rule is applicable. In (73b), however, lexical properties of the verb *promise* state that the subject variable must be the controller of the subject of the infinitive. Since



there is no subject in the structure, the control rule is not applicable and the derivation blocks.

This illustrates that the analysis, in which NP-movement is factored out of the passive rule, makes it possible to apply theories of control which make a distinction between obligatory and non-obligatory control, in a description in which no traces are used.<sup>4</sup>

**Auxiliary Verb** The auxiliary verb is not introduced in the rule *Rpassive* because there are passive structures where no such auxiliary occurs. This is the case in small clauses with a verb as their head, e.g. *Hij kreeg het boek afgeleverd* (*He had the book delivered*), *Hij wist zich door hem gesteund* (*He knew himself by him supported*), *They had him killed* .

So, the effects of passivization are achieved by the interaction of a number of rules, some of which are necessary on independent grounds. The analysis is very much in the spirit of analyses within the *Move  $\alpha$*  program. The following properties of these analyses have been incorporated: (1) passivization is *syntactic*, not lexical; (2) no reference to a *by*-phrase is made inside a rule specific to passive structures; (3) no reference to a direct-object NP or any NP whatsoever is made; (4) movement of NP is necessary in certain cases, but (5) this movement is independent of passivization; (6) the presence of the auxiliary is independent of passivization.

In the *Move  $\alpha$*  program, movement of NP is usually obligatory because the passive participle does not assign case, and each NP must have case. This approach has not been adopted. In the Rosetta grammars passive participles and ergative verbs can assign (nominative) case to direct objects. Thus, absence of case does not force movement. Instead, in such configurations no subject is present, but a subject must be present. There are several ways to get a subject into such structures. One way is moving a VP-internal NP into the subject position. Another way is inserting an expletive pronoun. Expletive pronouns impose restrictions on VP-internal arguments. These restrictions account for the fact that not just any NP can remain VP-internal in all passive structures (cf. *\*er werd het boek gekocht* ‘there was bought the book’). The relation between expletives and VP-internal arguments must be made anyway, and not only when there are VP-internal NPs, and this approach avoids complex case-transmission mechanisms needed in other approaches.

---

<sup>4</sup>Henk van Riemsdijk suggests that the following contrasts weaken the argument given:

- (i) a I asked her to leave / to be allowed to leave  
 b She was asked to leave/ \* to be allowed to leave

These contrasts, however, only show that the conditions under which complement control or subject control is appropriate cannot simply be stipulated, but must be derived from deeper properties, perhaps along the lines of [Růžička, 1983]. The argument given in the text would only be weakened if the position of *she* in (ib) were the crucial factor to account for the contrast. But there is no reason to assume this. In the analysis of [Růžička, 1983], the position of controllers plays no role whatsoever.

Note that the lexicon contains no passive verbs (passive participles) as entries. Each verb is in the lexicon in its base form. The passive forms of a verb are created in syntax. Even verbs which can occur in passive only (e.g. Dutch *achten* ‘expect’, English *rumour*) are in the lexicon in their base forms, though their inherent properties specify that they can occur in passive voice only. The analysis deviates considerably from analyses of the passive construction in most other computational frameworks (LFG, HPSG), in which the passive construction is usually accounted for by a lexical rule that creates a new lexical entry with new complementation properties from an existing lexical entry. In these latter analyses, many of the facts taken into account here, are usually not considered at all (e.g. ergative verbs, arguments of passive verbs that remain in direct object position) or dealt with in a completely different manner (e.g. small clauses). The operations which turn the subject into a *by*-phrase, or a complement NP into the subject are not factored out as separate operations in these analyses. I therefore judge these lexical analyses as inferior to the syntactic analyses.

### 5.3 Verb Second in Dutch

In Dutch, the finite verb can occupy two different positions in a clause. It can occur in a relatively final position in certain constructions, but must occur in a relatively initial position in other constructions. This is illustrated in (74):

- (74) a ...omdat hij de jongen een boek *gaf*  
       ...because he the boy a book gave  
       ‘...because he gave the boy a book’  
       b hij *gaf* de jongen een boek  
       he gave the boy a book  
       ‘he gave the boy a book’

In (74a) the finite verb *gaf* follows the indirect object and the direct object, but in (74b) it precedes these phrases. I will refer to the fact that the verb can be in a relatively initial (the ‘second’) position as the *Verb-Second phenomenon* (abbreviated as the *V2-phenomenon*), a misleading (but generally accepted) term, as we will see below.

In Rosetta, an analysis has been adopted which is based directly on the analysis of the V2-phenomenon as sketched by [Bierwisch, 1963], [Koster, 1975], Den Besten ([den Besten, 1983]), and others. Though there are several variants of this analysis, the essential features of the analysis are quite generally accepted. Only the analysis by [Travis, 1984] modifies the standard analysis slightly, for sentences containing a sentence-initial subject. She justifies the distinction by pointing out a number of subject-non-subject asymmetries with regard to topicalization. The facts pointed out by her are well-known at least since Koster ([Koster, 1978a]). Though Travis’s analysis constitutes a real alternative, I did not adopt her analysis, because ad-hoc stipulations are required to prevent verb movement in clauses

containing a subordinate conjunction and because there are many alternative accounts for the subject-non-subject asymmetries which do not have such drastic consequences for the structure of the sentence. In the actual system a special rule has been assumed for the topicalization of elements in the subject position. It differs from the normal rule of topicalization in that it can prepose meaningless elements such as idiom chunks, expletives and clitic pronouns.

A completely different analysis for the V2-phenomenon in German is presented by Hans Uszkoreit ([Uszkoreit, 1982],[Uszkoreit, 1987]), in a totally different framework, viz. GPSG. Uszkoreit's analysis can be described briefly as follows. It is assumed that finite verbs can have any of the two possible values (+ or -) for the attribute MC.<sup>5</sup> The order of constituents, hence also the order of finite verbs, is determined by Linear Precedence Rules. These Linear Precedence Rules state that a constituent with the specification +MC precedes all other material in a sentence. A [-MC] constituent, however, follows all material (if one abstracts away from extraposition phenomena). A separate rule states that a sentence introduced by a subordinate conjunction can only be combined with a sentence containing a [-MC] verb.

In my view, Uszkoreit's analysis is inferior to the Koster-Den Besten analysis. There are several arguments to support this point of view. The most important ones will be outlined below.

Uszkoreit must assume a completely ad-hoc feature, MC, to be able to account for the facts. Verbs in the V2-position cannot be distinguished from verbs in a clause-final position except by their position. The fact that a feature is required to distinguish verbs that must invert and verbs that must not invert, is a clear indication that a 'trick' is applied in the framework to get the facts correct and that in fact the formalism assumed might not be suited to deal with the V2-phenomenon. The major objection is that an ad-hoc feature is proposed for verbs which do not differ in any way other than their position in the sentence: this is a very unnatural way to deal with such facts, and it indicates that the wrong formal means are used to describe such facts. This is acknowledged as a problem in [Uszkoreit, 1982, 147].

In Uszkoreit's analysis the complementary distribution between subordinate conjunctions and V2 is not accounted for, but derived indirectly by means of a stipulation. Of course, Den Besten's arguments to show this complementary distribution were set up for Dutch, but it is quite evident that these arguments might be relevant for German as well, so that they should at least be explicitly discussed.<sup>6</sup> Uszkoreit, however, does not discuss it in any way: the claimed complementary distribution does not follow from Uszkoreit's analysis without a separate stipulation to this effect.

Uszkoreit argues that a separable prefix and a verb when adjacent do not form a word. He is forced to such an analysis, because otherwise his account of V2 will

---

<sup>5</sup>This is an acronym for *main clause*, but this name should not be taken too seriously, as pointed out by Uszkoreit, since the V2-phenomenon also occurs in certain subordinate clauses.

<sup>6</sup>Henk van Riemsdijk pointed out to me that Den Besten's arguments can be replicated for German.

not work properly: it would predict that the verb must occur in the V2-position accompanied by its separable prefix, which is incorrect (cf. *\*Er anruft mich* v. *Er ruft mich an*).

Even if we accept the hypothesis that such separable prefix-verb combinations do not form a word, Uszkoreit's analysis imposes a stronger requirement: they cannot even be a constituent. It is unclear whether this can be maintained. In Dutch, such combinations must be a constituent in certain configurations. The particle and the verb can be put to the right of a verb raiser together, as in (75):

- (75) ...omdat hij de man wilde *opbellen*  
 ...because he the man wanted up call  
 '...because he wanted to call up the man'

This shows that the sequence particle-verb must sometimes be a constituent. It is unclear whether Uszkoreit's analysis is compatible with a proper analysis of such separable prefixes. However, we will not discuss this anymore, since the syntax of separable prefixes is a very complicated issue for which no uniformly accepted analysis exists.<sup>7</sup>

Uszkoreit assumes that German VPs are left-branching, i.e. he assumes that the verbs inside a VP are the rightmost heads in VPs and that infinitival complements (VPs, in his analysis) precede the verb. Though he briefly brings up the subject of how to correctly analyze auxiliary verbs, he does not in any way discuss the arguments Evers ([Evers, 1975]) adduced to argue that the verbs in the kind of structures considered by Uszkoreit form verbal clusters (both in German and in Dutch). If the structures postulated by Evers are correct, Uszkoreit's account of the V2-phenomenon does not work correctly. A discussion of Evers's analysis would thus be appropriate.

Uszkoreit is forced to assume that a finite sentence has a completely flat structure, i.e. there is no VP, and no further internal structure to a VP. There is evidence, however, from topicalization facts, that German does have a VP, which is in fact a structure in which each NP argument is introduced by binary branching. (see Den Besten and Webelhuth ([den Besten and Webelhuth, 1987]), [Netter, 1992] and [Engelkamp et al., 1992]).

Finally, there are — to my knowledge — no facts that can be described in a more adequate manner than in the standard analysis. I thus rejected Uszkoreit's analysis as a basis for Rosetta.<sup>8</sup>

I will now illustrate the relevant facts in more detail and outline the analysis adopted in Rosetta. As pointed out before, the position of the finite verb can vary in a clause between a final position and a more initial ('second') position. This can be illustrated with the following examples:

<sup>7</sup>Henk van Riemsdijk suggested a 'lexoid' analysis for verbs with separable prefixes. See [Model, 1991b] for one possible elaboration of this proposal. For an alternative analysis, see [Neeleman and Weerman, 1992].

<sup>8</sup>In 1990, long after the implementation of V2 in Rosetta, Pollard ([Pollard, 1990]) proposed an analysis of V2-phenomena in German in the HPSG framework. His analysis resembles Uszkoreit's to a high degree, and is — as far as I can see — at best equivalent to Uszkoreit's analysis, but probably worse. As a reaction to this analysis, see also [Netter, 1992].

- (76) a \*hij het boek gekocht *heeft*  
 b hij *heeft* het boek gekocht ‘he has bought the book’
- (77) a ..dat hij het boek gekocht *heeft*  
 b \*..dat *heeft* hij het boek gekocht ‘that he has bought the book’

In the first sentence there are two verbs, a participle form (*gekocht*) and a finite form (*heeft*). In declarative main clauses such as in (76) the inflected verb must occur in second position. Hence (76a) is ill-formed and (76b) is well-formed. In subordinate clauses such as in (77) the inflected verb cannot occur in second position. Hence (77a) is well-formed, and (77b) is ill-formed. It is assumed that the verb is initially generated in a clause-final position, in accordance with the hypothesized SOV-character of Dutch (see [Koster, 1975]).

Though the phenomenon dealt with here is called *Verb Second*, the verb sometimes occurs in sentence-initial position (e.g. in main yes-no questions, *Heeft hij het boek gekocht?*), sometimes in sentence-second position (in normal main declarative sentences) and sometimes even in sentence-third position, e.g. in main declarative sentences with left-dislocated elements such as (78):

- (78) Dat boek, dat heeft hij niet gekocht  
 That book, that has he not bought  
 ‘That book, he hasn’t bought’

The correct characterization of the position of the inflected verb has been established by Den Besten ([den Besten, 1989, 24]):<sup>9</sup> the inflected verb is in the position of subordinate conjunctions.

This generalization has been incorporated directly into Rosetta by reserving only one position which must be used both by inflected verbs and by subordinate conjunctions. This single position is identified by the name of a grammatical relation, *conj* (see chapter 2 for the role of ‘grammatical relations’ as indicators for qualified positions).

It is often said (informally) that Verb Second occurs in main clauses, but not in subordinate clauses. In fact, however, there is no correlation at all between Verb Second and main clause status: there are main clauses and subordinate clauses where Verb Second occurs, and there are main clauses and subordinate clause where Verb Second does not occur, as illustrated in the following examples:

**Main, +V2** Hij is ziek geweest (He has ill been)

**Main, -V2** Ziek dat hij geweest is! (Ill that he been has)

**Subord, +V2** Is hij ziek, dan gaan we niet weg (Is he ill, then go we not away)

**Subord, -V2** Als hij ziek is, dan gaan we niet weg (If he ill is, then go we not away)

<sup>9</sup>Den Bestens paper was available in unpublished form since 1978.

In Rosetta, Verb Second is taken care of by transformations<sup>10</sup> (*Tv2*) that move the inflected verb into the special position reserved for inflected verbs and subordinate conjunctions.

Den Besten's analysis claims that the presence of a subordinate conjunction and V2 mutually exclude each other, i.e. if a subordinate conjunction occurs in the clause, there is no V2, and, conversely, if there is V2 in a clause, then no subordinate conjunction occurs. This leaves the possibility of clauses in which there is neither V2 nor a subordinate conjunction.

In the analysis in Rosetta a slightly stronger position was adopted by excluding the possibility of there being neither V2 nor a subordinate conjunction. Therefore, the transformations accounting for V2 must apply if no subordinate conjunction is present.

There are two constructions that require some additional clarification. In finite relative and subordinate interrogative clauses Verb Second cannot occur and no conjunction is present either. For these cases it is assumed that at the point of application of *Tv2*, conjunctions are present. They are deleted after the application of *Tv2*. For subordinate interrogatives this is the subordinate conjunction *of*, which can appear in surface trees but is usually deleted (cf. *Ik weet niet wie (of) het gedaan heeft lit. I know not who (if) it done has* 'I don't know who did it'); for relative clauses it is the subordinate conjunction *dat*, which is deleted obligatorily.

Thus we see that the stronger position adopted is actually not completely correct, and some auxiliary assumptions are required. One of these auxiliary assumptions, viz. the postulated presence of the subordinate conjunction *of* in subordinate interrogative clauses can be justified independently, but there is no independent evidence for the presence of conjunctions in relative clauses. Nevertheless, with this auxiliary assumption one can adopt the simplest possible analysis of V2 phenomena, so that it receives some indirect confirmation.

The analysis accounts for all the examples mentioned above, and also all other cases in Dutch. E.g. the method adopted is immediately able to deal with the fact that Verb Second occurs in clauses such as *Al is hij ziek, we gaan toch* (*though is he ill, we go anyway* 'Though he is ill, we'll go anyway') and *Hij mag dan ziek zijn, we gaan toch* (*he may be ill, we go anyway* 'He may be ill, we'll go anyway'). These constructions have not been implemented yet, but when they are, nothing special will have to be done to take care of the fact that Verb Second occurs in these clauses. That is automatically guaranteed by the adopted analysis of Verb Second.

The analysis outlined incorporates a number of crucial properties of Den Besten's original analysis, i.e. (1) the V2-position is identical to the position of subordinate conjunctions, and (2) there is a movement rule which moves the finite verb into this position. It is much simpler than the complicated gap-threading mechanism used by Van Noord *et al.* ([van Noord et al., 1990]) to account for the distribution of verbs in Dutch and Spanish (see also chapter 3 and Van Noord ([van Noord, 1991]), in which a different approach to these phenomena is proposed, which does not make

<sup>10</sup>Due to limitations of the M-rule notation Verb Second cannot be accounted for by a single transformation.

use of gap-threading techniques), and it is superior (with respect to descriptive adequacy) to the analysis presented by Uszkoreit ([Uszkoreit, 1982],[Uszkoreit, 1987]) in which linear precedence rules are used to account for comparable facts in German.

Note that none of the objections that apply to the other analyses hold for the analysis in Rosetta. Since the *V2*-phenomenon is accounted for by movement rules, it is relatively independent of the analysis of the rest of the sentence and the VP. The complementary distribution between subordinate conjunctions and finite verbs is dealt with in an adequate manner, no ad-hoc feature to distinguish finite verbs that are to be preposed from finite verbs that are not to be preposed is required, and verbs with separable prefixes can be dealt with even if the prefix and the verb form a constituent.<sup>11</sup>

In the analysis adopted, the movement rules to account for the *V2*-phenomenon are sensitive to syntactic structure only: the rules apply if the special position reserved for subordinate conjunctions and finite verbs is unoccupied. Though the *V2*-phenomenon plays a role in a wide variety of constructions (declarative sentences, yes-no interrogatives, conditional subordinate clauses, etc.), the analysis proposed makes it possible to factor out the *V2*-phenomenon in separate construction-independent rules.

## 5.4 Unbounded Dependencies

Unbounded dependencies, as found, for instance in WH-questioning and in relative and topicalisation constructions, are dealt with by transformations that move constituents to other positions. These movements are applied in a successive cyclic manner, as is customary in transformational grammar. To illustrate, a sentence such as

(79) What do you think that Peter would say that she had bought?

is derived in the following manner. First, the most embedded clause is derived. The WH-pronoun *what* is the direct object of the verb *buy*, and prior to the relevant movement rule the clause looks as follows: [*that she had bought what*]. The movement transformation preposes *what*, yielding: [*what that she had bought*]. In a second step, this clause is used as a complement to the verb *say* on the next cycle:

(80) [that Peter would say [what that she had bought]].

The movement transformation preposes *what* out of the embedded clause:

(81) [what that Peter would say [that she had bought]].

<sup>11</sup>Actually, in Rosetta it is assumed that the verb and the separable prefix form a word, and a special rule has been formulated to extract the verb from this word. This was decided because it was the only way to form derivative words such as *onafleidbaar* ‘underivable’, *onoplosbaar* ‘unsolvable’ uniformly.

Finally, this clause is used as a complement to the verb *think* in the next cycle:

- (82) you think [what that Peter would say [that she had bought]]

In this structure *what* is preposed again, yielding the relevant sentence (after the application of transformations dealing with subject-auxiliary inversion and *do*-support).

In all constructions where unbounded dependencies play a role, the same movement transformation can be used. WH-Questioning, relativization and topicalization have been implemented and some other constructions may be dealt with in this way (comparatives, tough-movement constructions, complements to the adverb *too*), see [Chomsky, 1977b].

The fact that unbounded dependencies are dealt with by means of transformations that perform actual movements, yields interesting results in interaction with certain idiomatic expressions (For the treatment of idiomatic expressions, see Schenk ([Schenk, 1986],[Schenk, 1989]) and [Landsbergen et al., 1989]). In particular, this makes it comparatively simple to deal with movement of idiom chunks. Consider a Dutch sentence such as

- (83) Welke man heeft zij in de maling genomen?  
Which man has she in the grind taken?

Assuming that the Dutch idiomatic expression *iemand in de maling nemen* must be translated into *to pull someone's leg*, this sentence should be translated into

- (84) Which man's leg did she pull?

in which a part of the idiomatic expression *to pull someone's leg* is in sentence-initial position. The system is able to yield this translation, and even the more complex example *Welke man denk jij dat zij in de maling genomen heeft?* is translated correctly into *Which man's leg do you think that she pulled?*, where one part of the idiom (*leg*) is in the main clause and the other part (*pull*) is in the subordinate clause.

In the Dutch sentence the idiomatic expression is in the subordinate clause as a whole, whereas in the English sentence the idiom is split over the main clause and the subordinate clause. This discrepancy can be dealt with without any special measures. In both languages a WH-phrase must be preposed. In Dutch this WH-phrase is simply a direct object, in English the WH-phrase is part of a larger NP containing a part of an idiom. The fact that the whole NP must be preposed, and not the WH-phrase on its own (cf. \**Which man's did she pull leg?*) is a consequence of the normal syntactic rules of English.

[Torrego, 1984] shows that inversion of verb and subject in Spanish is triggered by the presence of a certain class of WH-phrases in a clause-initial position. This is illustrated in (85). The relevant subjects have been italicized. Note that inversion must apply both in main and in subordinate clauses (see (85c,d)).



- (85) a Con quién vendrá *Juan* hoy?  
 (With whom will Juan come today?)  
 b \*Con quién *Juan* vendrá hoy?  
 c Es impredecible con quién vendrá *Juan* hoy  
 (It is impossible to predict with whom Juan will come today)  
 d \*Es impredecible con quién *Juan* vendrá hoy

These facts can be accounted for by assuming a rule which puts WH-phrases in a clause-initial position (WH-movement), and an inversion rule which inverts the subject and the finite verb when preceded by a WH-phrase from the relevant class. If the WH-movement rule applies before inversion, the relevant sentences can be derived immediately. Torrego shows that extraction of a WH-phrase out of a deeply embedded clause requires inversion in each intermediate clause and in the clause where the WH-phrase ends up. This is illustrated in (86).

- (86) a *Juan* pensaba que *Pedro* le había dicho que *la revista* había publicado ya el artículo  
 (Juan thought that Pedro had told him that the journal had published the article already)  
 b Qué pensaba *Juan* que le había dicho *Pedro* que había publicado *la revista*?  
 (What did Juan think that Pedro had told him that the journal had published?)  
 c \*Qué pensaba *Juan* que *Pedro* le había dicho que *la revista* había publicado?

Torrego accounts for these (and many other) facts by, crucially, adopting a successive cyclic application of WH-movement, in combination with a specific assumption relating to bounding nodes in Spanish. By assuming the simple rules which are required to account for inversion in simple sentences, in combination with successive cyclic application, the relevant facts follow directly from this analysis. If other approaches are adopted to deal with WH-movement, e.g. gap-threading techniques, or rules which relate WH-phrases to intermediate traces, the rule for inversion must be made more complex to take inversion in intermediate clauses into account. This is a very clear example in which the adoption of a transformation to account for WH-movement makes it possible to keep other rules (inversion) simple, while other techniques to account for WH-movement require that the inversion rule is complicated to account for all relevant facts. For these reasons an analysis which incorporates these transformations is descriptively more adequate than the alternative analyses.

Torrego's analysis has been implemented directly, and the relevant pattern of facts follow directly from it. The WH-phrase is moved step by step from one clause to the next, where it creates the environment for inversion, and after the application of inversion the process is repeated on the next cycle.

Adopting a movement transformation to deal with unbounded dependencies gives a simple and elegant account of very complex facts. In many frameworks the

role of movement rules has been taken over by rules relating a phrase and a gap (its trace). This is equivalent in many respects, provided that the relation involves the whole phrase and not just its top node (cf. the idiom chunks), and provided that there are intermediate gaps with the right properties in the right positions (cf. the Spanish inversion facts). In many (computationally-oriented) frameworks the relevant relation is implemented by techniques to pass features from one node to another (e.g. gap-threading techniques, [Pereira and Shieber, 1987]). However, using powerful operations such as movement rules generally yields simpler, more elegant and descriptively more adequate analyses. Using powerful M-rules which allow movement makes it possible to really factor out movement from all other rules. One can simply write one (set of) rule(s) dealing with this phenomenon, and let it interact with all the other rules. Furthermore, given such a (set of) rule(s), the formulation of certain other rules can be kept simple as well. Using feature-passing techniques, however, will require the addition of features to virtually every node occurring in a rule, and each node in the syntactic objects (trees, or directed acyclic graphs, etc.) will — so to speak — be ‘contaminated’ by the presence of features only required to pass on information concerning other parts of the structure. So, many other rules are actually made more complex. In addition, if rules are made sensitive to information contained in these nodes (e.g. to account for successive cyclicity effects), analyses using feature-passing techniques will predict that successive cyclicity effects will occur everywhere along the path to the antecedent of the gap. However, these effects have been observed only around the position for subordinate conjunctions.

In certain frameworks the treatment of long distance dependencies is factored out of other rules, though feature-passing techniques are used. An example is GPSG ([Gazdar et al., 1985]), where general principles govern the distribution of features over nodes. Similarly, in the MiMo machine translation system (Van Noord *et al.* ([van Noord et al., 1989])) separate rules account for the anaphoric relations holding for *wh*-phrases and their traces. Even though long distance dependencies have been properly factored out in these frameworks, they still suffer from other deficiencies. Neither of these frameworks can deal adequately with the facts of Spanish inversion. Van Noord *et al.* ([van Noord et al., 1989]) claim that they can deal adequately with Spanish inversion, but this is yet to be conclusively shown.

The movement analysis outlined above makes it possible to account for the facts of Spanish inversion in subordinate clauses by using the same rule which is required for inversion in main clauses. This is not the case in the analysis sketched by Van Noord *et al.* ([van Noord et al., 1989]). They characterize the anaphoric relations to account for long distance dependencies as relations holding between comp-positions. Their rule of inversion is sensitive to properties of the comp-position:

The Spanish synthesis component can check whether the comp-position of a clause is **either filled or bound**. If so, the clause is inverted.  
Van Noord *et al.* ([van Noord et al., 1989, 305]) [emphasis mine, JO]

Though this formulation might work, one should notice that the relevant condition for inversion contains a disjunction (as emphasized by me), in itself an indication that a generalization is missed. In addition, the second part of this disjunction is required exclusively to account for inversion in subordinate clauses which are not introduced by a *wh*-phrase. Clearly, the rule of inversion must be made more complex to deal with inversion in such subordinate clauses. This is not required under the movement analysis.

Furthermore, idiom chunks can only be dealt with properly if the full antecedent tree is related to the trace, but in these frameworks only the properties of the top node are.<sup>12</sup> I conclude that the movement analysis is superior to other analyses in which feature-passing techniques are used.

The feature-passing techniques are no more than very complicated techniques to partially simulate movement rules, and they usually do not have the properties mentioned above which are crucial for an adequate description of idioms and inversion in Spanish, and I know of no evidence which favors analyses using feature-passing techniques over movement analyses.

It is quite clear that the analysis given here makes it possible to account for unbounded dependencies independently of the construction in which they occur. The fact that the relevant rules are transformations makes it possible to factor the unbounded dependencies out of individual constructions, as desired.

## 5.5 Crossing Dependencies in Dutch

This section deals with the treatment of crossing dependencies in Dutch. Crossing dependencies are dependencies in which drawing lines between all the dependencies results in crossing lines. Such dependencies occur in a Dutch construction in which several verbs are clustered in a group which is preceded by the arguments of these verbs. A typical example is:

- (87) ..dat de man de vrouw de boeken zag kopen  
 ..that the man the woman the books saw buy  
 ‘..that the man saw the woman buy the books’

The following is a more complex, but fully grammatical and acceptable example:

- (88) ..dat de man de leraar zijn kinderen Frans had  
 ..that the man the teacher his children French had  
 kunnen willen laten leren  
 can want let teach  
 ‘..that the man could have wanted to make the teacher  
 teach his children French’

<sup>12</sup>One way out of this problem, adopted in [Gazdar et al., 1985], is to take the position that all idioms which are syntactically flexible are built up compositionally. This position, however, appears to be untenable. See [Abeillé and Schabes, 1989] and [Schenk, 1992].

In this example the verbal cluster *had kunnen willen laten leren* is preceded by the arguments *de man* ‘the man’, *de leraar* ‘the teacher’, *zijn kinderen* ‘his children’ and *Frans* ‘French’.

When we draw lines to indicate the dependencies between verbs and arguments, the lines cross:

- (89) ...dat hij die kinderen dat boek zou willen laten lezen
- 
- ...that he those children that book would want let read  
 ‘...that he would like to let those children read that book’

The phenomenon of crossing dependencies gives rise to structures which can (very abstractly) be characterized as being of the form  $\omega\omega$  and it has played an important role in the discussion whether the syntax of natural languages is context-free (see e.g. Huybregts ([Huybregts, 1976], [Huybregts, 1984]), [Shieber, 1985], and for an overview of this discussion [Pullum, 1991]).

The non-context-free character of these constructions is no problem whatsoever for the grammar formalism, since M-rules are powerful rules. It does, however, require special measures for the surface parser, since the surface parser uses a grammar which in essence is context-free. The part of the surface grammar which deals with the structures involved in crossing dependencies is discussed in section 5.5.4.

One of the most interesting and problematic properties of these structures is their hybrid character: on the one hand they appear to behave like biclausal structures, but on the other hand they appear to behave like monoclausal structures. As indicated by [Evers, 1975] their biclausal aspect is evident from their behavior with respect to S-pronominalization, lexicalization, reflexivization and passivization, though their monoclausal character is evident from their behavior with respect to gapping, nominalization, a number of extraposition phenomena, the placement of clitic elements, sentential negation and ‘quantifier hopping’.

This section is organized as follows: in section 5.5.1 I will describe how crossing dependencies are dealt with in the grammar. First a global description is supplied. In section 5.5.2 the relevant M-rules are described in more detail. The pruning operation, which plays an important role in the derivation of crossing dependencies, is presented in section 5.5.3. The role of the surface grammar with respect to these structures is characterized in section 5.5.4. Section 5.5.5 lists a number of remaining problems and summarizes the conclusions.

### 5.5.1 Outline of the Analysis

The treatment of *crossing dependencies* in the Dutch syntax in Rosetta is an almost direct implementation of the analysis given by [Evers, 1975].

Evers accounts for these structures in the following way. First, he assumes that Dutch is an SOV language and that infinitival complements start out on the left of the verb. Second, he assumes that infinitival complements are sentences. And third, he assumes the existence of a transformation, called *Verb Raising*,

which moves an infinitival verb out of its clause and adjoins it to the right of the embedding verb. Infinitival sentences which have lost their verb by Verb Raising are pruned, i.e. the node S is removed from the structure.

This analysis accounts for the hybrid character of the construction by assuming that there are two stages in the derivation of this construction. In the first stage these structures are biclausal, in the second stage they are monoclausal. These two stages are related to one another by the rule of Verb Raising in combination with pruning.

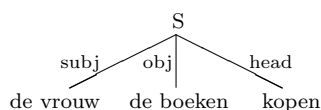
In Rosetta, these structures are accounted for in almost the same way: all Evers's assumptions mentioned have been implemented in Rosetta, though some slight modifications have been made, especially with respect to the pruning operation.

Within the Rosetta grammars, however, other considerations also play a role. First of all, because of the isomorphic grammar method, the grammars of the languages dealt with must be attuned. As a consequence, the analysis of Dutch Verb Raising structures is co-determined by the analysis of the translations of these sentences in English and Spanish. When one takes these translations into consideration, one can see that Dutch Verb-Raising structures usually correspond to biclausal structures in English, though in Spanish, biclausal structures occur as do structures with a hybrid character (i.e. they show a monoclausal character with respect to some phenomena, but a biclausal character with respect to others).

Given these facts, an analysis of Verb-Raising structures as biclausal structures which are turned into monoclausal structures is natural, since it can be justified both on monolingual and multilingual grounds.

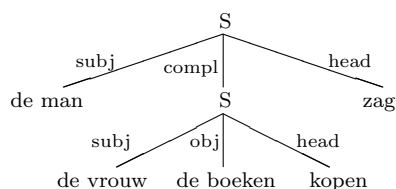
I will illustrate the analysis by going through the derivation of sentence (87). First, the embedded complement is generated as an infinitival sentence, as illustrated in (90):

(90)



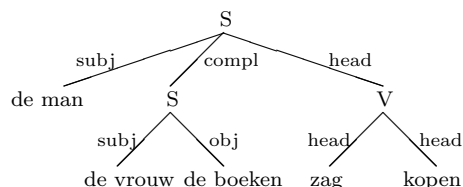
This infinitival sentence is used as a complement to the verb *zien* and put to the left of it (in accordance with the SOV status of Dutch):

(91)



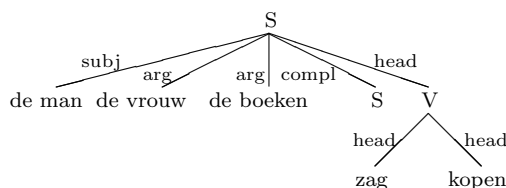
The Verb-Raising transformation is applied to this structure. It raises the verb *kopen* out of its clause and adjoins it to the verb *zag*:

(92)



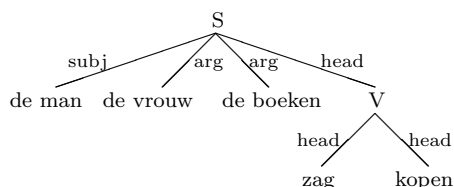
Next, the embedded S-node must be pruned. The S is emptied first, i.e. all nodes dominated by S are moved from under this S-node. In addition, the relations *subj*, *obj* and *indobj* (the last mentioned does not occur in the example structures) are turned into the relation *arg*. This is not really necessary in generation, but it is done to avoid ambiguities in the surface parsing process (see below). This yields:

(93)



When S no longer dominates any nodes, it is removed from the structure:

(94)



The subordinate clause can be derived by adding the subordinate conjunction *dat* to the sentence.

### 5.5.2 Verb-Raising Transformations

Verbs have the attribute **verb-raiser** with as possible values *noVR*, *optionalVR* and *obligatoryVR*. The value *noVR* is intended for verbs that do not allow Verb Raising, e.g. *zich schamen* 'to be ashamed', *afspreken* 'to make an appointment'.

The value *OptionalVR* is intended for verbs that allow Verb Raising but do not require it, e.g. *proberen* ‘to try’. The value *ObligatoryVR* is intended for verbs that require Verb Raising, e.g. *laten* ‘to make, let’, *schijnen* ‘to seem’.<sup>13</sup> When a verb takes a sentential complement, either extraposition of this sentential complement, or Verb Raising must take place. The values of the attribute *verb-raiser* co-determine which rule must apply.

Verb-Raising transformations apply after the introduction of sentential complements and after the application of rules dealing with the temporal and aspectual properties of the clause. These temporal and aspectual rules might have introduced all kinds of auxiliary verbs (*zijn*, *hebben*, *worden*), so that a structure that is input to the Verb-Raising transformations can be represented as follows:

...pred/VP[.compl/S[...r/V<sub>1</sub>] head/V<sub>2</sub> (progaux/V<sub>3</sub>) (aux/V<sub>4</sub>)...]...

where *r* is a variable over relations, and the round brackets indicate optional presence. Auxiliary verbs are introduced either as a verb bearing the grammatical relation *progaux* (e.g. *zijn* when accompanied by *aan het* + infinitive), or as a verb bearing the grammatical relation *aux* (the verbs *hebben*, *zijn*, *worden* when used as perfect or passive auxiliaries).

The transformation VR1 raises V<sub>1</sub> (the rightmost verb cluster) out of the embedded complement and adjoins it to V<sub>2</sub>, if V<sub>2</sub> has the value *optionalVR* or *obligatoryVR* for the attribute *verb-raiser*. This results in a structure of the form:

...pred/VP[...compl/S[...] r/V<sub>5</sub>[head/V<sub>2</sub> adjoin/V<sub>1</sub>] (progaux/V<sub>3</sub>) (aux/V<sub>4</sub>)...]...

where the node V<sub>5</sub> is created which gets the relation *r*. If V<sub>2</sub> is a past participle in this structure in an active sentence, then it is turned into an infinitive.<sup>14</sup>

I will supply some concrete examples to illustrate:

...[...zwemmen] kan... →...[...] [kan zwemmen]...  
 ...[...gezwommen hebben] kan... →...[...gezwommen] [kan hebben]...  
 ...[...te zwemmen] [schijnt]... →...[...] [schijnt te zwemmen]  
 ...[...zwemmen] gewild heeft →...[...] [willen zwemmen] heeft

The cluster V<sub>1</sub> may already have incorporated several verbs in earlier cycles.

...[...[te willen laten lezen] schijnt ...→...[...] [schijnt [te willen laten lezen]] ...  
 ...[...[hebben gezwommen]] kan ... →...[...] [kan [hebben gezwommen]] ...

<sup>13</sup>There are certain generalizations between the types of complements a verb can take and whether it allows, disallows or requires Verb Raising. Such generalizations should be expressed in the grammar, and they partly have been.

<sup>14</sup>See [Hoeksema, 1988] for a suggestion of why the participle is excluded here. To my knowledge, there is no account for the fact that, instead of the participle, the infinitival form can appear.

In this way, verb clusters of — in principle — unlimited size can be created.

Apart from the transformation VR1, there is a transformation VR2 that raises a verb that is part of a combination of *particle + (te)+ verb* and adjoins it to a higher verb:

...[...opbellen ] kan... →...[...op] [kan bellen]...  
 ...[...op [te bellen] ] schijnt... →...[...op] [schijnt [te bellen]]...

After the Verb-Raising transformations other transformations follow to put the auxiliary verbs in their proper positions. Auxiliaries such as *hebben* ('have'), *zijn* ('be') and *worden* ('be') optionally form a cluster with a preceding participle. This transformation can be formulated informally as follows:

...[...head/V<sub>ptc</sub> aux/V<sub>1</sub>... ]... →...[...head/V<sub>2</sub>[aux/V<sub>1</sub> head/V<sub>ptc</sub>]... ]...

Some examples:

(dat hij hem) gedood heeft →(dat hij hem) heeft gedood  
 (dat hij door hem) gedood is →(dat hij door hem) is gedood  
 (dat hij door hem) gedood werd →(dat hij door hem) werd gedood

The combination *auxiliary + participle* forms a verb cluster, but the combination *participle + auxiliary* does not. This can be shown by the fact that the first combination can, though the second one cannot, be subject to VR1:

...dat hij [ hem gedood hebben ] kan →\*...dat hij [ hem ] kan gedood hebben  
 ...dat hij [ hem hebben gedood ] kan →...dat hij [ hem ] kan hebben gedood

If the auxiliary verbs *hebben* ('have') and *zijn* ('be') are preceded by an infinitival verbal cluster, then the formation of a new verbal cluster is obligatory:

...[... ] [willen zwemmen] heeft →...[... ] [heeft [willen zwemmen]]

Such structures can be created after the application of VR1 or VR2 in which a participle has been changed into an infinitive.

Furthermore, there is a transformation that leaves a particle behind (*opgebeld heeft* →*op [heeft gebeld]* ).

### 5.5.3 Pruning

When the Verb-Raising transformations have applied, the embedded S-node must be pruned, i.e. the S-node is removed but the nodes it dominates remain. In Evers's analysis, the S-node is pruned directly when the verb is raised. This is in Evers's analysis a consequence of Kuroda's *Guillotine Principle* (see [Ross, 1967,



56]), which states that when a head is deleted, its projections are immediately pruned.

In the Rosetta analysis, the sentence node is always pruned. This appears to be Evers' intention as well, though he does not specify in sufficient detail how this will work in constructions involving auxiliary verbs. Do auxiliary verbs take sentential complements in his analysis? If so, will this sentence node be pruned when a participle is not raised out of it? If not, what kind of complement do auxiliary verbs take? Is the auxiliary the head of the sentence? Such questions remain unanswered in Evers' analysis.

In the Rosetta grammars the main verb is always the head of the VP. Auxiliary verbs are introduced syncategorematically as a kind of modifiers. Given this, the *Guillotine principle* cannot be adopted if the sentential nodes are always to be pruned (since the head is not always raised).

The S-node from which Verb Raising has been applied is always pruned, but it is emptied first, i.e. all nodes dominated by S are moved from under this S-node (except verbs and particles). This is done in this manner because a reversible rule which simply prunes S will be enormously complex.

When S no longer dominates any nodes apart from verbs and particles, it is removed from the structure. Removing this node requires that it be recovered in analysis. The node itself can be recovered easily, but recovering all the values of the attributes of the S-node in an acceptable way<sup>15</sup> requires care. Furthermore, once the S-node has been recovered, the correct phrases must be put back under it in analysis.

Some of the attributes of the S-node are *voice* (encodes whether the sentence is in active or passive voice), *actvps* (encodes which verb pattern of the possible verb patterns of the verb has actually been used in this structure), *mood* (encodes whether the sentence is declarative, interrogative, imperative etc.), *inftype* (encodes whether we are dealing with a bare infinitive or a combination of *te* + infinitive), *PROsubject* (encodes whether it is a control or raising structure), etc. The values of all these attributes must be recovered in analysis. As a consequence, the pruning transformations are accompanied by a complex calculation to allow deletion of the proper values for the attributes (in generation) and to make their recovery possible (in analysis). These calculations are performed on the basis of the information available in the structure. The need for such calculations is a direct consequence of the requirement that the grammar be reversible.

In order to empty the S-node, transformations remove constituents that are not verbs or particles out of S one by one (retaining their mutual order). In analysis this process is reversed: constituents are put into S one by one, in the correct order. This process is governed by information of the main verb in order to make it as deterministic as possible.

In this process, the relations are retained, except for the relations *subj*, *obj* and *indobj*, which are turned into the relation *arg* in order to facilitate the surface

---

<sup>15</sup>All attributes have a finite number of possible values, so trying out all combinations of all possible values is possible in principle. However, it is not acceptable, because the number of combinations is very large.

parsing process (see the following section).

#### 5.5.4 Surface Grammar

As pointed out before, the constructions involving crossing dependencies cannot be described by a context-free grammar. The surface grammar used by S-PARSER, however, is in essence a context-free grammar.

This problem is solved as a specific instance of the general strategy with respect to the surface grammar. The surface grammar describes a superset of the structures defined by the M-grammar. This is allowed, because the surface grammar has no principled status in the system. See chapter 2.

In the particular case of crossing dependencies there are basically two sets of rules. One set of rules defines which verb cluster can appear in surface structures. Another set of rules specifies that a VP can consist of an unlimited sequence of complements followed by a verb cluster and possibly auxiliary verbs. The relation between the complements and their selecting verbs is *not* made. Since it is this relation which is the cause of the non-context-freeness of these constructions, this problem is solved.

Without additional measures this would lead to an explosion of ambiguities for structures of the form NP\* V\*, because the relation that each NP bears must be chosen randomly. A sentence such as:

- (95) ..dat hij de kinderen dat boek zou willen laten lezen  
 .. that he the children that book would want let read  
 ‘..that he would like to let the children read the book’

would be analyzed in all the following ways if only the two NPs *de kinderen* en *dat boek* are taken into account:

|           |         |           |        |           |           |
|-----------|---------|-----------|--------|-----------|-----------|
| subj/NP   | subj/NP | subj/NP   | obj/NP | subj/NP   | indobj/NP |
| obj/NP    | subj/NP | obj/NP    | obj/NP | obj/NP    | indobj/NP |
| indobj/NP | subj/NP | indobj/NP | obj/NP | indobj/NP | indobj/NP |

In order to avoid these explosions of ambiguities, the correct relation of the NPs is kept vague. Every NP in a Verb-Raising structure that does not belong to the least embedded verb gets the relation *arg*. We make use of the following generalization:

Verb Raisers take (apart from a sentential complement) only NPs as arguments<sup>16</sup>

Because subjects are always NPs, non-NP arguments (PPs, APs, Ss) can only belong to the most embedded verb. By ‘percolating’ the verb pattern of the most embedded verb upward in a verbal cluster the legitimacy of the presence of non-NP arguments of this verb can be checked. Furthermore, the presence of arguments

<sup>16</sup>Henk van Riemsdijk pointed out to me that this generalization might follow in part from the linear locality of Verb Raising, i.e. the verb raised and the triggering verb must be adjacent.

of the least embedded verb can be checked by the verb pattern of this verb. This method reduces the number of ambiguities drastically (for the example given, all ambiguities disappear).

The whole analysis process will be illustrated in more detail with the example *..dat hij de kinderen het boek liet lezen* lit. *...that he the children the book let read* ‘..that he let the children read the book’. I will simplify the example (e.g. no mention is made of VP and the computation of the attributes of S), and omit all irrelevant details in the example structures.

The surface parser yields the following surface tree (many nodes have been assigned indices to be able to identify them):

$$(96) \quad S_1[ \text{subj/NP}_1 \text{ arg/NP}_2 \text{ arg/NP}_3 \text{ head/V[ head/V}_1 \text{ adjoin/V}_2 ] ]$$

Since  $V_1$  (*laten* ‘to let’) triggers Verb Raising, the pruning operation has to be undone in analysis. An S-tree headed by  $S$  and with relation *compl*, but dominating nothing, is inserted before the verbal cluster. The result is represented in (97).

$$(97) \quad S_1[ \text{subj/NP}_1 \text{ arg/NP}_2 \text{ arg/NP}_3 \text{ compl/S}_2[ ] \text{ head/V[ head/V}_1 \text{ adjoin/V}_2 ] ]$$

Next, information about the complementation properties of  $V_2$  is used when putting NPs into this sentential complement.  $\text{NP}_3$  can be put into  $S_2$ , either as a direct object (since  $V_2$  (*lezen* ‘to read’) can be used transitively), or as a subject (each sentence allows a subject). Both are attempted, resulting in the structures (98a,b), respectively:

$$(98) \quad \begin{array}{l} \text{a} \quad S_1[ \text{subj/NP}_1 \text{ arg/NP}_2 \text{ compl/S}_2[\text{obj/NP}_3] \text{ head/V[ head/V}_1 \text{ adjoin/V}_2 ] \\ \quad \quad ] \\ \text{b} \quad S_1[ \text{subj/NP}_1 \text{ arg/NP}_2 \text{ compl/S}_2[\text{subj/NP}_3] \text{ head/V[ head/V}_1 \text{ adjoin/V}_2 ] \\ \quad \quad ] \end{array}$$

At this point there are several possible paths in the analysis. One possibility is to apply Verb Raising in reverse to these structures. Another possibility is to put  $\text{NP}_2$  under  $S_2$ , and to undo Verb Raising only after this.

If the first option is taken, undoing Verb Raising yields the structures (99a,b) respectively:

$$(99) \quad \begin{array}{l} \text{a} \quad S_1[ \text{subj/NP}_1 \text{ arg/NP}_2 \text{ compl/S}_2[\text{obj/NP}_3 \text{ head/V}_2 ] \text{ head/V}_1 ] \\ \text{b} \quad S_1[ \text{subj/NP}_1 \text{ arg/NP}_2 \text{ compl/S}_2[\text{subj/NP}_3 \text{ head/V}_2 ] \text{ head/V}_1 ] \end{array}$$

The rule which introduces sentential complements removes the tree headed by  $S_2$  from these structures. This tree is subjected to M-PARSER again. In (99a) the sentential complement lacks a subject (this subject is not created by control transformations either, since *laten* ‘to let’ is not a control verb), so no successful parse results. In (99b) the complement sentence is syntactically well-formed (the verb *lezen* can be used intransitively as well), but the main clause contains one

NP too many for the verb *laten* ‘to let’, which can only combine with a subject and a complement sentence.

We return to the structures (98), and take the second option: an attempt is made to put NP<sub>2</sub> under S<sub>2</sub>. In (98a) NP<sub>2</sub> can be put under S<sub>2</sub> as a subject, as represented in (100).

$$(100) \quad S_1[ \text{subj/NP}_1 \text{ compl/S}_2[\text{subj/NP}_2 \text{ obj/NP}_3] \text{ head/V[ head/V}_1 \text{ adjoin/V}_2] ]$$

Other options are not available. NP<sub>2</sub> cannot be put into S<sub>2</sub> as an indirect object, since the verb *lezen* ‘to read’ does not allow indirect objects. NP<sub>2</sub> cannot be put under S<sub>2</sub> as a direct object in any of these structures: (98a) already contains a direct object, and in (98b) this would result in an object preceding the subject, which will not lead to a successful parse. NP<sub>2</sub> cannot be put into S<sub>2</sub> in (98b) as a subject, since a sentence cannot contain two subjects.

Now Verb Raising can be undone again. This yields:

$$(101) \quad S_1[ \text{subj/NP}_1 \text{ compl/S}_2[\text{subj/NP}_2 \text{ obj/NP}_3 \text{ head/V}_2 ] \text{ head/V}_1 ]$$

The sentential complement can now be analyzed correctly (the verb *lezen* can be used transitively and it co-occurs with a subject and a direct object in the correct order), and the main clause can also be analyzed correctly (the verb *laten* co-occurs with a subject and an infinitival complement containing an overt subject).

This illustrates how crossing dependencies are dealt with in analysis. In M-PARSER nothing special happens: the analytic functions associated with the rules are applied as usual. In the surface grammar the relation between the predicates and their arguments is not made, avoiding the problem of having to parse a non-context-free language with a context-free grammar: this relationship is established by M-PARSER. In addition, some measures have been taken to avoid unnecessarily inefficient parsing due to the absence of establishing the relation between predicates and arguments in Verb-Raising structures.

The method adopted resembles, to a certain extent, the method proposed by [Bresnan et al., 1982] in the LFG framework. In that method a c-structure is created on the basis of a context-free grammar. The relation between predicates and their arguments is not made at the level of c-structure, but at the level of f-structure, which is created by annotations associated with context-free rules which define the c-structures. Differences between the two methods are: (1) the surface grammar has no principled status in M-grammars, but the rules creating c-structure do have a principled status in LFG, and (2) all rules operate on S-trees, but in LFG two kinds of objects are distinguished: c-structures (which are syntactic trees which contain categorial information) and f-structures (which are directed acyclic graphs containing functional information).

### 5.5.5 Concluding Remarks

I have described how Verb-Raising structures of Dutch are dealt with in Rosetta.

The method adopted in Rosetta is based on the analysis of [Evers, 1975]; it is simple and completely general, and can be extended quite easily to deal with related phenomena in other languages (e.g. Verb-Raising in German, and so-called V-projection Raising as in Flemish and Swiss dialects, see e.g. [Haegeman and van Riemsdijk, 1986]).

There are some differences between the Rosetta treatment and Evers's analysis, which will be summarized here.

The first difference relates to reversibility. Evers's account is not reversible, but the analysis described here is. This had specific consequences for the formulation of the pruning operation (see the description given above). It is very difficult, if not impossible, to formulate by means of the M-rule notation a reversible version of pruning, in which the S-node is directly pruned and where this process is governed by the available elements in the structure to make it as deterministic as possible. In addition, a complex calculation is required to allow deletion or recovery of attribute-value pairs of the S-node.

The second difference concerns pruning. In Rosetta the S-node of the embedded infinitive is always pruned, even if a participle has remained in the structure. It is not clear whether this is the case in Evers's analysis, since Evers does not state explicitly how participles behave under Verb Raising.

A third difference is a specific instance of a more general difference: Evers works within a framework where grammatical relations are not represented. In S-trees, grammatical relations must be specified. These grammatical relations are changed into a unique relation (*arg*) for NP-arguments upon pruning of the sentential node, in order to speed up the analysis process, as described in section 5.5.4.

The fourth difference concerns descriptive detail. The Rosetta treatment has been worked out in full detail, though Evers's analysis leaves a number of aspects of this construction rather vague or completely unspecified. In particular, the behavior of participles and of particles in these constructions, the so-called *Infinitivus-Pro-Participio* effect, and some deviating behavior of certain auxiliary verbs and modal verbs are accounted for in full in Rosetta, but described only marginally or not at all by Evers. This reflects a more general difference between the purposes of theoretical linguistics and linguistic engineering, as described in chapter 1.

Obviously, there are still many problems to be solved. I will mention some problems that require further investigation. First, there are some indications that pruning of sentential nodes might not be the correct way to describe the fact that certain rules act as if they operate in monoclausal structures. Clitics, e.g. can move to a clitic position of a superordinate verb (suggesting a monoclausal structure), but they can also move to intermediate positions (suggesting a multiclausal structure) in certain cases (see [Kroch and Santorini, 1991, 284]).

A second problem area concerns the so-called *Third Construction*. This construction has properties in common with Verb Raising and with extraposition structures. NP-arguments of the embedded verb can precede the embedding verb (as in VR-constructions), but there are no IPP-effects (which suggest that extraposition is involved). An example is given in (102) (taken from [Broekhuis, 1992, 71]):

- (102) ...dat Jan dat boek heeft geprobeerd te lezen  
 ...that Jan that book has tried to read  
 ‘...that Jan has tried to read that book’

When we developed our analysis, no consistent interpretation of the facts which later motivated the postulation of the ‘Third Construction’ was available. In addition, many examples are rather marginal (and not part of the standard language). Nevertheless, we noticed that certain verbs did not participate in the IPP-effect though they apparently did participate in Verb Raising. In the system, we added the possibility to specify verbs as lexical exceptions to the IPP-effect. This is clearly insufficient, because it should at least be accompanied by a different structure of the verb cluster in order to account for inversion of auxiliary and main verb in the ‘Third Construction’. It has been suggested that the Third Construction is created by extraposing the complement sentence, followed by preposing the NP-complements of the embedded verb into the embedding clause (Den Besten and Rutten ([den Besten and Rutten, 1989])). An alternative analysis of the Third Construction is considering it ‘Verb Raising without the IPP effect’, as suggested by [Model, 1991a]. Such an analysis would account for the fact that the Third Construction yields the best results with verbs which are not, or only with great difficulty, subject to IPP (e.g. *zeggen* ‘to say’, *beloven* ‘to promise’, *denken* ‘to think’, etc.). In this analysis certain verbs are marked (perhaps due to their word-internal make-up) to form exceptional verb cluster structures (left-branching instead of, or in addition to, right-branching). Adopting such a structure would destroy the configuration required for IPP, it would allow for inverted orders (*..dat hij dit probleem zo had gedacht / gedacht had op te kunnen lossen* lit. *that he this problem this way had thought / thought had to be able to solve*), and it would account for the fact that these constructions are equal to Verb-Raising structures in all other respects. But this analysis cannot account for sentences such as (103) from [Broekhuis, 1992, 187]:

- (103) ...dat hij het hek geprobeerd heeft groen te verven  
 ...that he the gate tried has green to paint  
 ‘...that he has tried to paint the gate green’

Such sentences are apparently judged to be well-formed by many people. It is clear that the Third Construction requires further work.

In conclusion, it can be stated that the formalism makes it possible to incorporate a variant of the analysis of crossing dependencies fairly directly, and it is relatively easy to extend this analysis to obtain better coverage of the relevant facts. The analysis of crossing dependencies presented is based on Evers’s analysis, but it has been adapted so that it becomes reversible. In addition, it has been extended to achieve observational adequacy. The analysis has been implemented in full. The fact that Evers’s analysis can be incorporated into grammars in the compositional M-grammar framework directly makes it possible to describe a very complex array of facts concerning the distribution of infinitival verbs and their arguments and accompanying adjuncts with a minimum of machinery.

## Chapter 6

# R-pronouns

### 6.1 Introduction

R-pronouns are a special class of pronouns in Dutch consisting of the seven pronouns *er*, *hier*, *daar*, *ergens*, *nergens*, *overal*, *waar*. The members of this class share a number of syntactic properties which other pronouns do not possess. They have been coined *R-pronouns* by Van Riemsdijk ([van Riemsdijk, 1978]) who proposed distinguishing them from other pronouns by the feature *R*. The name of this feature was inspired by the fact that each of the relevant pronouns contains the letter *r*.

The syntax of R-pronouns is a notoriously difficult area of Dutch grammar. The problematic aspects of these words involve their function(s) within a sentence (most R-pronouns can serve more than one function in a sentence) and their distribution (when and where, can or must they occur?). Of these R-pronouns, the R-pronoun *er* shows the most complex behavior. For an extensive discussion of these and related phenomena, see Bech ([Bech, 1952]), Van Riemsdijk ([van Riemsdijk, 1978]), Bennis ([Bennis, 1980],[Bennis, 1986]) and Model ([Model, 1991a]).

In this chapter I will describe how the syntax and semantics of R-pronouns are dealt with in the Rosetta grammar for Dutch. The analysis presented takes Van Riemsdijk ([van Riemsdijk, 1978]) as a starting point, and an attempt is made to extend this analysis to other facts. Furthermore, an alternative is presented to the (in my view) unsatisfactory deletion analysis given by Bennis ([Bennis, 1980],[Bennis, 1986]) to account for the fact that one occurrence of *er* can have several functions at the same time.

The resulting analysis is, to my knowledge, the only analysis of R-pronouns (except the one given by Bennis) which explicitly accounts both for the distribution of R-pronouns and for the fact that certain R-pronouns can have multiple functions. I will argue that the analysis outlined here is superior to Bennis's analysis.

Section 6.2 introduces some relevant facts; in section 6.3 the different functions of R-pronouns are discussed in detail. Section 6.4 takes up the problem of the distribution of R-pronouns, and a general approach to this problem is given. In

section 6.5 this general approach is worked out in a number of specific assumptions, each of which is explained and justified. In section 6.6, I will illustrate in detail, how these assumptions account for the relevant facts. In section 6.7 some remaining problems are discussed and the conclusions are presented.

## 6.2 Some Relevant Facts

I will first briefly introduce some of the phenomena that must be accounted for. First, the word *er* can occur in several functions: as an expletive element (as in (104a)), as a locative adverbial (as in (104b)), as an element that must appear if quantified headless count NPs occur (quantificational use), as in (104c), and as a form that must appear instead of the pronoun *het* as a complement to prepositions (prepositional use), as in (104d):

- (104) a *Er* werd gedanst  
           *There* was danced  
           There was dancing going on  
       b Hij woont *er*  
           He lives *there*  
           ‘He lives there’  
       c Hij ziet *er* twee  
           he sees *there* two  
           ‘He saw two of them’  
       d Hij kijkt *ernaar*  
           he looks *there* at  
           ‘He is looking at it’

These functions will often be indicated by means of subscripts on the R-pronoun or its gloss: *X* for the expletive use, *L* for the locative use, *P* for the prepositional use, and *Q* for the quantificational use.

Many of these functions can be combined in one occurrence of *er*<sup>1</sup> (though not all combinations are allowed), e.g. expletive and prepositional use (105a), quantificational and prepositional use (105b), expletive and locative use (105c), expletive, quantificational and prepositional use (105d), etc. (for an overview see Bennis ([Bennis, 1986]), Model ([Model, 1991a])):

- (105) a *Er* werd naar gekeken  
           *There* was at looked  
           ‘It was being looked at’  
       b Hij beschuldigde *er* twee van  
           He accused *there* two of  
           ‘He accused two (of them) of it’

---

<sup>1</sup>In such cases several subscripts will often be used to indicate which functions an R-pronoun fulfills.



- c *Er* woont iemand  
*There* lives someone  
 ‘Someone lives there’
- d Werden *er* twee van beschuldigd?  
 Were *there* two of accused  
 ‘Have two (of them) been accused of it?’

The following sentences illustrate combinations that are not possible:

- (106) a \**Er*<sub>XQ</sub> waren twee  
 There were two  
 ‘There were two of them’
- b \*Hij legde *er*<sub>LP</sub> een artikel over  
 He put there an article about  
 ‘He put an article about it there’

Furthermore, though the appearance of two occurrences of *er* in one sentence cannot be excluded in general (cf. (107a)), and though they can even occur adjacent to one another in certain sentences (cf. (107b)), many sentences do not allow the presence of two occurrences of *er*, not only when they are adjacent (as in (107c)), but also when they are at a certain distance (as in (107d)):

- (107) a *Er* werden *er* twee gekocht  
*There* were *there* two bought  
 ‘Two (of them) were bought’
- b Hij keek *er* *erna* naar  
 He looked *there* *thereafter* at  
 ‘He looked at it afterwards’
- c \*Werden *er* *er* twee gekocht?  
 Were *there* *there* two bought?  
 ‘Have two (of them) been bought?’
- d \**Er* werd gisteren *ernaar* gekeken  
*There* was yesterday *thereat* looked  
 ‘It was looked at yesterday’

These, and many more problems, should be accounted for adequately.

I will begin by discussing each of the separate functions which R-pronouns can fulfill, and the way these have been accounted for in the Rosetta system. Then, I will turn to a discussion of their interaction, and describe how their distribution and the fact that several functions can coincide are accounted for.

## 6.3 The Functions of R-pronouns

### 6.3.1 Expletive Function

The pronoun *er* can function as an expletive in several sentence types, illustrated in (108):

- (108) a Er werd gedanst  
 There was danced  
 ‘There was dancing going on’  
 b Er werd een boek gekocht  
 There was a book bought  
 ‘A book was bought’  
 c Er danste iemand  
 There danced someone  
 ‘Someone was dancing’  
 d Er kocht iemand een boek  
 There bought someone a book  
 ‘Someone bought a book’  
 e Er wordt beweerd dat dat niet waar is  
 There is claimed that that not true is  
 ‘It is claimed that that is not true’

Expletive *er* occurs in impersonal passives (108a), in passive and active sentences where the ‘subject’ occurs more to the right in the clause, in the *postsubject* or *object* position (see (108b,c,d)), and in passives with a sentential complement.

In the analysis to be outlined, expletive *er* is not a subject at any point in the derivation. Expletive *er* occupies either a special position, the so-called R-position (encoded by the relation *erpos*), or a topicalized position (encoded by the relation *shift*), but never the subject position. See Hoekstra ([Hoekstra, 1984]:220ff), Schermer-Vermeer ([Schermer-Vermeer, 1985], [Schermer-Vermeer, 1986b], [Schermer-Vermeer, 1986a]), [Paardekooper, 1986] and Bennis ([Bennis, 1986]:212) for discussion of this issue. This point of view will be justified extensively below. However, I make the additional assumption that expletive *er* is related to the subject position.

In the derivation of a sentence which contains an occurrence of expletive *er*, three phases can be distinguished: (1) a phase where the subject position is not occupied by any element, (2) a phase where the subject position is occupied by an abstract expletive element, and (3) a phase where the abstract expletive element is removed from the subject position under appropriate conditions; if necessary and possible, the pronoun *er* is introduced into the R-position.

In the first phase the subject position is not occupied by any element, either because there simply is no subject (in passives and ergatives) or because the subject has been placed in a different position by a rule creating so-called *postsubjects*. At a certain point (phase 2), however, either the direct object is put into the subject position (in the case of ergative and passive verbs), or an abstract expletive element dominated by NP is inserted into the subject position. The abstract expletive element can either be the element EC, or the element EREC. In both cases a relation is established between the elements in the *postsubject* position (if there is one), and with the element in the *object* position (if there is one), or with no element at all. The number feature of the element in *postsubject* position if there is one, or otherwise the number feature of the element in *object* position, is copied onto the NP-node dominating the abstract expletive element. If there

are no elements in the *postsubject* or *object* position the number attribute is set to *singular*. If the abstract element EREC is introduced, several restrictions on the definiteness of the related elements (described below) are checked. No such restrictions hold if the element EC is introduced, but the introduction of this element is subject to different requirements (described below). Finally, in the third phase, the pronoun *er* is introduced into the R-position, and the subtree of category NP which contains EREC is deleted. A special transformation deletes the subtree of category NP which contains EC if certain conditions described below are satisfied.

This approach can be justified as follows. Note that the finite verb agrees with an element in postsubject position (if any) or even with elements in object position (if any) in such constructions:

- (109) a Er kocht(\*en) een jongen een boek  
 There bought a boy a book  
 ‘A boy bought a book’  
 b Er kocht\*(en) jongens een boek  
 There bought boys a book  
 ‘Boys bought a book’  
 c Er werd(\*en) een boek verkocht  
 There was/\*were a book sold  
 ‘A book was sold’  
 d Er werd\*(en) boeken verkocht  
 There \*was / were books sold  
 ‘Books were sold’

In order to simplify the agreement rule, the abstract expletive elements are postulated in the subject position, and the number features of the postsubject or the object are copied to the expletive phrase. This makes it possible to keep the agreement rule simple: it need refer to the subject only. Furthermore, it is absolutely necessary to formulate the agreement rule in this manner for these cases, since in certain constructions the finite verb of a clause can agree with the direct object or the postsubject of embedded infinitival clauses, as in (110), where the NP *boeken/een boek* is the complement of the verb *verkocht* which itself is a complement of *kunnen* which itself is a complement of the verb *moeten* which is a complement of the verb *schijnen*; the NP must agree with this latter verb.

- (110) a Er \*schijnt/schijnen vandaag boeken verkocht te moeten kunnen worden  
 b Er schijnt/\*schijnen vandaag een boek verkocht te moeten kunnen worden  
 lit. There seem(s) to must can be sold a book/books  
 ‘It seems that it must be possible to sell books today’

Because such dependencies can extend over indefinite distances, it is impossible to write an M-rule which relates the finite verb to these elements directly. Instead, there is a rule relating an abstract expletive element to these elements locally (within one clause), and this will directly yield the relevant agreement facts. The

rule will interact with other independently required rules, such as the subject-verb agreement rule, and the subject-to-subject-raising rule. In essence, this is the analysis of similar phenomena involving expletive *there* in English as given by Chomsky ([Chomsky, 1981]:87).

The rule of subject-to-subject raising provides us with a second argument for the postulation of the abstract expletive elements: this rule raises *subjects* to other subject positions, but it does not raise elements from the *R-position* to any other position. Recall that expletive *er* never occupies the subject position. This, however, implies that *er* is never subject to the subject-to-subject-raising rule, though sentences such as (110) appear to indicate that it should be. This problem can be overcome by the postulation of the abstract expletive elements: *these* elements are related to the relevant argument, *these* elements occupy the subject position, and *these* elements undergo the subject-to-subject-raising transformation.

The rule which introduces expletive *er* introduces it only in finite clauses. This accounts for the interpretation of sentences such as (111), in which *er* can only be interpreted as a locative adjunct:

- (111) a    hij liet er iemand spelen  
           lit. he let there someone play  
           ‘He let someone play there’  
       b    hij zag er iemand binnenkomen  
           lit. he saw there someone enter  
           ‘He saw someone enter there’

Note the minimal contrast between (111b) and a sentence such as *hij zag dat er iemand binnenkwam* ‘he saw that there entered someone’, in which *er* is in a finite clause and can be interpreted as an expletive.

There are certain examples which appear to indicate that expletive *er* occupies a slightly different position than non-expletive uses of *er*. These phenomena will be discussed below.

It is now necessary to specify in more detail under which conditions the abstract elements EC and EREC can be inserted into the subject position. The element EC can be inserted into the subject position if the subject position is unoccupied and if there is an indirect object in the VP, or if the relevant verb is reflexive. The latter case can be illustrated with examples such as:

- (112) a    Gisteren heeft zich het grootste ongeluk voorgedaan  
           Yesterday has itself the largest accident occurred  
           ‘Yesterday the largest accident occurred’  
       b    Vandaag dienen zich deze nieuwe problemen aan  
           Today present themselves these new problems  
           ‘Today these new problems present themselves’

Recall that the abstract element EC is simply deleted later in the derivation, after subject-verb agreement, and that no restrictions whatsoever are imposed by EC on the definiteness of the direct object. An additional restriction on the

insertion of EC is that the direct-object position — if present — is not occupied by an NP headed by a personal pronoun. This restriction is required to prevent the generation of sentences such as (113b,d):

- (113) a Haar vader werden gisteren de kinderen toevertrouwd  
 Her father were yesterday the children put under care  
 ‘The children were put under her father’s care’  
 b \*Haar vader werden gisteren zij/hen/hun toevertrouwd  
 Her father were yesterday they/them put under care  
 c ...dat de meisjes alleen die man opviel  
 ...that the girls only that man struck  
 ‘...that only that man struck the girls’  
 d \*...dat de meisjes alleen hij/hem opviel  
 ...that the girls struck only he/him

Such facts — observed by Bennis ([Bennis, 1986]:161)<sup>2</sup> — do not have a principled explanation in any framework known to me, so I simply added conditions to the relevant rule to obtain observational adequacy.

Den Besten ([den Besten, 1985]) analyses these constructions in which I have postulated an abstract expletive element EC as involving movement of the indirect object into the subject position. In order to avoid confusion I will no longer refer to the relevant NP by its grammatical function (indirect object), or by its position, which is under discussion here, but by the term *dative NP*, because it is marked with dative case.

I have not assumed movement of the dative NP into the subject position for the following reasons. First, such an account would complicate the mechanism to deal with agreement between a finite a verb and a (possibly deeply embedded) direct object. The method adopted above will not work if the dative NP moves into the subject position, so an additional method must be worked out, or a new method must replace the one proposed above. In such an alternative method the subject-verb agreement rule cannot be formulated in such a way that only a reference to the subject and the verb is necessary, so that it probably will have to be more complex.

Second, if the dative NP actually occupies the subject position, one would expect a dative clitic pronoun in subject position to behave like a subject clitic. Subject clitics can, but object clitics cannot be topicalized (see Koster [Koster, 1978a]). Dative clitic pronouns cannot be topicalized in the constructions under discussion here either:

---

<sup>2</sup>Bennis suggests that the ill-formedness of these examples might perhaps be attributed to the fact that there are special rules for the distribution of (weak) pronouns. This, however, cannot be correct, since such special rules do not exist for emphatic pronouns (as in (113d)) which are also ill-formed in this construction.

- (114) a ..dat je mag komen  
           ‘..that you can come’  
       b je mag komen  
           ‘you can come’  
       c ..dat ik je een boek geef  
           ‘..that I give you a book’  
       d \*je geef ik een boek  
           you give I a book  
       e ..dat je dat boek niet was opgevallen  
           ..that you that book not was fallen-up  
           ‘that the book didn’t strike you’  
       f \*je was dat boek niet opgevallen  
           you was that book not fallen-up

The examples (114a,c,e) illustrate the non-topicalized clitic personal pronoun *je* ‘you’, in the subject position, in the indirect-object position, and as a dative NP in the construction relevant here. The corresponding examples in which the clitic personal pronoun is topicalized are given in (114b,d,f) respectively. A clitic personal pronoun can be topicalized from the subject position (114b), but not from the indirect-object position (114d). Topicalization of the clitic pronoun in (114f) is not possible either, which suggests that the clitic pronoun in (114e) is not in the subject position.<sup>3</sup>

A third argument can be constructed by taking the position of locative *er* into account. Locative *er* can occur either to the left or to the right of indirect-object NPs. Definite subjects, however, can only precede and never follow locative *er*. This is illustrated in (115):

- (115) a ..dat hij de man er een boek gaf  
           ..that he the man there a book gave  
       b ..dat hij er de man een boek gaf  
           ..that he there the man a book gave  
           ‘..that he gave the man a book there’  
       c ..dat de man er het boek gaf  
           ..that the man there the book gave  
           ‘.. that the man gave the book there’  
       d \*..dat er de man het boek gaf  
           ..that there the man the book gave

In (115a) the indirect object *de man* precedes locative *er*, and in (115b) the indirect object follows *er*. Both sentences are fully grammatical. In (115c) the

<sup>3</sup>This presupposes that the possibility of a clitic to topicalize depends on the position it occupies. It may be the case, however, that other factors are crucial here. The argument collapses if it can be shown that the relevant factor is not the position of the clitic, but whether the topicalized clitic agrees with the element in the *conj* position. But if this were the correct analysis, it would only diminish the number of arguments against the hypothesis that the dative NP is in the subject position. It does not provide evidence in favor of this hypothesis.

definite subject *de man* precedes locative *er*, in (115d) it follows. Only (115c) is well-formed. These facts supply us with the opportunity to test whether a certain NP is in the subject position or in the indirect-object position. Now consider the sentences in (116):

- (116) a ..dat de man er alleen het meisje was opgevallen  
 .. that the man there only the girl was fallen-up  
 b ..dat er de man alleen het meisje was opgevallen  
 ..that there the man only the girl was fallen-up  
 ‘..that only the girl struck the man there’

As shown in (116) the dative NP *de man* can either precede or follow locative *er*. This argument is a variant of an argument from De Haan ([Haan, 1979]:4.4.3), who uses it in a different context. From these facts it must be concluded that the dative NP is in the indirect-object position.

In principle, the two analyses make further different predictions. Thus, in Dutch, the subject position and the indirect-object position can be separated by all kinds of material, e.g. several types of clitic pronouns, adverbs, etc. However, I have not been able to find additional testable predictions. E.g. one possible argument could run as follows: when the direct object is the pronoun *'t*, then the normal order, indirect object precedes direct object, does not hold. Instead the direct object precedes the indirect object. The subject, however, precedes both. This fact makes it possible, in principle, to test the different hypotheses: if the order *dative-NP 't* is well-formed in the relevant constructions, this would be evidence for the assumption that the dative NP is in the subject position. If this order is impossible, it would constitute evidence for the assumption that the dative NP is in the indirect-object position. In fact, the pronoun *'t* cannot follow a dative NP in these cases, suggesting that the dative NP is in the indirect-object position and not in the subject position. Unfortunately, this test is not reliable, because, as we have seen above, the direct object in these constructions can never be a personal pronoun. Though a direct test to see whether this holds for *'t* as well is impossible, it is very likely that this will indeed be the case, so it is unclear whether the test designed above really distinguishes between the intended cases.

The arguments given by Den Besten ([den Besten, 1985])<sup>4</sup> in favor of the assumption that the dative NP moves into the subject position are rather weak. Den Besten acknowledges this in his own remarks on his paper in Den Besten ([den Besten, 1989]), and he concludes that ‘we might as well give up the idea that the dative NP ever shows up there (i.e. in the subject position, JO) at all’.

Summarizing, Den Besten argued that the dative NP can occupy the subject position. His arguments turned out to be insufficient. Several pieces of evidence were presented against Den Besten’s position, and Den Besten himself gave up this

<sup>4</sup>The paper mentioned appeared in 1985, but the proposal dates back at least to 1980 (see Den Besten ([den Besten, 1981])). The paper appears in Den Besten ([den Besten, 1989]) as well, where remarks concerning it have been added by Den Besten. These will be discussed below. I will from now on only refer to Den Besten ([den Besten, 1989]).

position. I conclude that the dative NP can never occupy the subject position. And this is what has been implemented in the Rosetta system.

The conditions that hold for the insertion of EREC are more complex. Bennis ([Bennis, 1986]) has an extensive discussion of the restrictions that hold. Bennis claims that the conditions on this rule are pragmatically motivated. Though his account may be correct for expletive *er* when no postsubjects or direct objects are present, it does not make much sense in accounting for the definiteness restriction imposed by *er*.

Bennis's account roughly runs as follows. He assumes that NPs are 'presuppositional' in varying degrees, with definite and specific indefinite NPs more 'presuppositional' than non-specific indefinite NPs. He also assumes a condition which states that there must always be at least one non-empty element in a sentence which is not presuppositional (the Empty Presupposition Condition). The R-pronoun *er* can be inserted to avoid violations of the Empty Presupposition Condition.

There are many problems with this account. Bennis claims, at several points, that *er* can only occur if all arguments in a sentence are indefinite: 'we may conclude that the acceptability of the insertion of expletive *er* decreases if the sentence contains a definite argument' (p. 214); 'expletive *er* appears if none of the arguments of the sentence in which *er* is contained is definite' (p. 216); 'expletive *er* appears if all arguments are [-spec] indefinite NPs' (p.227).

This claim, however, is incorrect. Expletive *er* imposes only definiteness restrictions on certain arguments, in particular on the postsubject and the direct object<sup>5</sup>, but not on a prepositional object (in the sense of [Quirk et al., 1972, 831]) or the NP in a *by*-phrase. It is unclear why in a pragmatic account only the definiteness of certain phrases would be relevant. Note that several minimal pairs can be made illustrating this:

- (117) a Er werd een / \*het programma bekeken  
 There was a / \* the program watched  
 'A / the program was watched'
- b Er werd naar een / het programma gekeken  
 There was at a / the programme looked  
 'A / The program was looked at'
- c Er danste een / \* de jongen  
 There danced a / \*the boy
- d Er werd door een / de jongen gedanst  
 There was by a / the boy danced  
 'A boy was dancing'

In (117a) the argument of the passivized transitive verb *bekijken* is realized as a direct object, and — due to the presence of *er* — it must be indefinite. However,

<sup>5</sup>The situation with regard to indirect objects is more complicated. For some people a sentence such as *er werd de jongen een prijs verleend* 'the boy was awarded a prize' is out. I think this sentence is ok, though its subordinate counterpart is significantly worse: *?? .dat er<sub>X</sub> de jongen een prijs werd verleend*. Because of the differing judgements, I will not use these sentences in the argumentation.



in the virtually synonymous sentence (117b) the argument of the verb *kijken* is realized as a prepositional object, as a complement to the preposition *naar*: here, this argument can be definite despite the presence of *er*. It is highly unlikely that other semantic factors are involved here since *kijken naar* and *bekijken* are virtually synonymous. In (117c) the postsubject must be indefinite due to the presence of *er*. In (117d), however, in which the same argument is realized inside a PP (a *by*-phrase) in a passive construction no such restriction holds.<sup>6</sup>

Bennis adduces the following examples to show that the definiteness of complements to a preposition is relevant (p. 214, example (81)):

- (118) a     dat er niemand op een cadeau rekende  
           that there nobody on a present counted  
           ‘that nobody counted on a present’  
       b     ??dat er niemand daar op rekende  
           that there nobody there on counted  
           ‘that nobody counted on that’

But this is hardly a minimal pair. If we construct real minimal pairs, we find that a definite NP instead of the indefinite NP *een cadeau* yields a perfect sentence (*dat er niemand op dat cadeau rekende*), and we find that the ill-formedness of (118b) has nothing to do with the definiteness of *daar*: replacing *daar* with the indefinite *ergens* yields a sentence which is at least as bad (*??..dat er niemand ergens op rekende*). See below, section 6.5, for a possible account of why (118b) is deviant.

In addition, Bennis supplies no independent criteria to determine the status of an NP as either presuppositional or not. A natural interpretation would be: a phrase is presuppositional if it belongs to the presupposition of a clause, and in fact Bennis appears to have something like this in mind, e.g. on p. 223 Bennis describes the extremes of his ‘Presuppositional Hierarchy’ (weak pronouns v. non-specific indefinite NPs) and he states: ‘it seems clear that non-specific indefinite NPs belong to the part of the sentence that conveys new information, i.e. Focus’. But this would imply that certain definite phrases can be non-presuppositional as well in certain sentences, e.g. if they are focussed, or if they supply the answer to a wh-question (e.g. *John* in: *Who left? John did*). But to my knowledge no definite object or postsubject ever allows the presence of expletive *er*. One cannot say: *\*Er kwam Jan binnen*, lit. *There came John in*.<sup>7</sup> Since Bennis does not supply any clear independent definitions or criteria to establish whether phrases are presuppositional or not, I can interpret this notion only as a different term for *specific definite*. For this reason, it makes no sense to say that the

<sup>6</sup>If the *by*-phrase is an adjunct, then this argument is not valid. However the other example remains valid, and there is no argument in favor of a pragmatic treatment. Why would pragmatics be sensitive to a distinction between arguments and adjuncts, which is a syntactico-semantic distinction?

<sup>7</sup>With the possible exception of structures in which the phrase following is an enumeration. The fact that such structures are well-formed only with enumerations does not follow from Bennis’s theory, so they cannot be used as an argument in favor of it.

fact that expletive *er* appears if all arguments are [-spec] indefinite NPs ‘follows from the assumption that expletive *er* is a pragmatic dummy that fills an empty presupposition’ (p. 227).

It may be that Bennis’s account of the pragmatic role of *er* is correct, but I do not see that this can account for the definiteness restriction. As far as I can see, this must be formulated as a purely syntactic condition, since the definiteness restriction is imposed on certain arguments but not on others, as we have seen above, and neither semantics nor pragmatics but only syntax can make the correct distinctions between these arguments.

Basically, the facts are the following: specific, generic and definite NPs can be inserted in the subject position, but not in the postsubject position. Non-specific indefinite NPs can be inserted in the postsubject position but not in the subject position. If the subject position is not occupied, EREC is inserted into the subject position, provided the following conditions hold:<sup>8</sup>

- if neither a postsubject nor an object is present.
- if a direct object is present, but no postsubject, then the properties of the direct object determine whether EREC can be inserted or not: the direct object must be indefinite, non-specific and non-generic.
- if a direct object and a postsubject are present, then EREC can be inserted if both are indefinite, non-specific and non-generic.

Later in the derivation, EREC will disappear when expletive *er* is inserted or by other means (see below).

### 6.3.2 Prepositional R-pronouns

R-pronouns in the prepositional function are always introduced syncategorematically, to replace their corresponding pronoun. The replacements carried out are given in table (119). The set of R-pronouns corresponding to the normal pronouns is exactly the set of R-pronouns which can be used as locatives when occurring independently. In analysis, both the pronouns *iets*, *niets*, *alles* and the R-pronouns *ergens*, *nergens*, *overal* are accepted when governed by a preposition.

(119)

| pronoun | R-pronoun | pronoun | R-pronoun |
|---------|-----------|---------|-----------|
| het     | er        | iets    | ergens    |
| dit     | hier      | niets   | nergens   |
| dat     | daar      | alles   | overal    |
| wat     | waar      |         |           |

<sup>8</sup>Additional distinctions must be made depending on one’s classification of expressions such as *alle*, *iedere*, *etc.* If one classifies these as indefinite (which might be justified by their behavior with respect to adjectival suffixes), an additional distinction is required. However, if one classifies them as definite (as in most semantic theories), no additional distinction is required.

The replacements are carried out only when the preposition specifies that this is possible. This specification is given by the attribute *postform*, which can have one of the following values:

**pre** the preposition can only be used with R-pronouns, e.g. *mee, toe, heen*: *ermee* v. *\*mee een hamer*

**post** the preposition cannot be used with R-pronouns, e.g. *met, tot, tijdens*: *\*ermet, \*ertot, \*ertijdens*

**both** the preposition can be used with an R-pronoun and with normal complements, e.g. *op, voor* : *op het huis, voor de oorlog* v. *erop, ervoor*

Certain prepositions have a different specification with respect to this attribute depending on their function. Examples: *naar* (postform=both if used in a prepositional object, postform=pre if used as a directional preposition, cf. *hij heeft er naar gekeken* v. *\*hij is er naar gegaan*); *door* in its directional interpretation does not allow *er* (cf. *hij is door de tunnel gelopen* v. *\*Hij is er door gelopen*), in other functions it does (cf. *ik werd door zijn komst verrast* v. *ik werd er door verrast*).

This analysis, in which R-pronouns are always the result of a replacement of a specific pronoun, is actually incorrect. This can be shown by the following examples:

- (120) a Zij kijkt altijd naar die programma's<sub>*i*</sub>. Maar ik kijk er<sub>*i*</sub> nooit naar.  
 b Zij bekijkt die programma's<sub>*i*</sub> altijd. Maar ik bekijk ze<sub>*i*</sub> / \*het<sub>*i*</sub> nooit.  
 'She always watches those programs<sub>*i*</sub>. But I never watch them<sub>*i*</sub>.'

The R-pronoun *er* can be co-referential with the NP *die programma's*, as indicated in (120a) in which co-subscripting indicates co-reference. The pronoun *het*, however, cannot refer to such a (plural) NP. Instead the pronoun *ze* must be used. This implies that the pronoun *er* must not only be derived from *het*, but also from the pronoun *ze*. And in a similar manner one can show that *er* must also be derived in certain cases from the pronouns *hem* and *haar*. These facts cast doubts upon the whole analysis in which R-pronouns are derived from non-R-pronouns, since the original motivation for this analysis (i.e. the complementary distribution between *het* and *er*) is now weakened. We noticed these problems in Rosetta, but we did not change the grammar, because in the relevant phase we were not interested in the interpretation of pronouns. Of course, when the interpretation of pronouns is taken into account, the grammar must be adapted, e.g. by independently generating R-pronouns in PPs, or in some other way.<sup>9</sup>

<sup>9</sup>In fact, this change can be performed by simply adding a few abstract pronouns to the lexicon, and making one change in the rules dealing with R-pronouns. The R-pronouns in their prepositional function are then still introduced syncategorematically, but not as a replacement of *het* etc., but as a replacement of a number of abstract pronouns. This approach might be preferable for independent reasons, since considering the R-pronouns basic expressions in their prepositional use would lead to ambiguities in analysis (between the prepositional and the loca-

### 6.3.3 Locative R-pronouns

Locative R-pronouns are pronouns to indicate locations. Locative R-pronouns are dealt with as basic expressions. They are generated as a special adverbial phrase. They occur either VP-internal (if they are arguments of a verb such as *wonen*, *doorbrengen* ‘live, spend’ and certain uses of verbs such as *zitten*, *liggen*, *staan*, *zijn* ‘be’), or VP-external (if they are locative adjuncts). Special rules move these R-pronouns in certain cases to special positions. This will be discussed in more detail below.

### 6.3.4 Quantificational R-pronouns

Quantificational *er* must occur in clauses which contain a particular kind of headless NP in certain positions. In (121) quantificational *er* must occur in the context of the headless NP *drie* ‘three’:

- (121) Hij heeft \*(er) drie gezien  
 He has there three seen  
 ‘He has seen three of them’

I will describe the distribution of quantificational *er* in more detail below. Quantificational *er* is always introduced syncategorematically. It is assumed in the grammar that quantificational *er* is an element that has only purely grammatical functions. It does not have a meaning of its own. It is clearly the case that in a typical example containing quantificational *er* (e.g. 121) it must be assumed that the sentence contains an element with the semantics of a pronoun, but I assume that the semantics of the pronoun are not associated with *er*, but rather with an abstract element which occupies the head position of the object NP. The assumption that quantificational *er* itself is meaningless, accounts for the fact why the other R-pronouns cannot be used in a quantificational function. If quantificational *er* had meaning itself, then one would expect it to be in opposition with R-pronouns such as *daar* and *hier* etc. (as is the case in the locative and prepositional functions of R-pronouns): but this is not so. Sentences such as (122) are completely out, though it would be quite clear what semantics could be assigned to them, had they been grammatical:

---

tive function) which cannot be solved by S-PARSER. The abstract pronouns could be viewed as corresponding to the empty objects proposed as complements to P when an R-pronoun is present (see Bennis ([Bennis, 1986])). An additional disadvantage of generating R-pronouns independently is that it will lead to a less deterministic process in generation: when generating the pronoun, one cannot know whether an R-pronoun or a non-R-pronoun will be appropriate in the syntactic structure, so both possibilities must be attempted. When pronouns are actually interpreted and linked to their antecedents, one can readily imagine an analysis in which all pronouns are introduced syncategorematically, as the spelling out of an abstract representation (e.g. a variable) which also carries the information on the antecedents etc., so that the problem of non-determinism will not arise.

- (122) a \*hij heeft daar drie gezien  
           ‘He has seen three of those’  
       b \*hij heeft hier drie gezien  
           ‘he has seen three of these’  
       c \*hij heeft overal drie gezien  
           ‘he has seen three of all’

Being meaningless, quantificational *er* behaves like meaningless expletive *er*: no form other than *er* is possible. The assumption that quantificational *er* itself is meaningless will also play a role below, when it will be discussed how *er* can have several functions.

The distribution of quantificational *er* is accounted for in the following manner. First, NPs are generated with abstract elements as their heads. An attribute on the top node of such NPs marks these NPs as containing these abstract heads. The variables in a clause for which these NPs will be substituted will now be marked for these properties as well. Depending on the syntactic context, and the status of the NP with respect to countability, transformations apply to introduce quantificational *er*. These transformations introduce quantificational *er* and mark the relevant variables as being licensed if the following conditions hold:<sup>10</sup>

- the variable is count, and its corresponding NP contains an abstract head
- the variable is in *object* position, *preadv* position, *postsubject* position, *subject* position

In addition, there are special substitution rules which substitute NPs headed by abstract elements and which delete these abstract heads. These substitution rules apply only when the variable for which they substitute is marked as being licensed by a rule which has inserted quantificational *er*. This method makes it possible to account for the fact that the distance between quantificational *er* and the NP containing the abstract head can be indefinitely long (cf. **Hoeveel** *dacht jij dat Piet dacht dat hij er gezien had?* ‘How many did you think Pete thought he had seen?’), while the rule relating *er* to this NP can be kept local. Rules which are independently necessary (e.g. wh-movement) can then account for the fact that the distance can stretch indefinitely in the normal manner.

## 6.4 The Distribution of R-pronouns

### 6.4.1 General Discussion

In order to account for the distribution of R-pronouns, and to account for the fact that R-pronouns can have more than one function under certain conditions, I adopted the point of view that the specific distribution of R-pronouns should follow from the assumption of a special position (the R-position) into which R-pronouns must often be moved. The occurrence of several R-pronouns would then

<sup>10</sup>See Bennis ([Bennis, 1986]:283-284) for discussion.

either result in a clash (hence ill-formedness), or could be solved by amalgamating R-pronouns under specific conditions. The assumption of a special position into which R-pronouns are to be moved has been taken directly from Van Riemsdijk ([van Riemsdijk, 1978]), who uses this to account for the interaction of locative and prepositional R-pronouns. The idea was, that by extending this analysis to the other R-pronouns, the facts with respect to expletive and quantificational pronouns could be accounted for in a principled manner. Van Hout ([van Hout, 1986]) started to work out a proposal to deal with R-pronouns in Rosetta along these lines.

With regard to the problem of R-pronouns having several functions, the leading idea has been the following. Assume that certain R-pronouns are meaningless. Since such R-pronouns are meaningless, their presence in a structure cannot be necessary for semantic reasons. They might be required in a structure because of syntactic reasons. Suppose now that the syntactic function of these R-pronouns can be fulfilled by other elements in the structure, which must be present anyway. Then there would be no need for these elements to be present at all. Now, if such an R-pronoun which need not be present, either for syntactic or semantic reasons, cannot occur in a certain structure (e.g. because its position is occupied by some other element), then non-insertion of this R-pronoun will not lead to ungrammaticality, but instead to a grammatical sentence in which one element appears to fulfill more than one function.

These were the two leading ideas on the basis of which an analysis of R-pronouns was made. The first idea can be seen as an attempt to generalize Van Riemsdijk's ([van Riemsdijk, 1978]) R-position hypothesis and to give the R-position hypothesis a wider range of application.

The point of view concerning the multiple functions of R-pronouns was developed to replace the unsatisfactory analysis of this phenomenon by Bennis ([Bennis, 1980]), which was slightly improved in Bennis ([Bennis, 1986]), which was available to us in a preliminary version at the moment that these ideas were developed. The major problems concerning Bennis's analysis are (1) that no account was given why certain collapses of functions are and others are not possible; (2) that no clear connection was made with the R-position hypothesis (the R-position analysis is in fact rejected by Bennis), and (3) that the relation of Bennis's deletion rules to the principle of recoverability of deletion was unclear. These objections will be discussed in more detail.

Bennis ([Bennis, 1980]) basically suggests that all cases where two occurrences of *er* come together lead to a deletion of one of these *er*-pronouns. But Bennis ([Bennis, 1980]) does not discuss the case where prepositional and locative *er* come together, where such a deletion is not possible. Though in Bennis ([Bennis, 1986]) a more thorough discussion of these phenomena is added, his remarks are absolutely unsatisfactory. Bennis ([Bennis, 1980]:67-68, fn. 11) does discuss the fact that his analysis is incompatible with the R-position hypothesis, but he makes remarks there which are completely beside the point. Bennis ([Bennis, 1986]) essentially rejects the R-position hypothesis on theoretical grounds ([Bennis, 1986]:205), but he is unable to account for the facts that it is supposed to handle (and he even partially denies these facts (implicitly, on p. 188, example (45b), and explicitly, p.

208)). Of course, it makes no sense to quarrel about facts, but I have assumed the correctness of Van Riemsdijk's ([van Riemsdijk, 1978]) facts, and made an analysis which could deal with them. The objections to Bennis's analysis are independent of this factual issue.

Finally, Bennis assumes that occurrences of *er* are deleted. He does not discuss the principle of recoverability of deletion in any way, though it clearly might be relevant in this connection. The rule which deletes *er* may be a syntactic rule or a PF-rule (a phonological rule) in the framework he adopts. If the deletion is a syntactic rule, there are two possibilities: either *er* is deleted only when there is an antecedent, or *er* belongs to a class of designated elements which can be deleted without there being an antecedent. The first situation certainly does not hold. As a consequence, *er* must belong to a class of designated elements which can be deleted without having an antecedent. But if that is so, it is a complete mystery why deletion of locative *er* in the context of a prepositional *er* is impossible (or the other way around). Suppose now, that the deletion rule is a PF-rule, which Bennis assumes. Then recoverability of deletion can play no role. But again, the problem with locative and prepositional uses of *er* appears.

In addition to these objections, the deletion analysis, i.e. an analysis in which an occurrence of *er* is deleted if some other occurrence of *er* is adjacent because of a phonological prohibition against the occurrence of certain phoneme sequences in Dutch, as suggested by Bennis, has further drawbacks. Such an analysis must be rejected for a number of reasons. First, this approach does not explain why the relevant phoneme sequence is avoided by deletion, and not by some other mechanism, e.g. insertion of *d* (cf. *raar* v. *\*rarer* v. *raarder*), or by dissimilation (cf. *tover-en* v. *\*toveraar* v. *tovenaar*), or by yet other mechanisms (see Bennis [Bennis, 1986], p. 184). Second, there are direct counter-examples, e.g. (107b), and sentences such as *..daar er niet over gesproken mocht worden* 'since it was not allowed to talk about it' pointed out by Henk van Riemsdijk, in which the sequence *daar er* can occur without any problem. Third, there are examples where apparently two non-adjacent *ers* cannot co-occur in a sentence (cf. (107d)). Under the deletion analysis it must be assumed that a different mechanism is operative here. Fourth, there are also examples where *er* cannot occur though no other occurrence of *er* is present, but a different R-pronoun (e.g. *daar*). This is illustrated in (123):

- (123) Daar<sub>P</sub> werd (\*er<sub>X</sub>) over gesproken  
 There was there about talked  
 'People were talking about that'

The deletion analysis can be ameliorated by assuming that deletion occurs only in a specific configuration, when two occurrences of *er* are put together in specific adjacent positions. This perhaps better describes what Bennis assumes, though Bennis ([Bennis, 1986]) is rather vague on this. It would remove the second and perhaps the third objection. However, the other problems remain, and what is more, if the rules for putting 'together' R-pronouns are carefully formulated — including as a crucial ingredient the R-position — the deletion rule is superfluous

for most cases, as will be shown below. In addition, such an analysis avoids the other problems of the deletion analysis. Since there is no deletion, there need be no phonological or morphological justification for deletion, and the account also extends to cases where occurrences of different R-pronouns cannot co-occur.

### 6.4.2 Global Characterization of the Results

Before outlining the analysis actually implemented in Rosetta in detail, I will consider it in view of the central ideas described in the preceding section which formed its basis.

The account for the distribution of expletive, prepositional and locative R-pronouns indeed makes crucial use of the R-position proposed by Van Riemsdijk. So, in the Rosetta analysis the application domain of the R-position has been extended to expletive *er*. The behavior of quantificational *er*, and in particular its interaction with other R-pronouns could not be accounted for by making use of the R-position. Some very simple facts actually indicate that this is impossible in principle. Quantificational *er* can occur in sentences which contain other R-pronouns (including *er*) which arguably are in, or have passed through, the R-position. Examples are given in (124):

- (124) a Hij schreef er daar twee over  
           He wrote there<sub>Q</sub> there<sub>P</sub> two about  
           he wrote two (of them) about it  
       b Er zijn er twee  
           There<sub>X</sub> are there<sub>Q</sub> two  
           ‘There are two of them’

Thus, it is necessary to adopt a second position in addition to the R-position. The distribution of quantificational *er* can now no longer be accounted for by assuming that no position is available for it in certain cases. For quantificational *er* I continue to adopt a deletion analysis. Note that none of the objections raised earlier against a deletion analysis to account for the distribution of all R-pronouns is applicable here. I do not assume that the deletion rule required here is justified by phonological or morphophonological considerations (which is insufficient anyway). To my knowledge, there are no counter-examples to the deletion rule. Deletion does not apply if two occurrences of *er* are not adjacent (see 124b). And deletion does not apply if quantificational *er* is adjacent to an R-pronoun other than *er* (see 124a).

In connection with these considerations a number of factual questions also arise. I have assumed throughout, that a single occurrence of *er* can serve a quantificational function in combination with any other function at the same time. I based this judgement on sentences such as:<sup>11</sup>

<sup>11</sup>I represented the functions of *er* in these sentences by the familiar subscripts. These subscripts only hold if the words between brackets are absent. These have been added to suggest a possible antecedent for the empty head.



- (125) a Kwamen  $er_{XQ}$  slechts twee (mensen) ?  
 Came there only two (persons)  
 ‘Did only two (persons) come’
- b Hij heeft  $er_{QP}$  slechts twee (artikelen) over gelezen  
 He has there only two (articles) about read  
 ‘he read only two (\*articles) about it’
- c Hij bracht  $er_{QL}$  slechts twee (dagen) door  
 He spent there only two (days)  
 He spent only two (days) there
- d Hij kocht  $er_{QL}$  slechts twee (boeken)  
 He bought there only two (books)  
 ‘He bought only two (books) there’

In addition, more than two additional functions can be collapsed with the quantificational function in one occurrence of *er* if these two additional functions can be independently collapsed:

- (126) a Werden  $er_{XQP}$  slechts twee (artikelen) over gelezen?  
 Were there only two (articles) about read  
 ‘Were only two (articles) read about it’
- b Lagen  $er_{XQL}$  slechts drie (boeken)?  
 Lay there only three (books)  
 ‘Were there only three (books) lying there’

Quantificational *er* can even serve several quantificational uses:

- (127) Hoeveel (mensen) hebben hem  $er_{QQ}$  twee (boeken) gegeven?  
 How many (persons) have him there two (books) given  
 ‘How many persons gave him two (books)’

And these can occur in combination with other functions:

- (128) a Hebben  $er_{XQQ}$  twee (mensen) slechts drie (boeken) gekocht?  
 Have there two (persons) only three (books) bought  
 ‘Have three (persons) bought only three (books)?’
- b Hoeveel (mensen) hebben  $er_{PQQ}$  twee (artikelen) over gelezen?  
 How many (persons) have there two (articles) about read  
 ‘How many (persons) have read two (articles) about it?’
- c Hoeveel (mensen) hebben  $er_{LQQ}$  slechts twee (dagen) doorgebracht?  
 How many (persons) have there only two (days) spent  
 ‘How many (persons) spent only two (days) there?’
- d Hebben  $er_{XPQQ}$  slechts twee (mensen) drie (artikelen) over gelezen?  
 Have there two (persons) only three (articles) about read  
 ‘Did two (persons) read only three (articles) about it?’
- e Hebben  $er_{XLQQ}$  slechts twee (mensen) drie (dagen) doorgebracht?  
 Have there only two (persons) three (days) spent  
 ‘Did only two (persons) spend three (days) there?’

The judgements w.r.t these sentences vary for different people. I think that all of them are rather good, although in certain cases, especially where two numerals end up adjacent the results deteriorate. Maybe this can be accounted for by a focus clash. The numerals in an empty-headed NP must have focus: if two numerals are adjacent, assigning focus to both becomes very difficult. In general, the results of such sentences get better if the numeral of an empty-headed NP is modified by such words as *slechts*, *maar* ‘only’, which might be related to focus. Bech ([Bech, 1952]) gives *Hij vond er<sub>QL</sub> vijf* ‘He found five of them there’ and *Toen lagen er<sub>XQL</sub> vijf* ‘Then, there were five of them over there’ as well-formed. Bennis ([Bennis, 1986]:179, 181) gives variants of (125c) and (128d) as fully grammatical. Other people judge at least some of these sentences differently, though the situation is not clear. Model ([Model, 1991a]:304) stars *zij bezit er<sub>QL</sub> drie*, and marks *Willen er<sub>XQL</sub> echt maar drie wonen* with two question marks (though there is additional discussion ([Model, 1991a, 305])). I agree that the first sentence is not perfect (and I cannot account for this in my analysis), but the second one is fine. The judgements for all of these sentences are rather subtle. The analysis presented has been designed with the judgements of the facts described above in mind. Model ([Model, 1991a]) arrives at a different analysis, which is based on partially different judgements.

Turning to the second central idea on which the analysis outlined below is based, it has been worked out in the manner indicated for all functions of R-pronouns, though for the expletive function an additional mechanism was required. It is assumed that locative and prepositional R-pronouns are meaningful. These elements are introduced — just as all other meaningful arguments and adjuncts — by first introducing variables, and by applying substitution rules to substitute the actual phrases later in the derivation. At the point in the derivation where the R-pronouns come together in the R-position, the locative and prepositional functions are still present in the form of variables. Suppose that a rule deletes one of these variables: the substitution rule which must substitute the actual phrases cannot apply in that case, and the derivation blocks. Thus locative and prepositional R-pronouns cannot be deleted.<sup>12</sup>

Quantificational and expletive R-pronouns are meaningless.<sup>13</sup> Non-insertion of them will therefore never lead to an incompatibility between form and meaning. Non-insertion, however, might have consequences for the form itself. Therefore, non-insertion is possible only if the syntactic function of the R-pronoun is fulfilled by other elements.

Let us illustrate this with quantificational *er*. The syntactic function of quan-

<sup>12</sup>This account will not work if there are two occurrences of a variable with the same index. Then one of them can be deleted, and the substitution rule can apply to the other variable. This situation, however, can only arise if pronouns are actually interpreted, phenomena which have not been dealt with in the version of Rosetta described here. The problem can be avoided by ensuring that pronominal interpretation has already applied before the rules discussed here (as has been done in the Rosetta system for pronouns which must be interpreted in the sentence grammar), so that at most one variable is present in the structure in all cases. For efficiency reasons, the relevant rules have been formulated in such a way that they never delete a variable.

<sup>13</sup>See also Bennis ([Bennis, 1986]:202,212) for a similar point of view.

tificational *er* is to license an empty head of certain types of NPs. It can license empty heads in such NPs only if these NPs are in certain positions (e.g. in the object position, but not in the prepositional-object position). The relevant condition can be formulated as follows:

- (129) An empty head of a count NP must be licensed; it can be licensed by *er*, provided *er* is ‘close enough’.

For a more precise description of the notion ‘close enough’, see above (section 6.3.4). Suppose now, that we have partially derived a sentence, and we have arrived at the point where quantificational *er* must be inserted. I assume that quantificational *er* can be inserted only in one position. Let’s consider two examples. (130) illustrates the results of part of their derivation, at the point where quantificational *er* must be inserted:

- (130) a Hij [twee EN] kocht ‘he [two EN] bought’  
 b [Hoeveel EN] hem [twee EN] gegeven hebben  
 ‘[How many EN] him [two EN] given have’

In (130a) an empty head EN occurs in the direct-object position. It must be licensed, and it can be licensed by inserting quantificational *er*. The position where quantificational *er* must be inserted is unoccupied, so quantificational *er* can be inserted. Non-insertion will lead to an unlicensed EN, hence to ungrammaticality. Exactly the same holds in (130b). Here the head EN of the NP [*twee EN*] must be licensed, and it can be licensed by inserting quantificational *er*. In this latter example, however, the EN of the NP [*hoeveel EN*] occurs as well. This EN must also be licensed. A second occurrence of quantificational *er*, however, cannot be inserted, because the position for quantificational *er* is occupied. Insertion of quantificational *er*, however, is not necessary for semantic reasons, and actually not for syntactic reasons either in this case, since the other occurrence of quantificational *er* can license EN of the NP [Hoeveel EN]. Thus, here we have a clear example where on the one hand it is not necessary to insert some element, either for semantic or for syntactic reasons, and on the other hand it is not possible to insert this element: sentences in which this constellation of facts holds are well-formed (with respect to the relevant phenomenon).

A similar account can be given for expletive *er*, though here an additional mechanism must be assumed. The basic idea is the same as before: expletive *er* is meaningless. Therefore it need not be present in a sentence for semantic reasons. If it cannot be inserted in a specific configuration, and its syntactic function can be fulfilled by some other element, then it need not appear. This account appears to work well at first sight. Consider the following partially derived structures:

- (131) a EREC  $er_P$  iemand naar keek. lit. EREC there someone at looked  
 b EREC  $er_L$  iemand woont. lit. EREC there someone lives

EREC is the abstract expletive element occurring in the subject position introduced earlier. The R-position is occupied by some R-pronoun. Now the abstract

expletive EREC must be licensed (and deleted). Normally this is done by inserting *er* in the R-position, but this is not possible here, since this position is occupied. Suppose now that any R-pronoun in the R-position can license EREC. Then insertion of *er* is not necessary. EREC is licensed, and we can derive sentences such as *keek er iemand naar* and *woont er iemand*, and we correctly exclude sentences such as *\*keek er er iemand naar* and *\*woont er er iemand*.

The problem with this account is that it predicts that these R-pronouns will now normally behave as prepositional or as locative pronouns in such constructions. This, however, is not the case. Normally, prepositional and locative *er* cannot occur sentence-initially in a main clause (they cannot be topicalized):

- (132) a De man woont er  
           ‘The man lives there’  
       b \*Er woont de man  
           There lives the man  
       c De man keek er gisteren naar  
           The man looked there yesterday at  
           ‘The man looked at it yesterday’  
       d \*Er keek de man gisteren naar  
           There looked the man yesterday at

But in the constructions mentioned *er* can be preposed:

- (133) a Er keek iemand naar  
           There looked someone at  
           ‘Someone looked at it’  
       b Er woont iemand  
           There lives someone  
           Someone lives there

And when expletive *er* occurs on its own, it can be preposed as well:

- (134) a Er wordt gedanst  
           There is danced  
           ‘There is dancing going on’  
       b Er keek iemand naar dat programma  
           There looked someone at that program  
           ‘Someone was watching the program’  
       c Er woont iemand in Amsterdam  
           There lives someone in Amsterdam  
           ‘Someone lives in Amsterdam’

This might appear to indicate that in the sentences of (133) the prepositional and locative *er* have been dispelled by expletive *er*. That conclusion, however, is not correct, because in all cases where expletive *er* can be distinguished from non-expletive R-pronouns, the non-expletive R-pronouns ‘survive’:

- (135) a Daar keek iemand naar  
 There looked someone at  
 ‘Someone looked at that’  
 b Daar werd naar gekeken  
 There was at looked  
 ‘That was watched’

and a sentence such as (133a) cannot have the meaning of (135a).

Let us consider the facts in (134) in more detail. Expletive *er* can be preposed to a sentence-initial position. This is actually quite strange, since generally only clitics in subject position can be preposed, but I assumed that *er* does not occur in the subject position. It must be the case then, that the relation between EREC (which is in subject position) and *er* makes preposing possible. In Rosetta we mark an occurrence of an R-pronoun which licenses EREC with the marker X. ‘Subject’-preposing is formulated in such a way that it also applies to clitics in the R-position which are marked with X. This account now also explains immediately why the sentences in (133) are well-formed and why the sentences in (132) are ill-formed: the occurrences of *er* in (133) are marked with X, the occurrences of *er* in (132) are not. We see then, that licensing of EREC must be formulated in the following way: there is a rule, which applies optionally and does not take any position other than the R-position into account. This rule inserts *er* into the R-position, provided this position is unoccupied. EREC licensing applies after application of this rule, and EREC is licensed if the R-position is occupied by an R-pronoun. This R-pronoun is marked with X. Therefore, the R-pronoun which has licensed EREC can be topicalized, even if it is a clitic. I have called the phenomenon that R-pronouns get properties of expletive *er* when they license EREC *amalgamation*, because properties of expletive *er* and non-expletive R-pronouns appear to have amalgamated in one form.

One might ask why arbitrary R-pronouns can license expletive elements. I have no real answer to this question, but would like to point out that this is not a property specific to R-pronouns, but a property of a certain class of (locative?) adverbial phrases generally. This is illustrated in (136):<sup>14</sup>

- (136) a In het stadion werd gevoetbald lit. In the stadium was soccer-played  
 b Werd in het stadion gevoetbald? lit. Was in the stadium soccer-played  
 c Hij zei dat in het stadion gevoetbald werd  
 lit. He said that in the stadium soccer-played was  
 d \*Werd gevoetbald lit. Was soccer-played  
 e Er werd gevoetbald lit. There was soccer-played  
 f Werd er gevoetbald? lit. Was there soccer-played  
 g Hij zei dat er gevoetbald werd lit. He said that there soccer-played was

In (136a) an (abstract) expletive element is licensed without there being any R-pronoun present. The licenser is the locative PP *in het stadion*. This can be

<sup>14</sup>See Bennis ([Bennis, 1986]:215-216,225-226), from which the discussion here is derived.

seen in (136d): if the locative PP is left out, the sentence is ungrammatical. One cannot say that (136d) is excluded because the topic position must be occupied, since that is only true for certain sentence types, e.g. for main declarative clauses, but not for yes-no interrogative clauses. (136d), however, is ill-formed under any interpretation. The fact that the PP is in topic position is not essential either, as shown by the well-formed (136b,c). If we compare the behavior of this PP with respect to the abstract expletive element with the behavior of *er* (as in (136d,e,f,g)) we see exactly the same behavior. Of course, PPs such as *in het stadion* do not occupy the R-position at any moment in the derivation. Therefore it is possible to have expletive *er* in these constructions as well:

- (137) a In het stadion werd er gevoetbald  
 lit. In the stadium was there soccer-played  
 b Werd er in het stadion gevoetbald?  
 lit. Was there in the stadium soccer-played?  
 c Hij zei dat er in het stadion gevoetbald werd  
 lit. He said that there in the stadium soccer-played was

And there are subtle semantic (or pragmatic) distinctions between the sentences of (136) and (137). See Bennis ([Bennis, 1986]:226) for discussion.

Let us now consider the interaction with quantificational *er*. Quantificational *er* is introduced in a position other than the R-position. Therefore, it never blocks insertion of expletive *er*. When expletive *er* and quantificational *er* end up adjacent, one of them is deleted. Otherwise two occurrences of *er* survive.

This concludes the general characterization of the analysis of R-pronouns in Rosetta. In the next section we will look in more detail at a number of specific assumptions made, and illustrate how they account for the relevant facts.

## 6.5 The Assumptions in More Detail

In this section the relevant assumptions will be presented in more detail.

The first assumption is that there are two special positions in a sentence, the  $er_Q$ -position and the R-position. This has been discussed above. The R-position is postulated because the analysis is crucially based on and in fact extends Van Riemsdijk's ([van Riemsdijk, 1978]) analysis. The separate position for quantificational *er* has been justified above.

Second, the  $er_Q$ -position can contain one occurrence of *er* in its quantificational use. This assumption simply states that the special position assumed for  $er_Q$  can only be used by (one occurrence of)  $er_Q$ .

Third, all other occurrences of *er* must be moved into the R-position if they can get there. This is a crucial assumption. It is also a natural assumption. *Er* is a clitic, and clitics must usually occupy special positions in clauses. The clitic *er* is no exception (see also Model([Model, 1991a])). Note that a sentence such as (138a) might appear to be an exception to this rule:

- (138) a Hij beschuldigde de man ervan  
 ‘He accused the man of it’  
 b Hij beschuldigde de man van moord  
 ‘He accused the man of murder’

At first sight, sentence (138a) appears to have the same syntactic structure as (138b), with a direct object followed by a prepositional object. But there are several reasons to assume that their structures actually differ considerably. First, direct objects can, in principle, precede the R-position, and second, *er* in (138a) is certainly not inside the PP (despite the misleading orthographic conventions which might suggest this). This is shown by (139), where the adverb *gisteren* can intervene between *er* and *van*.

- (139) Hij beschuldigde de man er gisteren van  
 He accused the man there yesterday of  
 ‘He accused the man of it yesterday’

Second, there are certain restrictions on direct objects preceding the R-position. In particular, if the direct object is indefinite, it must be interpreted as a *specific* NP. This is illustrated by the examples in (140):

- (140) a Hij beschuldigde iemand er van moord  
 He accused someone there of murder  
 ‘He accused someone of murder there’  
 b Hij beschuldigde er iemand van moord  
 He accused there someone of murder  
 ‘He accused someone of murder there’  
 c ??Hij beschuldigde een man er van moord  
 He accused a man there of murder  
 ‘He accused a man of murder there’  
 d Hij beschuldigde er een man van moord  
 He accused there a man of murder  
 ‘He accused a man of murder there’

Sentences (140a) and (140b) are both grammatical, but they have a different meaning: *iemand* in (140a) must be interpreted as a specific NP, and (140b) must be interpreted as a non-specific NP. If an indefinite NP cannot be interpreted as a specific NP for some reason, it cannot occur in a position preceding the R-position. This is illustrated in (140c,d): the NP *een man* (with the unstressed article *een*) cannot be interpreted as a specific NP, so (140c) is deviant. The sentences in (140) contain locative *er*, but if this element is removed and the complement of *van* is replaced by *er*, exactly the same pattern of facts arises:

- (141) a Hij beschuldigde iemand ervan  
 He accused someone there-of  
 ‘He accused someone of it’  
 b Hij beschuldigde er iemand van  
 He accused there someone of  
 ‘He accused someone of it’  
 c ??Hij beschuldigde een man ervan  
 He accused a man there-of  
 ‘He accused a man of it’  
 d Hij beschuldigde er een man van  
 He accused there a man of  
 ‘He accused a man of it’

This clearly shows that *er* in such sentences must occupy the R-position, and not be part of the PP of its governing preposition. Additional evidence can be derived from the following contrast:

- (142) a Hij legde het boek er vlak naast  
 He put the book there right next-to  
 ‘he put the book right next to it’  
 b ?Hij legde het boek vlak ernaast  
 He put the book right there next-to  
 ‘he put the book right next to it’

The sentence in which the pronoun *er* escaped from the PP (142a) is much more natural than sentence (142b), though the latter one is not completely ungrammatical for all native speakers.<sup>15</sup>

Note that expletive *er* also occupies the R-position, not the subject position, as stated before. This assumption plays a crucial role in accounting for the distribution of R-pronouns, as illustrated in detail below. There are some (rather subtle) facts which appear to indicate that expletive *er* and non-expletive *er* occupy different positions.

The R-pronoun *er* usually occurs to the right of all kinds of other clitics, cf. *\*..omdat hij er zich waarschijnlijk voor schaamt v. ..omdat hij zich er waarschijnlijk voor schaamt* ‘because he probably is ashamed of it’, *\*..omdat hij er ’m zag v. ..omdat hij ’m er zag* ‘because he saw him there’, etc.<sup>16</sup>

However, expletive *er* must occur to the left of *zich*. This can be seen in the following examples:

<sup>15</sup>In (142a) *er* can but need not be outside PP, according to Van Riemsdijk ([van Riemsdijk, 1978]).

<sup>16</sup>Many people, however, accept the order *er zich* in all sentences. For these people, some of the arguments will probably not be valid.



- (143) a ..omdat er zich iets voorgedaan heeft  
 .. because there itself something occurred has  
 ‘..because something occurred (there)’  
 b ..omdat zich er iets voorgedaan heeft  
 ..because itself there something occurred has  
 ‘..because something occurred there’

In the first example *er* can be interpreted as an expletive and as an expletive-locative, in the second example it can be interpreted only as a locative. These intuitions are rather subtle, but they are confirmed by the following observations. The contrast between the sentences in (144) can be accounted for if *er* in the second sentence must be interpreted as an expletive, so that it requires an indefinite subject, while *er* in the first sentence has only a locative interpretation and imposes no definiteness restrictions on the subject.

- (144) a Dit heeft zich er gisteren voorgedaan  
 This has itself there yesterday occurred  
 ‘This occurred there yesterday’  
 b \*Dit heeft er zich gisteren voorgedaan  
 This has there itself yesterday occurred

If we replace the definite subject *dit* by the indefinite subject *wat* both sentences are grammatical again, though they differ subtly in meaning: *wat* must be interpreted as specific in (145a) and as non-specific in (145b):<sup>17</sup>

- (145) a Wat heeft zich er gisteren voorgedaan?  
 What has itself there yesterday occurred  
 ‘What occurred there yesterday?’  
 b Wat heeft er zich gisteren voorgedaan?  
 What has there itself yesterday occurred  
 ‘What occurred (there) yesterday?’

A second observation which confirms my intuitions with regard to the first sentence pair comes from contrasting the following two sentences.

- (146) a ??Wat heeft zich er daar gisteren voorgedaan?  
 What has itself there there yesterday occurred  
 ‘What occurred there over there yesterday?’  
 b Wat heeft er zich daar gisteren voorgedaan?  
 What has there itself there yesterday occurred  
 ‘What occurred over there yesterday’

In the first sentence, *er* can be interpreted only as a locative. Because *daar* can also only be interpreted as a locative, the sentence is deviant because it contains

<sup>17</sup>And, of course, the locative interpretation of *er* is obligatory in (145a), but optional in (145b).

two locatives<sup>18</sup>. In the second sentence, *er* can be interpreted either as an expletive or as an expletive-locative. Since in this sentence *daar* is also a locative, the interpretation of *er* as an expletive-locative is impossible, but there is also an interpretation of *er* as an expletive. Under this latter reading the sentence is not deviant at all.

Note, however, that every use of *er*, provided one of its uses is expletive, shows the same distribution. Assuming a separate position for expletive *er* would destroy the account for the fact that the occurrences of several *ers* in one sentence is possible only in a limited number of circumstances, and it would destroy the account of amalgamation phenomena. Thus, these facts cannot be seen as evidence in favor of a special position for expletive *er*. The distributional phenomena described here have been accounted for by a late re-ordering rule which interchanges *zich* and the R-position. It appears that only expletive *er* is subject to this rule. It is impossible for *daar*, *hier*, *waar* (cf. *\*wat deed daar zich voor?*, *\*wat deed hier zich voor?*, *\*wat deed waar zich voor?*) ‘What occurred there/here/where?’.

Returning to the assumption proposed, it is clear that it cannot be implemented directly, since the part *if it can get there* presupposes a correct description of the language. In the actual implementation, the set of configurations from which *er* can get into the R-position have been enumerated. The formulation as given is nevertheless useful as a short and concise description of the relevant generalization. It predicts that there will be a correlation between several facts, e.g. the separability of *er* from its preposition (e.g. when there is an intervening adverb between the R-pronoun and the preposition it is a complement of) is predicted to correlate with its distribution, e.g. whether it can co-occur with other occurrences of *er* in one clause. This can be illustrated by the following examples:

- (147) a ?Hij keek gisteren ernaar  
He looked yesterday there-at
- b Hij keek er gisteren naar  
He looked there yesterday at  
‘He looked at it yesterday’
- c Er keek iemand (\*er) gisteren naar  
There looked someone (\*there) yesterday at  
‘Someone looked at it yesterday’
- d Hij keek erna naar dat programma  
He looked there-after at that program  
‘He looked at the program afterwards’
- e \*Hij keek er gisteren na naar dat programma  
He looked there yesterday after at that program
- f Er keek gisteren \*(er)na iemand naar dat programma  
There looked yesterday \*(there)after someone at that program  
‘Someone looked at that program afterwards’

<sup>18</sup>Two locatives in one sentence is not impossible, but additional conditions must be satisfied in such sentences. These additional conditions are not satisfied in the example sentences given.

In (147a) *er* is a complement of the preposition *naar*. As shown in (147b) this *er* and its governing preposition can be separated by an adverb. This means that *er* can reach the R-position. The marginal status of (147a) can be ascribed to the fact that adverbs such as *gisteren* preferably follow the R-position. Because prepositional *er* can reach the R-position, it must get there (by the assumption proposed here), and as a consequence it blocks the occurrence of other *ers* in the same clause: hence (147c) is out. In (147d) *er* is a complement of the preposition *na*. This occurrence of *er* and its preposition cannot be separated by an adverb such as *gisteren*, witness (147e), and this appears to be a general property of temporal PP adjuncts. This means that *er* cannot reach the R-position, and as a consequence it does not block other occurrences of *er* in the same clause (see (147f)).

We now return to the main theme of this section, a description of the relevant assumptions to account for the distribution of R-pronouns. The fourth assumption is that R-pronouns in their prepositional function must be moved into this R-position if they can. This is a more controversial assumption. Violating it yields somewhat marginal sentences. Some examples which can test the validity of this assumption are given in (148):

- (148) a  $Er_X$  werd daar<sub>P</sub>naar gekeken  
 b  $Er_X$  werd hier<sub>P</sub>naar gekeken  
 c  $Er_X$  werd ergens<sub>P</sub> naar gekeken  
 d  $Er_X$  werd nergens<sub>P</sub> naar gekeken  
 e  $Er_X$  werd overal<sub>P</sub> naar gekeken  
 There was there/here/somewhere/nowhere/everywhere at looked  
 f \*Daar<sub>P</sub> werd  $er_X$  naar gekeken  
 There was there at looked

According to my intuitions, all these sentences (under the intended interpretation) are ill-formed.<sup>19</sup> In Rosetta, assumption (D) has been incorporated, from which this follows immediately, because, in essence, the structure is equivalent to the structure required for (148f). However, I found later that some people accept the sentences (148a-e).

Another issue that is possibly relevant in this connection (which has been implemented only partially) concerns the behavior of R-pronouns in prepositional complements to adjectives. PP-complements to adjectives can either precede or follow the adjective, depending on the specific adjective. If a PP follows the adjective, and the PP contains an R-pronoun, then the structure is ill-formed or very marginal if the R-pronoun can be moved into the R-position. In that case, the R-pronoun must precede the adjective:

<sup>19</sup>Bennis ([Bennis, 1986]:214) marks the comparable sentence *..dat er niemand daar op rekende* lit. *..that there nobody there on counted*, ‘..that nobody counted on that’ with two question marks. However, he gives a completely different (and, in my view, incorrect) analysis of this sentence, which has been discussed above. Note that the sentences (148a-e) are well-formed when the the second R-pronoun is interpreted as a locative, and the first as a mixed expletive-prepositional.

- (149) a \*..omdat hij tevreden er/ daar/ hier/ nergens/ ergens/ overal over is  
 ..because he satisfied there / there/ here/ nowhere/ somewhere/ every-  
 where about is
- b \*..omdat hij blij er/ daar/ hier/ nergens/ ergens/ overal mee is  
 .. because he glad there/ there/ here/ nowhere/ somewhere/ everywhere  
 with is
- c ..omdat hij er/ daar/ hier/ nergens/ ergens/ overal tevreden over is  
 ..because he there/ there/ here/ nowhere/ somewhere/ everywhere satis-  
 fied about is
- d ..omdat hij er/ daar/ hier/ nergens/ ergens/ overal blij mee is  
 .. because he there/ there/ here/ nowhere/ somewhere/ everywhere glad  
 with is

It is not possible to account for this by assuming that the PP starts out in a pre-adjectival position and that the preposition is moved to a post-adjectival position, e.g. by the rule *ERVANDAAN* proposed by Van Riemsdijk ([van Riemsdijk, 1974],[van Riemsdijk, 1992]), though such an operation might be necessary independently.

The situation is different if the R-pronoun cannot reach the R-position, as shown in (150):<sup>20</sup>

- (150) a Tevredener erover kun je niet zijn  
 More satisfied thereabout can you not be
- b \*Er tevredener over kun je niet zijn  
 There more satisfied about can you not be  
 ‘More satisfied about it you cannot be’

Though the PP-complements to adjectives such as *blij*, *tevreden* can occur in a pre-adjectival position (cf. *..dat hij met dat bier blij is; dat hij over zijn medewerkers tevreden is*), it is clear that the pre-adjectival PPs are not in a complement position of the adjective, since the PP and the adjective can be separated by all kinds of material such as negation and sentential adverbs:

- (151) a Hij was over zijn medewerkers waarschijnlijk niet erg tevreden  
 ‘He was probably not very content about his employees’
- b Hij was met dat bier waarschijnlijk niet erg blij  
 ‘He was probably not very glad with that beer’

In addition, these PPs are not possible in pre-adjectival position in correlative coordinations with the adjectives given, though they are perfect in a post-adjectival position, which suggests again that the pre-adjectival PPs do not occupy a complement position of the adjective:

<sup>20</sup>These facts were pointed out to me by Henk van Riemsdijk.

- (152) a ??zowel blij als over zijn medewerkers tevreden  
 b    zowel blij als tevreden over zijn medewerkers  
       ‘both glad and satisfied with his employees’  
 c    ??zowel tevreden als met het bier blij  
 d    zowel tevreden als blij met het bier  
       ‘both content and glad with the beer’

I conclude that such PPs are probably preposed by independent rules. Assuming this to be correct, it remains to account for the position of the R-pronouns which must precede the adjectives. But this would follow immediately from assumption (D) proposed here. Note that all facts indicate that the R-pronoun preceding the adjective is actually in the R-position. In such structures no additional *er* can occur, and amalgamation takes place:

- (153) a \* $Er_X$  was daar<sub>P</sub> niemand blij mee  
       There was there noone glad with  
 b    \* $Hij$  was  $er_L$  daar<sub>P</sub> erg blij mee  
       he was there there very glad with  
 c    Daar<sub>XP</sub> was niemand blij mee  
       There was no one glad with  
       ‘No one was glad with that’

Turning to the fifth assumption, it states that prepositional R-pronouns cannot be preposed to the sentence-initial position directly, but must move via this R-position. This is an essential ingredient of the analysis proposed by Van Riemsdijk ([van Riemsdijk, 1978]). See section 6.6 for a more detailed discussion of the relevant facts and van Riemsdijk’s analysis. Van Riemsdijk tries to derive this assumption from independently required conditions (Subjacency), but in the actual system the assumption has simply been stipulated.

Sixth, if more than one R-pronoun ends up in the R-position, the derivation blocks. However, it is possible to give the R-pronoun in the R-position an additional syntactic function (amalgamation).

Seventh, amalgamation of two R-pronouns is possible only when one of them is meaningless (expletive or quantificational *er*), and has only a syntactic function: meaningful elements cannot remain unrealized. If an element is unrealized, its syntactic function must be fulfilled in some other way. See section 6.4.2 and the next section for illustration.

Eighth, if two occurrences of *er* (one in the  $er_Q$ -position and one in the R-position) remain adjacent after application of the transformations dealing with Verb Second, deletion transformations delete one of these occurrences. This is the residue of the deletion analysis proposed by Bennis. It must be emphasized again, that the objections against Bennis’s deletion analysis do not apply here.

For ease of reference, I will summarize the assumptions described above, and assign letters to them:

- A There are two special positions in a sentence, the  $er_Q$ -position and the R-position.

- B The  $er_Q$ -position can contain one occurrence of  $er$  in its quantificational use.
- C All other occurrences of  $er$  must be moved into the R-position if they can get there.
- D R-pronouns in their prepositional function must be moved into this R-position if they can get there.
- E Prepositional R-pronouns cannot be preposed to the sentence-initial position directly, but must move via this R-position.
- F If more than one R-pronoun ends up in the R-position, the derivation blocks. However, it is possible to give the R-pronoun in the R-position an additional syntactic function (amalgamation).
- G Amalgamation of two R-pronouns is only possible when one of them is meaningless (expletive or quantificational  $er$ ), and has only a syntactic function.
- H If two occurrences of  $er$  (one in the  $er_Q$ -position and one in the R-position) remain adjacent after application of the transformations dealing with Verb Second, deletion transformations delete one of these occurrences.

## 6.6 Illustration

The assumptions outlined have been implemented. They make it possible to account the facts mentioned above (and many more). I will briefly discuss the examples from the introductory section.

In (105a) expletive  $er$  and prepositional  $er$  are combined in one sentence. Both occurrences of  $er$  must be put into the R-position (cf. (C)), but that is not possible (cf. (F)). Amalgamation can save the structure. It is possible because expletive  $er$  is meaningless. Therefore, it need not appear in the structure for semantic reasons. Its syntactic function can be fulfilled by prepositional  $er$  in the R-position, which is marked with  $X$ . One can see that there is real amalgamation, because the remaining  $er$  can be topicalized, which is only possible if it has inherited properties of expletive  $er$  (expletive  $er$  can, but prepositional  $er$  cannot be topicalized).

In (105b) quantificational and prepositional  $er$  are combined. This is possible, since there are two positions, one for quantificational  $er$ , and one for other  $ers$  (cf. (A),(B),(C)). Since the two occurrences of  $er$  end up adjacent, they are subject to rule (H), hence one of them is deleted.

In (105c) locative and expletive  $er$  are combined. This is possible in the same way that prepositional and expletive  $er$  can be combined (see (5)).

In (105d) expletive, prepositional and quantificational  $er$  are combined. Expletive and prepositional  $er$  can be combined in the manner described for (105a). Quantificational  $er$  can occur in the sentence independently (cf. (A),(B)). Since the two occurrence of  $er$  end up adjacent, they are subject to the rule described in (H), so one of them is deleted.

In (107a) expletive and quantificational *er* are combined in one sentence. This is possible since there are positions to put these elements in (see (A,B,C)). Since they do not end up adjacent, they cannot be subject to the rule described under (H), hence two occurrences of *er* remain in the sentence.

Example (107b) has been discussed extensively above.

In (107c) we have the same structure as in (107b), though here the two occurrences of *er* end up adjacent, so the rule described in (H) must apply and delete one of them. Since this has not been done, the sentence given is ill-formed.

In (107d) prepositional *er* must be in the R-position (cf. (C)), because it can get there (cf. *hij keek er gisteren naar*). Expletive *er*, however, must also be in this position. So, only one of these words can occur in this sentence. Since there are two occurrences of *er* in this sentence, the sentence is ruled out.

The assumptions made also account for the facts pointed out and analyzed by Van Riemsdijk ([van Riemsdijk, 1978]), in a way similar to Van Riemsdijk's analysis. The R-pronoun *er* and the R-pronoun *waar* can both be used either as a locative adverbial (cf. (154a,b)), or as a complement to a preposition (cf. (154c,d)):

- (154) a Hij woont er 'He lives there'  
 b Waar woont hij 'Where does he live'  
 c Hij keek er gisteren naar 'He looked at it yesterday'  
 d Waar keek hij gisteren naar? 'What did he look at yesterday?'

When both R-pronouns are combined in one sentence (as in (155)), one would expect the sentence to be ambiguous, but it is not: in (17) the R-pronoun *waar* can only be interpreted as a locative adverbial and not as complement to the preposition *naar*.

- (155) Waar keek hij er gisteren naar?  
 'Where did he look at it yesterday?'

In the analysis presented, this fact can be accounted for in the following manner: if one tries to derive the sentence in its incorrect interpretation, locative *er* must be in the R-position (according to assumptions (A) and (C)). Prepositional *waar* must be preposed, but cannot be preposed in one step (see assumption (E)). Hence, it also gets into the R-position. Since this position is already filled (by locative *er*), the derivation blocks. Amalgamation is not possible, since both elements are meaningful (in accordance with assumption (G)). Under the other interpretation locative *waar* is simply preposed, and prepositional *er* is put into the R-position (in accordance with assumption (C)), so the sentence can be derived in this way. As a consequence, the system translates the sentence correctly as *Where did he look at it yesterday?* and not as *What did he look at there yesterday?*

The same holds if we take the R-pronoun *daar* instead of *waar*. The system also deals with it adequately, cf. (156), in a similar manner:

- (156) a Daar keek hij er gisteren naar  
 Over there looked he there yesterday at  
 ‘He looked at it there yesterday’  
 b Hij keek er daar gisteren naar  
 He looked there over there yesterday at  
 ‘He looked at it there yesterday’

In order to appreciate more fully the complexity of these phenomena, I would like to point out that these examples show that neither the mutual relative order of the R-pronouns *daar* and *er* nor their distance to the preposition *naar* is a relevant factor in determining these interpretations: in (156a) *daar* precedes *er*, in (156b) *daar* follows *er*; in (156a) *er* is closer to *naar* than *daar*, in (156b) the reverse situation holds.

## 6.7 Concluding Remarks

In this section I will first outline a number of remaining problems and then present my conclusions.

Though a certain degree of succes has been achieved in describing the complex syntax of R-pronouns in Dutch, there is still a number of remaining problems.

One area where problems still remain concerns the occurrence of R-pronouns fulfilling a multiple prepositional function. The analysis of R-pronouns developed in this chapter does not allow R-pronouns with a multiple prepositional function, and this is correct for many examples. The sentences in (157) contain R-pronouns with a multiple prepositional function. All these sentences are ill-formed, in accordance with the analysis developed in this chapter:

- (157) a \*Hij keek er gisteren mee naar  
 Hij looked there yesterday with at  
 b \*Waar keek hij gisteren mee naar?  
 Where looked he yesterday with at?  
 c \*Hij keek daar gisteren mee naar  
 He looked there yesterday with at  
 d \*De mensen waar hij gisteren mee over sprak  
 The people where he yesterday with about spoke

However, there are also well-formed sentences which contain R-pronouns fulfilling a multiple prepositional function:<sup>21</sup>

---

<sup>21</sup>Sentence (158b) is judged ill-formed by [Model, 1991a, 320], but I find it acceptable.



- (158) a Hij ging er gisteren mee naar toe  
 he went there yesterday with to  
 b We moesten er toen mee tegen op kletteren  
 We had-to there then against climb  
 c Er stond een artikel over in  
 There was an article about in  
 d Hij zette er een artikel over in  
 He put there an article about in

These examples appear to be problems for the analysis. Though these constructions require further research, I have the impression that they should not be accounted for in the same way that other R-pronouns fulfilling multiple functions have been analyzed. I suspect that these examples should be analyzed in a completely different manner, so that the analysis for R-pronouns with multiple functions developed here remains unaffected. Of course, it remains to be shown if this is the case, but I will present two observations which point in this direction.

First, examples of constructions containing R-pronouns with multiple prepositional functions are possible for the R-pronoun *er*, but not for other R-pronouns:

- (159) a \*Hij ging daar gisteren mee naar toe  
 He went there yesterday with to  
 b \*Waar ging hij gisteren mee naar toe?  
 Where went he yesterday with to  
 c \*Daar stond een artikel over in  
 There stood an article about in  
 d \*Het tijdschrift, waar een artikel over in stond  
 The journal, where an article about in was  
 e \*De schep, waar hij gisteren mee op af ging  
 The shovel, where he yesterday with towards went  
 f \*Hij zette daar een artikel over in  
 he put there an article about in

Second, it appears that the relevant construction is only possible if one of the prepositional phrases occurs inside a complement of the verb or in a circumpositional PP. The well-formed examples given contain either a PP inside an NP complement, or a PP containing a circumposition. No such restrictions hold in other cases where R-pronouns coincide. Perhaps it is possible to analyze circumpositional PPs as containing a post-position which takes a PP as its complement, so that we are dealing with PPs inside complements to the verb in both cases, though Van Riemsdijk ([van Riemsdijk, 1978]) argues that this is only true for some circumpositional PPs.

These properties, if real, are completely unexpected if these sentences are analyzed in the same way as the sentences containing other R-pronouns fulfilling multiple functions. This suggests that a completely different phenomenon is involved here. Further investigation will have to reveal whether these preliminary impressions can be upheld or not.

Sentences such as in (160) from [Model, 1991a, 314-5] form more serious problems for the analysis sketched:

- (160) a Waar moesten we er toen mee tegen op kletteren?  
Where had-to we there then with against climb?  
b Waar moesten we hier toen mee tegen op kletteren?  
Where had-to we here then with against climb

In these sentences two R-pronouns with prepositional function outside their PPs occur. Such sentences cannot be generated in the analysis proposed, since both prepositional R-pronouns must pass through the single R-position. Clearly, if these sentences are well-formed, a revision of the analysis is in order, perhaps along the lines of [Model, 1991a].

A third area where problems remain concerns an example adduced and analyzed by Van Bart and Kager ([van Bart and Kager, 1984]) and [Model, 1991a]:

- (161) \*Waar<sub>L</sub> heeft zij hier<sub>P</sub> vaak over gesproken  
Where has she here often about talked

If the sentence given is indeed ill-formed (the judgement is rather subtle), it forms a problem for my analysis, which allows the sentence.

An additional unsolved problem concerns the distribution of expletive *er*. It has been pointed out by [Drewes et al., 1984] that expletive *er* cannot occur in certain constructions (topicalization, relativization) where a PP has been preposed:

- (162) a Daarover wordt (\*er<sub>X</sub>) niet gepraat  
About that is (\*there) not talked  
'One does not talk about that'  
b Het onderwerp, waarover (\*er<sub>X</sub>) niet gepraat wordt  
The subject, about which (\*there) not talked is  
'The subject about which one does not talk'

If the PP contains the interrogative pronoun *waar*, however, expletive *er* can occur:

- (163) Waarover wordt (er) niet gepraat?  
About what is (there) not talked?

The same set of facts holds for locative R-pronouns, when preposed:

- (164) Waar woont (er) iemand met een Jaguar?  
Where lives (there) someone with a Jaguar?  
Daar woont (\*er) iemand met een Jaguar?  
There lives (\*there) someone with a Jaguar  
Het huis, waar (\*er) iemand met een Jaguar woont  
The house, where (\*there) someone with a Jaguar lives

And the facts hold for PPs not containing or consisting of R-pronouns as well:

- (165) In welke straat woont (er) iemand met een Jaguar?  
 In which street lives (there) someone with a Jaguar?  
 In deze straat woont (\*er) iemand met een Jaguar  
 In this street lives (\*there) someone with a Jaguar  
 De man, in wiens huis (\*er) iemand met een Jaguar woont  
 The man, in whose house (\*there) someone with a Jaguar lives

Henk van Riemsdijk pointed out that definiteness of the locative phrase might be a relevant factor in these constructions, but I have no explanation for these facts, which require further investigation.

A final area relates to implementing a practical system which deals with facts where there is a certain amount of speaker variation. Apart from the core facts analyzed here, a practical system must have measures to deal with speaker variation and with typical errors made in real texts. Of course, this problem is not limited to R-pronouns, but concerns all constructions. In the current system these problems have been dealt with only marginally.

Despite these problems, it can be stated that quite a succesful analysis of the complex syntax of R-pronouns in Dutch has been given and implemented. The approach adopted to deal with R-pronouns is elegant and systematic. A fairly complex array of facts has been accounted for by adopting only a few simple assumptions. The analysis avoids the ill-motivated deletion analysis presented by [Bennis, 1986], who deals with the same array of facts in a theoretical framework, and replaces it with a more principled account of why certain R-pronouns can and others cannot coincide. In addition, the analysis given is also empirically more correct. Therefore, my analysis provides a description which is superior to Bennis's analysis in many respects. The analysis given in [Model, 1991a] developed much later is approximately equivalent (perhaps superior) with respect to the basic facts concerning the distribution of R-pronouns, but has nothing to say about R-pronouns fulfilling multiple functions. I know of no analyses or real implementations of the syntax of R-pronouns in computational frameworks which cover the same amount of facts.

R-pronouns play an important role in a wide variety of constructions in Dutch: they appear as locative arguments and adjuncts, as complements of prepositions and in constructions in which headless NPs occur. They play a role in existential sentences, in impersonal passives and in several sentence types to induce non-generic and/or non-specific readings of subject NPs. I have argued that in most cases (all but the case of locative R-pronouns) the apparently associated semantics should be associated with other elements or rules. Most rules accounting for the syntax of R-pronouns are not associated to a meaning and are not relevant for translation. This makes it possible to account for the syntax of R-pronouns and for their interaction by fairly simple and general rules which are valid for almost all uses of R-pronouns.



## Chapter 7

# Concluding Remarks

### 7.1 Conclusions

The main purpose of this dissertation was to attempt to incorporate syntactic generalizations in the compositional framework of controlled M-grammar. I believe that I have shown that compositional M-grammar forms an excellent framework for achieving this goal.

In chapter 3 it was shown that the controlled M-Grammar formalism makes it possible to find a proper balance between purely syntactic requirements and the requirements of a compositional grammar in which form and meaning must be associated in a fairly direct way. Though the strong semantic bias of a compositional grammar could easily lead to syntactically inadequate analyses, there is compensation by allowing syntactic transformations: rules which have the identity operation as their meaning. In addition, we found that it was necessary to allow that rules which are normally associated with a meaning are applied without this associated meaning, to form structures which are introduced syncategorematically. In this way it is possible to express syntactic generalizations adequately within this framework.

The combination of meaningful rules and syntactic transformations, each of which may be relatively simple, increases the amount of rule interaction and makes the relation between form and meaning somewhat indirect in the following sense: many rules which account for form differences should be transformations in a syntactically adequate description, even though the form differences appear to correlate with meaning differences. Other rules, which often apply to more abstract representations of the sentence and therefore do not have directly visible formal effects in the surface, account for the meaning differences.

Chapter 4 discussed how predicate-argument relations are dealt with in the Rosetta grammars. It was shown that the compositional nature of controlled M-grammar has immediate and far-reaching consequences for the treatment of predicate-argument relations: the very design of the grammar immediately derives a variant of the  $\theta$ -criterion in a very strict way; it implies that only the notion

'semantic argument' is relevant, and it has immediate consequences for the treatment of optional arguments. The treatment of arguments within the framework has been dealt with extensively. To avoid the undesirable consequences of this approach for a certain class of phrases, it has been proposed that this class should not be dealt with as arguments, but as *bound adjuncts*, which are intermediate between arguments and adverbials (free adjuncts). The treatment of small clauses, which play an important role throughout the grammar, has been illustrated in detail.

Chapter 5 illustrated how a number of constructions have been incorporated in the Rosetta grammars. The treatment of passivization (section 5.2), Verb Second in Dutch (section 5.3), unbounded dependencies (section 5.4) and crossing dependencies in Dutch (section 5.5) have been discussed. For all these constructions a particular existing analysis has been adopted as a starting point, and an attempt has been made to incorporate these analyses in the framework adopted here. This required adaptation of the original analyses in certain cases, though in general, the existing analyses could be incorporated fairly directly.

The resulting analysis for passives is one in which a conglomerate of rules, only one of which is particular to passives, derives passive constructions. Most of the rules postulated are also used in the derivation of completely different constructions, and are for this reason not construction-specific. The operations which are common to the derivation of several constructions have been identified and isolated (factored out) in separate rules, so undesirable duplication has been avoided.

The analysis of Verb Second in Dutch incorporates the analysis developed originally by Den Besten ([den Besten, 1983]) in a fairly direct manner. It was argued that this analysis, which requires a movement rule, is superior to alternative analyses developed in the GPSG and HPSG frameworks. The analysis of Verb Second as developed by Den Besten states that the application of the relevant rule is conditioned by syntactic structure: no other factors are relevant and there is no direct link with semantic properties. Though the position of the finite verb in a sentence may appear to have a direct influence on the interpretation of a sentence, the rule accounting for the position of the verb has been formulated as a transformation which is sensitive to syntactic structure only. It plays a role in the derivation of many different kinds of sentences and is therefore not construction-specific. The fact that Den Besten's analysis can be incorporated directly forms an illustration of one of the main virtues of the controlled M-grammar framework: the transformations allowed in this framework make it possible to factor out common properties in separate rules.

For unbounded dependencies, an analysis is adopted in which phrases are actually moved from one position in the tree to another position in a successive cyclic manner, as proposed originally in [Chomsky, 1973]. It was argued that such a movement analysis is superior to analyses using feature-transportation mechanisms to describe unbounded dependencies.

The movement rule involved in describing these unbounded dependencies is also a transformation which is sensitive to syntactic structure only. It is not associated to a specific meaning aspect. This makes it possible to use this rule in

all constructions where unbounded dependencies play a role, so that this rule is also not construction-specific.

Finally, I discussed the treatment of crossing dependencies in Dutch. I adopted the analysis by [Evers, 1975] as a starting point, and incorporated it with a few minor modifications into the grammar. It was shown that the controlled M-grammar formalism makes it possible to incorporate this syntactic analysis rather directly, so the relevant syntactic generalizations concerning the distribution of verbs in Dutch can be expressed with a minimum of machinery.

In chapter 6 quite a successful analysis of the complex syntax of R-pronouns in Dutch has been given. The approach adopted to deal with R-pronouns is elegant and systematic. A fairly complex array of facts has been accounted for by adopting only a few simple assumptions. It was argued that the analysis provides a description which is superior to other existing analyses.

It was shown that R-pronouns play an important role in a wide variety of constructions in Dutch. I have argued that in most cases (apart from locative R-pronouns) the apparently associated semantics should be associated with other elements or rules. Most rules accounting for the syntax of R-pronouns are not associated to a meaning and are not relevant for translation. This makes it possible to account for the syntax of R-pronouns and for their interaction by fairly simple and general rules which are valid for almost all uses of R-pronouns.

These conclusions can be summarized in the following way: (1) the framework of controlled M-grammars supplies — due to its compositional nature — a firm framework to deal with syntactic phenomena which are directly related to semantics (e.g. predicate-argument relations); (2) many complex constructions can be dealt with adequately in this framework by means of a conglomerate of rules. Most of these rules are also used in the formation of other constructions, and are therefore construction-independent. In this way it is possible to capture several syntactic generalizations. Syntactic transformations played a crucial role in this. As a consequence, the relation between form and meaning becomes rather indirect in many cases, as explained above. (3) The framework makes it possible to incorporate and extend existing syntactically adequate descriptions based on insights from theoretical linguistics into a compositional framework in a fairly direct manner.

In short, the framework makes it possible to integrate some of the best work done in semantics (the work done in the Montague tradition) with some of the best work done in syntax (the generative-grammar tradition).

## 7.2 Topics for Further Research

Though considerable success in expressing syntactic generalizations in a compositional framework was achieved, a number of problems that have not been addressed sufficiently remain. A few of them have been mentioned before. I will repeat these here and add some additional areas where further improvement can be achieved. They can serve as topics for further research.

I mentioned the fact that the rule notation, which was fixed after a certain point, is not fully optimal yet. The limitations of the rule notation make it impossible to adequately express the relevant generalizations in certain cases. I did not deal with this in this book, but it certainly is a real problem. An initial attempt to ameliorate the rule notation has been made already by [Rous and Jansen, prep] and [Jansen, 1992].

I also discussed syncategorematically introduced structure. In the current grammar this is dealt with in a way that is not satisfactory, but in chapter 3 I introduced an alternative method which makes it possible to capture the relevant generalizations.

Sometimes two different rules (e.g. rules from different subgrammars) must perform exactly the same task. In the current situation it is not possible to turn them into one rule, and it is also not possible to let them share their bodies. It would certainly be desirable if one of these restrictions were relieved, so that either the same rule can occur in several subgrammars or different rules can share the same body.

Finally, the *Move  $\alpha$*  program claims that there are no construction-specific and no language-specific rules in grammars, at least for a core set of facts. Though I extensively discussed how we attempted to incorporate construction-independent rules in controlled M-grammar, I did not discuss in any way this second part of the *Move  $\alpha$*  program and no attempts have yet been made to attain it. Nevertheless, it is often the case that rules from different languages share many properties, and it would certainly be appropriate to isolate these and represent them only once. One might imagine that a ‘universal grammar’ is developed, which contains the aspects of the individual grammars which are shared by all grammars. The language-specific grammars should then contain only properties which are specific to the language at hand. One way in which this might be done would be by deriving a full grammar by superimposing a language-specific grammar on the ‘universal grammar’, perhaps by a unification-like operation.

It is clear that such an approach would be a substantial improvement upon the current situation, although it introduces certain complexities as well.

The advantages of such an approach would be that the grammars of individual grammars would be shorter, easier and better to maintain. Furthermore, the grammars of the individual languages will be more uniform, i.e. they will be identical wherever possible, while in the current situation the grammars of the individual languages can differ in uninteresting aspects such as notational conventions, choice of attribute-value names, choice of category names, etc.

The complexities that this approach would introduce basically all relate to making changes to a grammar: for each change one must determine whether it should be a change of ‘universal grammar’ or of the grammar of the language at hand. One should also be careful not to let ‘universal grammar’ be the somewhat accidental ‘intersection’ of the grammars of the individual languages. If that were the case, extending the system with additional languages, or even extending the grammars of the languages dealt with might become more complex. Finally, a change in ‘universal grammar’ may require adaptations in all grammars of the



individual grammars.

Though these additional complexities are introduced by this approach, they cannot count as real arguments against the approach: it will be necessary to consider carefully how changes are made, and this will perhaps be more time-consuming than in the current system, but the eventual result is a much better, more modular grammar which expresses syntactic generalizations more adequately than is possible in the current situation.



# Bibliography

- [Abeillé and Schabes, 1989] Abeillé, A. and Schabes, Y. (1989). Parsing idioms in lexicalized TAGs. In *Proceedings of the European ACL*, pages 1–9, Manchester.
- [Appelo, 1993] Appelo, L. (1993). *Categorial Divergences in a Compositional Translation System*. PhD thesis, Institute for Perception Research (IPO)/ University of Utrecht.
- [Appelo et al., 1987] Appelo, L., Fellingner, C., and Landsbergen, J. (1987). Subgrammars, rule classes and control in the Rosetta translation system. In *Proceedings of the 3rd ACL Conference, European Chapter*, pages 118–133, Copenhagen. Philips Research M.S. 14.131.
- [Arnold et al., 1986] Arnold, D., S.Krauwer, Rosner, M., des Tombe, L., and Varile, G. (1986). The  $\langle C, A \rangle, T$  framework in Eurotra: A theoretically committed notation for machine translation. In *Proceedings of COLING 86, Bonn, August 25th-29th*, pages 297–303.
- [Bech, 1952] Bech, G. (1952). Über das niederländische Adverbialpronomen *er*. In *Travaux du Cercle Linguistique de Copenhague*, 8, pages 5–32. Copenhagen/Amsterdam. also appeared in Hoogteijling ([Hoogteijling, 1969], 147-174).
- [Belletti and Rizzi, 1988] Belletti, A. and Rizzi, L. (1988). Psych-verbs and  $\theta$ -marking. *Natural Language and Linguistic Theory*, 6(3):291–352.
- [Bennis, 1980] Bennis, H. (1980). *Er*-deletion in a modular grammar. In Daalder, S. and M.Gerritsen, editors, *Linguistics in the Netherlands 1980*, pages 58–68. North-Holland, Amsterdam.
- [Bennis, 1986] Bennis, H. (1986). *Gaps and Dummies*. Foris Publications, Dordrecht.
- [Bierwisch, 1963] Bierwisch, M. (1963). *Grammatik des Deutschen Verbs*, volume II of *Studia Grammatica*. Akademie-Verlag, Berlin, 1967 edition.
- [Bresnan, 1982] Bresnan, J. (1982). Control and complementation. In Bresnan, J., editor, *The Mental Representation of Grammatical Relations*, pages 282–390. The MIT Press, Cambridge, Massachusetts/ London, England.

- [Bresnan et al., 1982] Bresnan, J., Kaplan, R., Peters, S., and Zaenen, A. (1982). Cross-serial dependencies in Dutch. *Linguistic Inquiry*, 13(4):613–635.
- [Broekhuis, 1992] Broekhuis, H. (1992). *Chain-Government: Issues in Dutch Syntax*. PhD thesis, University of Amsterdam.
- [Burzio, 1981] Burzio, L. (1981). *Intransitive Verbs and Italian Auxiliaries*. PhD thesis, MIT.
- [Chomsky, 1957] Chomsky, N. (1957). *Syntactic Structures*. Mouton, The Hague.
- [Chomsky, 1964] Chomsky, N. (1964). *Current Issues in Linguistic Theory*. Mouton, The Hague / Paris.
- [Chomsky, 1973] Chomsky, N. (1973). Conditions on transformations. In Anderson, S. R. and Kiparsky, P., editors, *A Festschrift for Morris Halle*. Holt, Rinehart and Winston, New York. Reprinted in [Chomsky, 1977a].
- [Chomsky, 1977a] Chomsky, N. (1977a). *Essays on Form and Interpretation*. North-Holland, New York.
- [Chomsky, 1977b] Chomsky, N. (1977b). On wh-movement. In Culicover, P., T. Wasow, and Akmajian, A., editors, *Formal Syntax*, pages 71–132. Academic Press, New York.
- [Chomsky, 1981] Chomsky, N. (1981). *Lectures on Government and Binding*. Foris, Dordrecht.
- [Chomsky, 1986] Chomsky, N. (1986). *Knowledge of language. Its Nature, Origin and Use*. Convergence. Praeger, New York.
- [Chomsky, 1991] Chomsky, N. (1991). Some notes on economy of derivation and representation. In Freidin, R., editor, *Principles and Parameters in Comparative Grammar*, pages 417–454. MIT Press, Cambridge, Mass.
- [Chomsky, 1992] Chomsky, N. (1992). A minimalist program for linguistic theory. *MIT Occasional papers in Linguistics*, 1.
- [Chomsky and Lasnik, 1991] Chomsky, N. and Lasnik, H. (1991). Principles and parameters theory. MS. MIT. Also appears in J. Jacobs *et al.* (1993), *Syntax: An International Handbook of Contemporary Research*, pp. 506–569, Berlin: Walter de Gruyter.
- [Cinque, 1990] Cinque, G. (1990). Ergative adjectives and the lexicalist hypothesis. *Natural Language and Linguistic Theory*, 8(1):1–39.
- [den Besten, 1981] den Besten, H. (1981). Government, syntaktische Struktur und Kasus. In Kohrt, M. and Lenerz, J., editors, *Sprache: Formen und Strukturen. Akten des 15. Linguistischen Kolloquiums, Münster, 1980, vol. 1*, pages 97–107. Niemeyer, Tübingen. *Linguistische Arbeiten* 98.

- [den Besten, 1982] den Besten, H. (1982). Some remarks on the ergative hypothesis. *Groninger Arbeiten zur Germanistischen Linguistik*, 21:61–82.
- [den Besten, 1983] den Besten, H. (1983). On the interaction of root transformations and lexical deletive rules. In Abraham, W., editor, *On the Formal Syntax of the Westgermania. Papers from the "3rd Groningen Grammar Talks", January 1981*, volume 3 of *Linguistik Aktuell*, pages 97–107. John Benjamins, Amsterdam/Philadelphia.
- [den Besten, 1985] den Besten, H. (1985). The ergative hypothesis and free word order in Dutch and German. In Toman, J., editor, *Studies in German Grammar*, volume 21 of *Studies in Generative Grammar*, pages 23–64. Foris Publications, Dordrecht.
- [den Besten, 1989] den Besten, H. (1989). *Studies in West Germanic Syntax*. Rodopi, Amsterdam / Atlanta, Georgia.
- [den Besten and Rutten, 1989] den Besten, H. and Rutten, J. (1989). On verb raising, extraposition and free word order in Dutch. In *Sentential Complementation and the Lexicon: Studies in Honour of Wim de Geest*, pages 41–56. Foris, Dordrecht.
- [den Besten and Webelhuth, 1987] den Besten, H. and Webelhuth, G. (1987). Remnant topicalization and the constituent structure of VP in the Germanic SOV languages. *GLOW Newsletter*, 18:15–16.
- [Drewes et al., 1984] Drewes, H., Jacobs, E., van den Maagdenberg, F., Veld, J., and de Wolff, M. (1984). Crossing R-graphs. In de Haan, G., Trommelen, M., and Zonneveld, W., editors, *Van Periferie naar Kern*, pages 29–37. Foris Publications, Dordrecht.
- [Emonds, 1976] Emonds, J. (1976). *A Transformational Approach to Syntax*. Academic Press, New York.
- [Engelkamp et al., 1992] Engelkamp, J., Erbach, G., and Uszkoreit, H. (1992). Handling linear precedence constraints by unification. In *Proceedings of the 30th Annual Meeting of the ACL*, pages 201–208.
- [Evers, 1975] Evers, A. (1975). *The Transformational Cycle in Dutch and German*. PhD thesis, University of Utrecht.
- [Flickinger et al., 1987] Flickinger, D., Nerbonne, J., Sag, I., and Wasow, T. (1987). Toward evaluation of NLP systems. Hewlett-Packard Laboratories, Palo Alto, California.
- [Gazdar et al., 1985] Gazdar, G., Klein, E., Pullum, G., and Sag, I. (1985). *Generalized Phrase Structure Grammar*. Blackwell.
- [Grimshaw, 1979] Grimshaw, J. (1979). Complement selection and the lexicon. *Linguistic Inquiry*, 10(2):279–326.

- [Groenendijk and Stokhof, 1982] Groenendijk, J. and Stokhof, M. (1982). Semantic analysis of wh-complements. *Linguistics and Philosophy*, 5:175–233.
- [Haan, 1979] Haan, G. J. d. (1979). *Conditions on Rules*. Foris Publications, Dordrecht.
- [Haegeman and van Riemsdijk, 1986] Haegeman, L. and van Riemsdijk, H. (1986). Verb projection raising, scope, and the typology of rules affecting verbs. *Linguistic Inquiry*, 17(3):417–466.
- [Hoeksema, 1988] Hoeksema, J. (1988). A constraint on governors in the West Germanic verb cluster. In Everaert, M., Evers, A., Huybregts, M., and Trommelen, M., editors, *Morphology and Modularity*. Foris, Dordrecht.
- [Hoekstra, 1984] Hoekstra, T. (1984). *Transitivity. Grammatical Relations in Government-Binding Theory*, volume 6 of *Linguistic Models*. Foris Publications, Dordrecht, Holland.
- [Hoekstra, 1988] Hoekstra, T. (1988). Small clause results. *Lingua*, 74:101–134.
- [Hoogteijling, 1969] Hoogteijling, J., editor (1969). *Taalkunde in Artikelen: een Verzameling Artikelen over het Nederlands*. Wolters-Noordhoff, Groningen.
- [Huybregts, 1976] Huybregts, M. (1976). Overlapping dependencies in Dutch. *University of Utrecht Working Papers in Linguistics*, 1:24–65.
- [Huybregts, 1984] Huybregts, M. (1984). The weak inadequacy of context-free phrase structure grammars. In de Haan, G., Trommelen, M., and Zonneveld, W., editors, *Van Periferie naar Kern*, pages 81–100. Foris, Dordrecht.
- [Isabelle, 1989] Isabelle, P. (1989). Toward reversible MT systems. In *MT Summit II*, Munich.
- [Jansen, 1992] Jansen, P. G. (1992). Reversible programming in  $4_2$ . Master's thesis, University of Amsterdam / Institute for Perception Research (IPO). IPO Report no. 856.
- [Janssen, 1986] Janssen, T. (1986). *Foundations and Applications of Montague Grammar: part 2, Applications to Natural Language*, volume 28 of *CWI tract*. CWI, Amsterdam.
- [Kaplan and Bresnan, 1982] Kaplan, R. M. and Bresnan, J. (1982). Lexical-functional grammar: A formal system for grammatical representation. In Bresnan, J., editor, *The Mental Representation of Grammatical Relations*, pages 173–281. The MIT Press, Cambridge, Massachusetts.
- [King, 1983] King, M. (1983). Transformational parsing. In King, M., editor, *Parsing Natural Language*, pages 19–34. Academic Press, London/New York.

- [Koster, 1975] Koster, J. (1975). Dutch as an SOV language. *Linguistic Analysis*, 1:111–136.
- [Koster, 1978a] Koster, J. (1978a). *Locality Principles in Syntax*. Foris, Dordrecht.
- [Koster, 1978b] Koster, J. (1978b). Why subject sentences don't exist. In Keyser, S. J., editor, *Recent Transformational Studies in European Languages*. MIT Press, Cambridge, Massachusetts.
- [Koster, 1984] Koster, J. (1984). On binding and control. *Linguistic Inquiry*, 15(3):417–459.
- [Kroch and Santorini, 1991] Kroch, A. S. and Santorini, B. (1991). The derived constituent structure of the West Germanic verb-raising construction. In Freidin, R., editor, *Principles and Parameters in Comparative Grammar*, volume 20 of *Current Studies in Linguistics*, pages 269–338. The MIT Press, Cambridge, Massachusetts/ London, England.
- [Landsbergen, 1981] Landsbergen, J. (1981). Adaptation of Montague grammar to the requirements of parsing. In Groenendijk, J., Janssen, T., and Stokhof, M., editors, *Formal methods in the Study of Language Part 2*, number 136 in MC Tract, pages 399–420. Mathematical Centre, Amsterdam. Philips Research Reprint 7573.
- [Landsbergen, 1987] Landsbergen, J. (1987). Isomorphic grammars and their use in the Rosetta translation system. In King, M., editor, *Machine Translation Today: the State of the Art*, pages 351–372. Edinburgh University Press. Philips Research M.S. 12.950.
- [Landsbergen et al., 1989] Landsbergen, J., Odijk, J., and Schenk, A. (1989). The power of compositional translation. *Literary and Linguistic Computing*, 4(3):191–199.
- [Model, 1991a] Model, J. (1991a). *Grammatische Analyse: Syntactische Verschijnselen van het Nederlands en het Engels*. ICG Publications, Dordrecht, Holland.
- [Model, 1991b] Model, J. (1991b). Incorporatie in het Nederlands. *Gramma*, 15(1):57–88.
- [Montague, 1974a] Montague, R. (1974a). English as a formal language. In Thomason, R. H., editor, *Formal Philosophy: Selected Papers of Richard Montague*, pages 188–221. Yale University Press, New Haven and London.
- [Montague, 1974b] Montague, R. (1974b). The proper treatment of quantification in ordinary English. In Thomason, R. H., editor, *Formal Philosophy: Selected Papers of Richard Montague*, pages 247–270. Yale University Press, New Haven and London.

- [Neeleman and Weerman, 1992] Neeleman, A. and Weerman, F. (1992). The balance between syntax and morphology: Dutch particles and resultatives. Ms., University of Utrecht. to appear in NLLT.
- [Netter, 1992] Netter, K. (1992). On non-head non-movement: An HPSG analysis of finite verb position in German. ms. Deutsches Forschungszentrum für Künstliche Intelligenz, GmbH, Saarbrücken.
- [Odiijk, 1989] Odiijk, J. (1989). The organization of the Rosetta grammars. In *Fourth Conference of the European Chapter of the Association for Computational Linguistics. Proceedings of the Conference*, pages 80–86. ACL, European Chapter, ACL. 10-12 April, University of Manchester.
- [Odiijk, 1992] Odiijk, J. (1992). The Rosetta machine translation project. In ter Stal, W., Nijholt, A., and op den Akker, H., editors, *Proceedings of the second Twente Workshop on Language Technology: Linguistic Engineering: Tools and products*, number 92-29 in Memoranda Informatica, pages 87–91, Enschede. Faculteit Informatica, University of Twente.
- [Paardekooper, 1986] Paardekooper, P. (1986). Er was es /ldots. *Onze Taal*, 55:134.
- [Partee, 1973] Partee, B. (1973). Some transformational extensions of Montague grammar. *Journal of Philosophical Logic*, 2:509–534. (Also appeared in Partee (/citeyearPartee76)).
- [Partee, 1977] Partee, B. H. (1977). Constraining transformational Montague grammar: A framework and a fragment. In *Conference on Montague Grammar, Philosophy, and Linguistics*, pages 51–102. University of Texas Press, Austin, Texas.
- [Partee, 1979] Partee, B. H. (1979). Montague grammar and the well-formedness constraint. In Heny, F. and Schnelle, H. S., editors, *Selections from the third Groningen round table*, number 10 in Syntax and Semantics, pages 275–314. Academic Press, New York.
- [Partee et al., 1990] Partee, B. H., ter Meulen, A., and Wall, R. E. (1990). *Mathematical Methods in Linguistics*, volume 30 of *Studies in Linguistics and Philosophy*. Kluwer Academic Publishers, Dordrecht/Boston/London.
- [Pereira and Shieber, 1987] Pereira, F. C. and Shieber, S. M. (1987). *Prolog and Natural-Language Analysis*. Number 10 in CSLI Lecture Notes. Center for the Study of Language and Information.
- [Perlmutter, 1977] Perlmutter, D. (1977). Toward a universal characterization of passive. *Proceedings of the Annual Meeting of the Berkeley Linguistics Society*, 3:394–417. also appeared in /mycitePerlmutter83.



- [Perlmutter, 1978] Perlmutter, D. (1978). Impersonal passives and the unaccusative hypothesis. *Proceedings of the Annual Meeting of the Berkeley Linguistics Society*, 4:157–189.
- [Pollard and Sag, 1987] Pollard, C. and Sag, I. A. (1987). *Information-based Syntax and Semantics*. Number 13 in CSLI Lecture Notes. Center for the Study of Language and Information, Stanford, California.
- [Pollard, 1990] Pollard, C. J. (1990). On head non-movement. Paper presented at the Conference on Discontinuous Constituency, Tilburg.
- [Pullum, 1991] Pullum, G. K. (1991). *The Great Eskimo Vocabulary Hoax and Other Irreverent Essays on the Study of Language*. The University of Chicago Press, Chicago/London.
- [Quirk et al., 1972] Quirk, R., Greenbaum, S., Leech, G., and Svartvik, J. (1972). *A Grammar of Contemporary English*. Longman, Essex. (Eleventh Impression, 1985).
- [Rosetta, 1993] Rosetta, M. T. (1993). Compositional translation. M.S. 924, Institute for Perception Research (IPO).
- [Ross, 1967] Ross, J. (1967). *Constraints on Variables in Syntax*. PhD thesis, MIT, Cambridge, MA.
- [Rous and Jansen, prep] Rous, J. and Jansen, P. (in prep.). Reversible programming in  $\lambda_2$ . Philips Research Laboratories.
- [Růžička, 1983] Růžička, R. (1983). Remarks on control. *Linguistic Inquiry*, 14(2):309–324.
- [Sag, 1991] Sag, I. (1991). Linguistic theory and natural language processing. In Klein, E. and Veltman, F., editors, *Natural Language and Speech*. Springer-Verlag, Berlin.
- [Schenk, 1986] Schenk, A. (1986). Idioms in the Rosetta machine translation system. In *Proceedings of the 11th Conference on Computational Linguistics*, Bonn.
- [Schenk, 1989] Schenk, A. (1989). The formation of idiomatic structures. In Everaert, M. and van der Linden, E.-J., editors, *Proceedings of the First Tilburg Workshop on Idioms*, pages 145–158. ITK proceedings, Tilburg University.
- [Schenk, 1992] Schenk, A. (1992). The syntactic behaviour of idioms. In Everaert, M., van der Linden, E.-J., Schenk, A., and Schreuder, R., editors, *Proceedings of Idioms, International Conference on Idioms, Tilburg, The Netherlands, 2-4 september 1992*, pages 97–110. ITK, Tilburg University.
- [Schermer-Vermeer, 1985] Schermer-Vermeer, E. (1985). De onthullende status van *er*. *Spektator*, 15(2):65–84.

- [Schermer-Vermeer, 1986a] Schermer-Vermeer, E. (1986a). Er was eens een misverstand. *Onze Taal*, 55:134–135.
- [Schermer-Vermeer, 1986b] Schermer-Vermeer, E. (1986b). Er was eens . . . . *Onze Taal*, 55:48–49.
- [Shieber, 1985] Shieber, S. (1985). Evidence against the context-freeness of natural language. *Linguistics and Philosophy*, 8:333–343.
- [Shieber, 1987] Shieber, S. M. (1987). Separating linguistic analyses from linguistic theories. In Whitelock, P., Wood, M., Somers, H., Johnson, R., and Bennet, P., editors, *Linguistic Theory and Computer Applications*, pages 1–36. Academic Press, London.
- [Stowell, 1981] Stowell, T. (1981). *Origins of Phrase Structure*. PhD thesis, MIT.
- [Stowell, 1991] Stowell, T. (1991). Small clause restructuring. In Freidin, R., editor, *Principles and Parameters in Comparative Grammar*, volume 20 of *Current Studies in Linguistics*, pages 182–218. The MIT Press, Cambridge, Massachusetts/ London, England.
- [Thomason, 1974] Thomason, R. H., editor (1974). *Formal Philosophy: Selected Papers by Richard Montague*. Yale University Press, New Haven.
- [Torrego, 1984] Torrego, E. (1984). On inversion in Spanish and some of its effects. *Linguistic Inquiry*, 15(1):103–129.
- [Travis, 1984] Travis, L. (1984). *Parameters and Effects of Word Order Variation*. PhD thesis, MIT, Cambridge, Mass.
- [Uszkoreit, 1982] Uszkoreit, H. (1982). German word order in GPSG. In D. Flickinger, M. and N. Wiegand, editors, *Proceedings of the First West Coast Conference on Formal Linguistics*, Stanford, California. Stanford University.
- [Uszkoreit, 1987] Uszkoreit, H. (1987). *Word Order and Constituent Structure in German*. Number 8 in CSLI Lecture Notes. Center for the Study of Language and Information, Stanford, California.
- [van Bart and Kager, 1984] van Bart, P. and Kager, R. (1984). Er is hier - een alternatief voor een diagnose. In de Haan, G., Trommelen, M., and Zonneveld, W., editors, *Van Periferie naar Kern*, pages 1–10. Foris Publications, Dordrecht.
- [van Gestel, 1989] van Gestel, F. (1989). Idioms and X-bar theory. In Everaert, M. and van der Linden, E.-J., editors, *Proceedings of the First Tilburg Workshop on Idioms*, pages 103–126. ITK proceedings, Tilburg University.
- [van Hout, 1986] van Hout, A. (1986). *ER-peculiarities in Rosetta: an analysis of er and its translations in English and Spanish*. Master's thesis, KUB Tilburg.

- [van Noord, 1991] van Noord, G. (1991). Head corner parsing for discontinuous constituency. In *Proceedings of the 29th Annual Meeting of the Association for Computational Linguistics, Berkeley*, University of Texas at Austin.
- [van Noord, 1993] van Noord, G. (1993). *Reversibility in Natural Language Processing*. PhD thesis, University of Utrecht.
- [van Noord et al., 1989] van Noord, G., Dorrepaal, J., Arnold, D., Krauwer, S., Sadler, L., and des Tombe, L. (1989). An approach to sentence-level anaphora in machine translation. In *Proceedings of the 4th ACL Conference, European Chapter, Manchester, 10-12 april 1989*, pages 299–307, University of Manchester. University of Manchester, Institute of Science and Technology.
- [van Noord et al., 1990] van Noord, G., Dorrepaal, J., van der Eijk, P., Florenza, M., and des Tombe, L. (1990). The MiMo2 research system. In *The Third International Conference on Theoretical and Methodological Issues in Machine Translation of Natural Language*, pages 213–233, University of Texas at Austin. Linguistic Research Center.
- [van Riemsdijk, 1974] van Riemsdijk, H. (1974). À propos de l'extension du principe A-sur-A aux syntagmes prépositionnels. In Rohrer, C. and Ruwet, N., editors, *Actes du Colloque Franco-Allemand de Grammaire Transformationnelle. Vol I, Études de Syntaxe*, pages 206–215. Niemeyer, Tübingen.
- [van Riemsdijk, 1978] van Riemsdijk, H. (1978). *A Case Study in Syntactic Markedness*. The Peter de Ridder Press, Lisse.
- [van Riemsdijk, 1992] van Riemsdijk, H. (1992). Terug bij ERVANDAAN. In Bennis, H. and de Vries, J. W., editors, *De Binnenbouw van het Nederlands. Een bundel artikelen voor Piet Paardekooper*, pages 311–317. ICG Publications, Dordrecht.
- [Williams, 1980] Williams, E. (1980). Predication. *Linguistic Inquiry*, 11:203–237.
- [Williams, 1981] Williams, E. (1981). Argument structure and morphology. *The Linguistic Review*, 1:81–114.
- [Zubizarreta, 1985] Zubizarreta, M. L. (1985). The relation between morphophonology and morphosyntax: The case of Romance causatives. *Linguistic Inquiry*, 16(2):247–289.
- [Zwarts, 1988] Zwarts, J. (1988). An analysis of genericity and its translation in Rosetta. Master's thesis, Instituut A.W. de Groot voor Algemene Taalwetenschap, University of Utrecht.



# Samenvatting

## Compositionality en Syntactische Generalisaties

Het doel van dit proefschrift is aan te geven hoe syntactische generalisaties uitgedrukt kunnen worden in een compositioneel grammaticaal formalisme voor natuurlijke taal dat ontworpen is om een automatisch vertaalsysteem te ontwikkelen.

Het grammaticaal formalisme dat als uitgangspunt is genomen, ‘gecontroleerde M-grammatica’, is compositioneel van aard, d.w.z. het is een kader waarin vorm en betekenis op een zeer directe manier aan elkaar gerelateerd worden door de zogenaamde ‘regel-op-regel’ benadering. Een compositionele grammatica bestaat uit basisexpressies, d.w.z. vormeenheden die een basisbetekenis hebben, en grammaticaregels die geassocieerd zijn met een betekenisoperatie. Het compositionele karakter van de grammatica leidt tot een beschrijving waarin de vormopbouw in hoge mate bepaald wordt door de betekenisopbouw. Hoewel dit gewenst is vanuit het oogpunt van het leggen van de relatie tussen vorm en betekenis, kan het ook leiden tot beschrijvingen die vanuit zuiver vormoogpunt minder adequaat zijn.

Dit probleem wordt in gecontroleerde M-grammatica opgevangen door het aannemen van een speciaal type regels, *syntactische transformaties*. Syntactische transformaties zijn regels die geassocieerd zijn met de identiteitsoperatie als betekenis. In dit proefschrift wordt beargumenteerd dat transformaties noodzakelijk zijn om syntactische generalisaties tot uitdrukking te brengen.

Het vangen van syntactische generalisaties met behulp van transformaties leidt tot beschrijvingen waarin een conglomeraat van betekenisvolle regels en syntactische transformaties constructies afleidt. Er worden analyses voorgesteld waarin zorgvuldig is afgewogen welke aspecten specifiek zijn voor een bepaalde constructie, en welke aspecten algemenere toepassing hebben. De niet aan een specifieke constructie gebonden aspecten worden uitgefactoriseerd in transformaties en op zo’n manier in de gehele grammatica opgenomen dat zij ook een rol spelen bij het afleiden van andere constructies. Deze methode wordt geïllustreerd aan de hand van een groot aantal constructies uit de talen die het vertaalsysteem behandelt (Nederlands, Engels, Spaans), met de nadruk op het Nederlands. De voorgestelde analyses worden ook vergeleken met analyses uit andere computationeel georiënteerde formalismes. De nadruk ligt op monolinguale aspecten.

Aan de constructie-onafhankelijke regels kan gewoonlijk geen betekenis toegekend worden. Hier brengt het bestaan van syntactische transformaties uitkomst. Een consequentie van het type analyse dat voorgesteld wordt, is dat de relatie tussen vorm en betekenis minder direct is: veel zinnen vertonen onderling zowel vormverschillen als betekenisverschillen, maar in veel gevallen worden de vormverschillen door de ene verzameling regels verantwoord, en de betekenisverschillen door de andere.

De geïllustreerde constructies zijn, samen met een groot aantal andere verschijnselen, daadwerkelijk geïmplementeerd in een groot, consistent, experimenteel automatisch vertaalsysteem, Rosetta3, dat ontwikkeld is op het Philips Natuurkundig Laboratorium en dat bedoeld was als een prototype op basis waarvan echte applicaties ontwikkeld zouden kunnen worden. Het ontwikkelen van dit systeem vereiste dat bepaalde assumpties vanaf een zeker punt gefixeerd werden. Dit had zowel voordelen als nadelen. Een nadeel is dat bepaalde aannames niet meer gewijzigd konden worden zelfs als duidelijk was dat ze complicaties van bepaalde beschrijvingen opleverden. Een voordeel was, dat bepaalde problemen juist aan het licht konden komen doordat een groot aantal assumpties gefixeerd waren en de beschrijvingen binnen dit kader gegeven moesten worden. Dit proefschrift identificeert een aantal van dit soort problemen of tekortkomingen van het formalisme en er worden verscheidene suggesties en voorstellen gedaan om deze problemen in de toekomst te vermijden.

Hoofdstuk 1 beschrijft het doel van het proefschrift en beschrijft het karakter van de ontwikkelde grammatica's: het zijn artefacten, waarin de nadruk ligt op beschrijvende adequaatheid en grote overdekking van de feiten. Ze zijn effectief omkeerbaar, zodat ze zowel voor generatie als voor analyse gebruikt kunnen worden. Het belang van preciese grammatica's, zowel vanuit theoretisch als vanuit praktisch oogpunt, wordt benadrukt.

Hoofdstuk 2 beschrijft het gebruikte formalisme, gecontroleerde M-grammatica, in detail en illustreert hoe het als basis kan dienen voor een daadwerkelijk automatisch vertaalsysteem waarin de methode van compositionele vertaling toegepast wordt. Bovendien worden enkele substantiële kenmerken van de ontwikkelde grammatica's behandeld.

Hoofdstuk 3 geeft aan dat de analyse van bepaalde constructies in een compositioneel formalisme zonder syntactische transformaties leidt tot beschrijvingen die vanuit syntactisch oogpunt minder adequaat zijn. Geïllustreerd wordt hoe adequatere beschrijvingen van deze constructies verkregen kunnen worden en hoe constructie-specifieke en constructie-onafhankelijke aspecten afgescheiden kunnen worden in aparte regels. Het gedrag van hulpwerkwoorden in het Engels, verschillende zinstypes in het Nederlands, volgordevarianten, ongebonden afhankelijkheden en generische zinnen worden vanuit dit perspectief behandeld. Bovendien wordt een specifiek voorstel gedaan om de regelmatigheid van de syntactische structuur van syncategorematisch geïntroduceerd materiaal op principiële wijze te verantwoorden.

Hoofdstuk 4 geeft een beschrijving van de behandeling van een verschijnsel waar syntaxis en semantiek nauw aan elkaar raken, namelijk predicaat-argument

relaties. Geconcludeerd wordt dat gecontroleerde M-grammatica een stevig fundament verschaft voor een principiële behandeling van deze relaties. Er wordt een onderscheid gemaakt tussen argumenten en bepalingen. Drie kenmerken van de behandeling van argumenten, die onmiddellijk voortvloeien uit het compositionele karakter van de grammatica's, worden beschreven en toegelicht. Er wordt beargumenteerd dat een ordeningsconventie voor argumenten en het onderscheid tussen externe en interne argumenten noodzakelijk zijn. Het hoofdstuk bevat een uitgebreide beschrijving van de representatie van deze informatie bij basisexpressies

In detail wordt de analyse van verschillende soorten 'verborgen argumenten' (d.w.z. argumenten die niet aan de oppervlakte direct zichtbaar zijn) geschetst. Voor een bepaalde verzameling constructies wordt beargumenteerd dat noch een beschrijving als argument noch een beschrijving als bepaling adequaat is, en een nieuwe categorie, *gebonden bepalingen*, wordt geïntroduceerd om de problemen die verbonden zijn aan de beschrijving van deze constructies te vermijden. De status van zogenaamde 'small clauses' in de grammatica's wordt uitgebreid toegelicht. Tot slot wordt voorgesteld dat allerlei regelmatige relaties tussen predicaten m.b.t. hun argumenten door regels, bijvoorbeeld door lexicale regels, uitgedrukt dienen te worden.

Hoofdstuk 5 behandelt de beschrijving van een aantal syntactische constructies, met name lijdende-vormconstructies, de distributie van de persoonsvorm in het Nederlands, kruisende afhankelijkheden tussen werkwoorden en hun argumenten in het Nederlands, en onbegrensde afhankelijkheden. In detail wordt aangegeven hoe in de beschrijvingen van deze constructies constructie-specifieke en constructie-onafhankelijke aspecten enerzijds, en betekenisvolle en betekenisloze aspecten anderzijds, geïdentificeerd en geïsoleerd zijn, zodat syntactische generalisaties tot hun recht kunnen komen. Verder wordt aangegeven wat de relatie is tussen de geïmplementeerde analyses en de analyses uit de meer theoretisch georiënteerde grammaticale formalismes waarop zij gebaseerd zijn.

Hoofdstuk 6 bevat een gedetailleerde beschrijving van de analyse binnen het gecontroleerde M-grammatica formalisme van een van de meest complexe verschijnselen van het Nederlands, namelijk de syntaxis van R-pronomina (d.w.z. de pronomina *er*, *hier*, *daar*, *ergens*, *nergens*, *overal*, *waar*). Er wordt beargumenteerd dat de twee centrale problemen m.b.t. de syntaxis van R-pronomina, nl. (1) de distributie van R-pronomina, en (2) de schijnbare 'samenvall' of versmelting van meerdere R-pronomina, beregeld kunnen worden door een beperkt aantal aannames. Deze aannames implementeren twee centrale ideeën, nl. (1) de distributie van R-pronomina wordt verantwoord door een beperkt aantal posities te postuleren waar R-pronomina op kunnen treden en (2) de schijnbare 'samenvall' van R-pronomina wordt verantwoord doordat het niet optreden van R-pronomina welgevormde resultaten oplevert indien het optreden noch op syntactische gronden noch op semantische gronden noodzakelijk is.

Hoofdstuk 7 vat de belangrijkste conclusies samen, en schetst onderwerpen voor verder onderzoek. Deze onderwerpen zijn voor het grootste deel rechtstreeks af te leiden uit de overblijvende problemen van deze dissertatie en betreffen o.a.

de beschrijving van syncategorematisch geïntroduceerde elementen en de taalafhankelijke formulering van bepaalde regels.



# Curriculum Vitae

Jan Odijk was born in Schiedam on 28 February 1956. In 1974 he obtained the Gymnasium  $\alpha$  certificate from Scholengemeenschap Spieringshoek in Schiedam and started to study Slavic languages and literature and general linguistics at the University of Utrecht. He graduated in 1981 cum laude in syntactic theory with subsidiaries in computational linguistics and Slavic linguistics. He was employed from 1982-1985 as a research assistant at the institute A.W. de Groot for general linguistics in Utrecht with a grant from the Dutch organization for scientific research (ZWO) and worked on Slavic syntax. Since 1985 he participated in the Rosetta machine translation project which was carried out at Philips Research Laboratories, first as an employee of the university of Utrecht (1985-1988), later as an employee of Philips (since 1988). He is currently working at Philips Research in the group Language of the Institute for Perception Research (IPO).