



**HAL**  
open science

## Sentence Generation by Pattern Matching.

Michael Zock

► **To cite this version:**

Michael Zock. Sentence Generation by Pattern Matching.. R. Mitkov & N. Nicolov. John Benjamins, pp.35, 1997, Recent Advances in Natural Language Processing. Series: Current Issues in Linguistic Theory. hal-00661773

**HAL Id: hal-00661773**

**<https://hal.science/hal-00661773>**

Submitted on 20 Jan 2012

**HAL** is a multi-disciplinary open access archive for the deposit and dissemination of scientific research documents, whether they are published or not. The documents may come from teaching and research institutions in France or abroad, or from public or private research centers.

L'archive ouverte pluridisciplinaire **HAL**, est destinée au dépôt et à la diffusion de documents scientifiques de niveau recherche, publiés ou non, émanant des établissements d'enseignement et de recherche français ou étrangers, des laboratoires publics ou privés.

## Contents

*Michael Zock*

Sentence generation by pattern matching: the problem of syntactic choice	1
<b>Index</b>	<b>36</b>

# Sentence Generation by Pattern Matching: The Problem of Syntactic Choice

MICHAEL ZOCK  
*LIMSI, CNRS*

## Abstract

This paper tries to account for verbal fluency, that is, the speed with which people compute syntactic structures. As we all know, people produce speech fluently without making too many mistakes. Given the known time constraints this is a remarkable performance. How is this possible? Verbal fluency, we believe, can be accounted for by the following two facts: people essentially use *pattern matching* and *mapping rules* as strategy and knowledge source. Rather than being confined to local strategies (strict incremental processing on a concept-to-concept basis) and formal grammars, they operate on larger chunks (global strategy) by using mapping rules. This is more economical, without being necessarily more error prone. Second, proficient speakers have learnt to recognise potential linguistic structures on the basis of the formal characteristics of the conceptual structures, that is, proficient speakers are able to make good guesses concerning the syntactic structures that best express the conceptual input.<sup>1</sup>

## 1 Introduction: The speaker's problem

Text or discourse production basically consists in determining, organising and translating content in order to achieve specific communicative goals. We shall be concerned here only with the last component, the translation of a conceptual structure (message) into its corresponding linguistic form. Looking at this problem from a psycholinguistic point of view we will try to

---

<sup>1</sup> This paper is a slightly revised version of a paper presented in 1988 at the 1st International Workshop on "Cognitive Linguistics", held in Bulgaria. It was meant to appear three years later in a book entitled "Explorations in Cognitive Linguistics". Unfortunately, though announced, this book never saw the market. While our views have evolved in the meantime, — this was then only preliminary work, — we believe that our basic premises concerning the process and the speaker's knowledge still hold: (i) natural language generation is basically pattern-matching; (ii) the speaker's expertise resides in knowing a set of *patterns* and a set of *mapping rules* for converting an input (message) into an output (linguistic form).

provide evidence for two claims, one dealing with *knowledge*, the other one dealing with the *process*. *Claim 1*: the speaker's expertise resides in knowing a set of patterns (conceptual, syntactic) and a set of rules (mapping rules) for converting a given structure (deep structure) into its corresponding form (surface structure).<sup>2</sup> *Claim 2*: the process of conversion is basically pattern-driven (pattern matching). People typically work on chunks (e.g., noun groups, propositions) rather than on atomic units (single words, concepts).

If one accepts the idea that messages are coded in terms of semantic networks or conceptual graphs (Sowa 1984), one will understand the text producer's problem. Having generated a message (what to say), s/he is left with a set of nodes and arcs (concepts and relationships) for which s/he must find words and adequate sentence patterns (how to say it).<sup>3</sup>

Surprisingly enough, this seems to be no real problem for most speakers, even in completely new settings (spontaneous discourse). Skilled speakers seem to have plenty of time, and they do not make many mistakes. One may well ask how they succeed in performing such a complex task given the known space and time constraints: human short-term memory is limited (Miller 1956), and speech is very fast (3–5 words per second).

The secret of skilled speakers, we believe, is that, rather than operating on small isolated units such as words or concepts (local strategy), they operate on larger chunks, that is, conceptual configurations. Put another way, skilled speakers do not proceed strictly word-by-word (concept-by-concept), rather they operate on larger conceptual patterns.<sup>4</sup> Having gained a large

---

<sup>2</sup> The fact that people use patterns is not in contradiction with the notion of a formal grammar. The latter is actually a device to generate them.

<sup>3</sup> While the arcs are conceptual or syntactic relations (agent, subject, etc.), the nodes of the graph can be words or concepts of various levels of abstraction (**animal** vs. **dog** vs. **four-legged carnivorous wild or domesticated animal**). The representation and function of abstract concepts and words is thus very much alike: both are very economical means, — a kind of short hand notation, — for larger conceptual chunks.

The fact that graphs allow for hybrid knowledge representation, and the fact that they can be manipulated easily (contraction/expansion), makes them excellent tools at the interface level (i.e., for a potential user in the case of applications) and for modelling the cognitive process: expansions of an abstract, underspecified message graph (conceptual level); contraction of this conceptual graph to a lexically specified, but syntactically unspecified graph (lexical level); visualisation of syntactic reflexes resulting from choices made at a higher level (pragmatic, conceptual, linguistic). The syntactic consequences may show up in changes of the names of the links (*agent* becoming *subject*, a *beneficiary* becoming an *indirect object*, etc.) and in the addition of morphosyntactically relevant information (type of auxiliary, type of preposition, etc.). For an example, see Zock (1994).

<sup>4</sup> The fact that word-by-word processing may lead into dead ends has been shown by

amount of experience in a given language, they recognise typical structures (patterns), i.e., they know straight away what conceptual structures match what linguistic forms.<sup>5</sup>

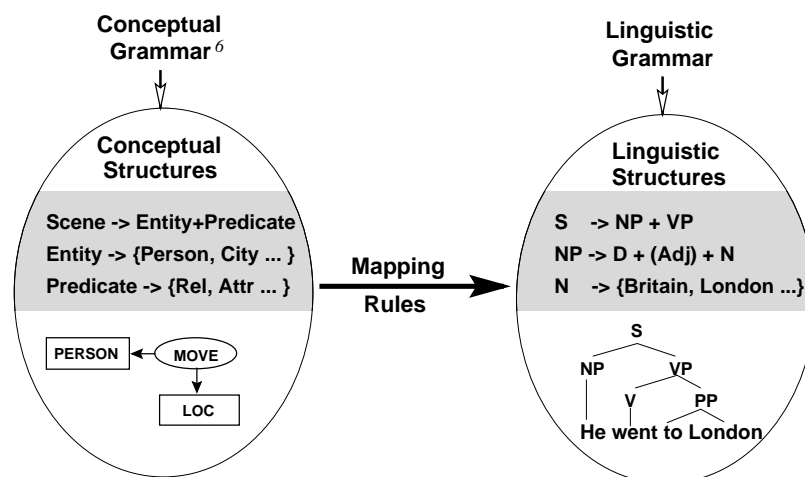


Fig. 1: *Mapping rules*,  
the missing link between conceptual structure and linguistic structure

Actually, the idea of patterns or schemata is not new. It has a long tradition in *philosophy* (Kant 1781), in *sociology* (Goffmann 1974), in *structural-*

Zock et al. (1986). Actually, clitics in French nicely illustrate the need for lookahead or preplanning. Suppose you were to pronominalize  $y$  and  $z$  of the following proposition give  $(x, y, z)$ . In this case it is not possible to determine their relative position, unless one knows the roles (person) of both objects.

If you compare (a-c) you'll notice that the position of the direct object ("it", that is "le" in French) depends on the person of the indirect object (3d person or not). Put differently, the positions of the two objects are interdependant, that is, their respective position cannot be determined unless the value of the attribute PERSON of both objects is known.

- (a) il me LE donne (he gives it to ME)
- (b) il LE lui donne (he gives it to HER)
- (c) il te LE donne (he gives it to YOU)

<sup>5</sup> Learning a language is thus learning a set of variably abstract *patterns*, a set of *mapping rules* and their respective *conditions of use*. Sentence generation can go either way, from abstract to specific patterns (refinement), or from specific to more general patterns (generalisation), lower level patterns becoming integrated into higher level patterns: (det + adj + noun) => NP; (verb + noun) => VP; (adv + verb) => AdvP; (NP + VP) => Sentence.

<sup>6</sup> The *conceptual grammar* controls the assembly, i.e., legal combinations of possible contents, that is, it specifies what is meaningful in a given culture, whereas the *linguistic grammar* specifies the possible forms.

and in *text linguistics* (Harris 1951; Fries 1952; Roberts 1962; van Dijk 1977), in *psychology* (Bartlett 1932; Koffka 1935; Piaget 1970; Bruner 1973; Rumelhart 1975; Mandler 1979; Ausubel 1980) and in *artificial intelligence* (Minsky 1975; Wilks 1975, Schank & Abelson 1977).<sup>7</sup> Nevertheless, despite its long standing tradition, schema approaches and formal grammars have two major shortcomings: with the exception of Mel'cuk's work, neither nor account for the *correspondences* (mappings) between the different structures or levels. Hence, they do not make explicit on what structural cues, i.e., formal characteristics of the message (conceptual structure), the speaker's decisions are based when s/he chooses a specific linguistic form (syntactic structure). Yet, structures are of little help if one does not know what to do with them, that is, what they stand for.<sup>8</sup>

Another problem with the schema approach lies in the fact that schemata are hard to constrain. Hence, lack of refinement, or lack of proper constraints may result in ambiguity (analysis — parsing) or overgeneration (production). If the user is not told where the limits of the schemata lie (explicitation of the schema constraints), s/he will use these patterns even in cases where they do not apply.

If one agrees with our point of view that natural language processing is basically schema-driven,<sup>9</sup> then the question arises of how people manage to recognise linguistic structures on the basis of conceptual structures. This is what this paper is about.

---

<sup>7</sup> While Bartlett, Schank/Abelson, van Dijk and Rumelhart identified patterns on the text or discourse level (schemata, scripts, macro-structures, story grammars), Harris and Fries dealt with sentence patterns. The idea of linguistic patterns has also been extensively used in the classroom, where pattern-drills have been a major teaching strategy especially during the sixties, when behaviorism was at its peak (Lado 1964; Rivers 1972). Things have changed radically after Chomsky's devastating critique of Skinner's book "Verbal Behavior" (Skinner 1957; Chomsky 1959).

<sup>8</sup> One can't but agree with Bock et al. when they write "In existing models of language production the first mapping from messages to linguistic relations involves linking non-linguistic cognitive categories to linguistic categories. However the categories themselves are variably specified, because there is little consensus of what the appropriate ones might be." (Bock et al. 1992:151)

<sup>9</sup> While early natural language systems like SIR (Raphael 1968), STUDENT (Bobrow 1968), ELIZA (Weizenbaum 1966) and SHRDLU (Winograd 1972) relied heavily on low-level schemata (syntactic patterns), more recent systems use high level schemata, i.e., text patterns (Mc Keown 1985; Rösner, 1987). For a criticism of the latter see Hovy (1990). See also Patten et al.'s use of the notion of knowledge compilation which is somehow akin to our notion of pattern matching (Patten et al. 1992).

## 2 What kind of evidence can we provide in favour of pattern-matching?

There are several good reasons for accepting such an approach, both structural and procedural:

**Structural evidence:** Human experience and social interactions are structured, regular, hence predictable to some extent. This regularity, of course, reflects in language. There are definite limits with regard to linguistic creativity: new, original thoughts have still to be cast in old patterns. Languages are schematic to a great extent, that is, every language has a fairly large set of patterns in order to express concepts, relations and events, etc. For example,

Table 1: *Schemata for a definition*

<X> is <Y> that <VP>	A computer is a machine that processes information.
<X> is a sort of <Y>	A bicycle is a sort of vehicle.

Table 2: *Schemata for comparison*

<X> is somehow like <Y>.	A cat is somehow like a tiger.
<A> is to <B> as <C> is to <D>.	Good is to light as evil is to darkness.

This is true not only on the higher levels (paragraph, text level), — stories, news, weather forecast, sport reports, etc. are clearly schematic — but also on the lower levels (phoneme, word, sentence level).

Table 3: *Mapping of ontological categories on syntactic categories*

Actions, events, states, processes	verbs	<i>build, happen, be, sleep</i>
Entities, names, places	nouns	<i>car, Paul, Tokyo</i>
Properties, attributes of entities	adjectives	<i>young, bright</i>
Manner, attributes of actions	adverbs	<i>slowly</i>
Intensifier, location, time	adverbs	<i>very, here, tomorrow</i>
Means	prepositions	<i>by, with</i>
Spatial relations:		
path, position, direction	prepositions	<i>from, in, on, towards</i>

When translating a message into discourse, the speaker maps a conceptual structure (deep structure) onto a linguistic form (surface structure). Thus, concepts are mapped on words, each of which have a specific categorial

potential, i.e., part of speech (Table 3),<sup>10</sup> *deep-case relations* are mapped on grammatical functions (Table 4), and *conceptual configurations*, i.e., larger conceptual structures, are mapped on syntactic structures (Table 5), etc.

Table 4: *Mapping of case relations on grammatical functions*

agent, cause	subject
object, patient	direct object
beneficiary, recipient	indirect object

Table 5: *Mapping of conceptual structures onto syntactic structures*

1. [PERSON:#]<-agnt-[PERFORM]-obj->[MUSIC:*] <sup>11</sup> D + N + V + D + N	The girl plays a song.
2. [MUSIC:*]<-attr-[QLTY: GOOD] D + Adj + N D + N + RelPr + Copula + Adj	a nice song a song that is nice
3. [PERSON:#]<-agnt-[PERFORM]-obj->[MUSIC:*]<-attr-[QLTY: GOOD] D + N + V + D + Adj + N	The girl plays a nice song.

This mapping can be done in various ways, directly or indirectly, that is, via various intermediate structures.<sup>12</sup> The units on which these processes operate may be single concepts (Tables 3), relations (Tables 3 and 4), or larger chunks (Table 5). Subordinate patterns may be integrated into superordinate patterns, etc.<sup>13</sup>

The following example reveals several interesting formal characteristics concerning patterns:

<sup>10</sup> Hence, we assume that there is no *one-to-one correspondence* between concepts and words, or between words and parts of speech. For example, a given concept (LOVE) may well map on several words (love, like, be fond of, . . .), each of each may be realized by several syntactic categories (noun, verb, adjective).

<sup>11</sup> We use the “#” and “\*” signs to signal the communicative status of the referent (determinate vs. indeterminate). The “#” sign signals the fact that the entity referred to is known, hence requires a definite article. In a similar vein “\*” signals indeterminacy.

<sup>12</sup> Swartout’s (1983) and Mel’cuk’s work (Mel’cuk & Zholkovskij 1970) lie at the extremes. The former used no intermediate structure at all, whereas the latter used no less than seven levels to get from a meaning representation to its surface form. ATN-approaches (Simmons & Slocum 1972) and semantic grammars (Burton 1976; Hendrix 1977) lie somewhere in between.

<sup>13</sup> If you look at Table 5, you will notice that the conceptualisation 2 [MUSIC:\*]<-attr-[QLTY: GOOD] is integrated into message 1, yielding message 3.



(TITLE) <PERSON> **has** (kindly) **brought to**  
 <DET> **attention the** (following) <FACT/PROBLEM>.

*Prof Joshi has kindly brought to my attention the following fact.*

- Patterns have a fixed and a variable part. The fixed parts are highlighted and written in small letters, while variable parts are written in capital letters and in angle brackets or normal parenthesis.
- The variables may be optional or obligatory parts of the pattern. <PERSON>, <FACT>, <DET> are all obligatory elements, while (TITLE) is, from a linguistic point of view, optional. If we generalised “kindly” into some variable (MANNER), — we would do so if we found out that other attributes, or synonyms can occur in this position, — then we would have an optional variable.
- The variables may be of different kind: semantic or syntactic (<PERSON> vs. <DET>).
- Patterns have optional parts. This is illustrated by the adverb “kindly” here above which appears in parenthesis.
- Certain elements of the patterns may need morphological adjustment: subject-verb agreement, determiner (*my* vs. *her*), etc.
- Patterns may contain information of different kind (conceptual, lexical, syntactic). Put differently, patterns can be hybrid.
- Patterns can be semantically equivalent, that is, they can be synonyms (“he brought to my attention” vs. “he drew to my attention”).
- Patterns can be embedded into each other (hierarchy of patterns).

A grammar could be seen as a list of patterns composed only of variables, and a hybrid approach like ours is only a cut through of this grammar at different levels of abstraction.

All the above mentioned patterns are *productive* in the sense that they allow for the creation of a wide range of linguistic forms. There may be even direct connections between situations and conceptual structures on one hand, and between these conceptual structures and their linguistic counterparts on the other. It should be noted, however, that the relationship between conceptual structures and linguistic structures is not one-to-one (see below). Hence, neither syntactic categories (e.g., part of speech), nor syntactic structures are *determined* by conceptual structures, that is, the former cannot be predicted solely on the basis of the latter. Nevertheless, there is a strong tendency to translate a given conceptual element or structure by a specific syntactic form, that is, there are default mappings (see Tables 3 and 5). For empirical evidence of how conceptual structures induce syntactic structures see Brown (1958).

As we have said, there are other reasons that plead in favour of schema-driven processing, namely, the speed and economy of learning and processing.

**Procedural arguments:** Speech is fast, yet it is slow compared to thought. As we all know, thoughts tend to get lost if not expressed in time. Word-by-word processing is thus not a good candidate for processing the data: it is not only slow, but also error prone. If words and syntactic structures were computed strictly incrementally, that is on a word-to-word, or concept-to-concept basis (local planning), we would never get the job done in time, and we would easily get stuck in dead ends, that is, talk ourselves into a corner. The order in which the concepts come into our minds, and the order in which they have to be expressed in the surface string are not necessarily the same. Global planning (*look-ahead* or global view) is thus necessary in order to speed up the process and to cut down on backtracking. It, should be noted though that, while pattern matching seems to be a good strategy at the initial stages of generation, it is not obvious at all that it can be used safely throughout the process. It provides generally only a first sketch (outline) which needs to be checked against the syntactic requirements, i.e., constraints of the lexical material actually used (subcategorisation features).

It seems that an approach whereby pattern matching is only seen as a first step allowing for refinements or changes, contains all the basic ingredients to go from canned text to full-blown sentence generation.

Another argument for pattern matching stems from the observation that people manage to communicate even very complex thoughts by using very few patterns. This is even more so if they speak in a foreign language. Many students use pattern matching as a basic strategy when producing language: their messages are structured the same way, that is, they are cast in the same kind of sentence- or discourse pattern. Actually, verbal skill can be measured in terms of the number of patterns a speaker is able to use adequately in a given communicative situation, and by his aptitude to vary and to make local adjustments to them.

From the above it should be clear that spontaneous discourse is only possible if the speaker is operating on larger units (chunks) than words or single concepts.<sup>14</sup> If this point of view is accepted, together with the point

---

<sup>14</sup> For an enlightening discussion concerning the use of larger units than words (most notably, idiomatic and fixed expressions) see Becker (1975). This view, just as ours, contrasts with the notion of strict incremental processing (see, for example, Kempen & Hoenkamp 1987; de Smedt 1990).

that proficient speakers are, above all, good pattern matchers, then the question arises how people manage to recognise these structures. In other words, is there a way to identify a good candidate for a *relative clause*, a *that-clause*, an *infinitive*, etc. on purely conceptual grounds, that is, on the basis of the formal characteristics of the underlying conceptual input? Before answering this question, we would like to show why this problem cannot be adequately addressed within the framework of structure-oriented linguistics.

### 3 Why cognitive linguistics, or, why study natural language in the realm of cognitive science?

Natural languages are both products and processes. Understanding the way they function thus requires the study of both. In other words, language, or language use cannot be adequately accounted for solely by studying the outputs (sentence-structure).

In contrast to *structure-oriented linguistics* which describes only the physical products (sentences), *cognitive linguistics* tries to account for the processes, i.e., the operations necessary to transform an input (for example, a visual scene) into an output (text: description of the scene). While the former are concerned with products, the latter are interested in the processes operating on data. Hence the following questions are relevant for cognitive linguistics:

- What are the different knowledge sources (pragmatic, conceptual, linguistic)?
- What are the input-output data?
- What kind of operations are performed on these data (transformations, mapping rules)?
- How do biological and cultural factors constrain the representations and processes?
- What are the functional relations between the components (hierarchical vs. heterarchical architecture)?
- How is the process decomposed (control of information flow)?
- How is the relevant information coded, stored, retrieved and processed?

The goal of cognitive linguistics is to describe and to explain linguistic competency and performance for natural systems (human beings). Obviously, structure and process vary with the restrictions of the information processor (man vs. machine).

Languages are systems for the coding, manipulation and communication of information. They are symbolic means for storing, processing (reason) and transmitting information. As with any tool, they are designed with respect to a goal (function) and with respect to the user-constraints. As these constraints are different for human beings and for machines (memory, attention span), we would expect natural languages to be different from artificial languages (algebra, logic, etc.).<sup>15</sup> Natural languages, as opposed to artificial languages, are very flexible. The different components (conceptual, lexical and syntactic) are highly interdependent, each component possibly influencing the others. The advantage of such a heterarchical architecture is that it allows for various orders of data-processing. For example, lexical choice may precede the choice of syntactic structure and vice versa. For more details see Zock (1990).

One could view the functioning of the mind, hence, the functioning of natural language somehow like the functioning of a complex society (oligarchy). The two systems are organised in a similar way: (i) problem solving is decomposed: the result is produced not by a superexpert, but by a team of specialists; (ii) the different agents (components) contributing to the solution have a certain amount of autonomy; (iii) the agents negotiate, that is, they do not only communicate their results and draw on the results produced by their colleagues, but they can also adapt their behavior to allow for accomodation of the results produced by the other components.

The advantages of such a heterarchical kind of organisation are multiple: (i) *freedom of processing*: various orders are possible to reach the solution; (ii) *time-sharing*: each agent can work on its own without having to wait for an order coming from a higher component; (iii) *flexibility*: information flow is bidirectional; (iv) *opportunistic planning*: as information becomes available at different moments and in unpredictable ways, and since the different components can accomodate, it is possible to have the different agents compete and to use the first result produced by any of them.

The major drawback with this kind of system, where everything is more or less interdependant, is that it becomes extremely difficult to see the dependency relationships, that is, it is hard to see what causes what, or what action has what outcome. This is quite obvious for covert activities like

---

<sup>15</sup> "Natural languages are ambiguous, imprecise and sometimes awkwardly verbose. These are all virtues for general communication, but something of a drawback for communicating concisely as precise a concept as the power of recursion. The language of mathematics is the opposite of natural language: it can express powerful formal ideas with only a few symbols." (Friedman & Feldeisen 1987: xi).

language, but it is also true for such complex activities as political decisions. Whichever the case, consequences of a choice may be far reaching and hard to predict. In this sense, there are many points in common between speaking well a language, hence communicate efficiently, and being a good politician. In both cases one has to make the right choice at the right moment.

Man and machine are not subject to the same constraints. Humans have a very limited working memory,<sup>16</sup> they are poor serial processors, and they are not very logical. But they are intuitive, creative and above all good pattern matchers. That is, they can spontaneously discover the right solution on intuitive grounds, they can conceive of efficient strategies to solve a problem (for example, how to access information stored in the long-term memory), and they can recognise complex patterns (configurations, *Gestalt*). Machines, on the other hand are logical, they are good serial processors, they have a perfect memory, but they have little capacity for creativity, for intuition, or for perceiving global structures.

Linguists who want to work out an *ecologically* valid theory — that is, provide a description of the data which is not only formally correct, but which is also computationally sound (processable) — need to take these factors into account. Otherwise their theory will remain mere description with limited explanatory power or relevance for practical purposes.

Now, if our point concerning pattern matching — that is the use of mapping rules applied to larger chunks (i.e., conceptual configurations) — is sound, one may ask several questions: Where do linguistic structures come from? (Section 4); What is the difference between conceptual and linguistic structures? (Section 5); What do syntactic structures depend upon? (Section 6); What do typical conceptual structures, patterns or configurations look like? (Section 7); How can the speaker recognise a specific syntactic structure? (Section 8).

#### 4 Where do linguistic structures come from?

One can try to answer this question from a phylogenetic point of view, or from the point of view of the process that takes place when translating thoughts (messages) into language (text). From a phylogenetic point of view it seems that linguistic structures reflect perceptual structures. This is the

---

<sup>16</sup> Due to memory constraints, sentences are built incrementally: planning and execution partially overlap. While uttering a partially planned conceptual structure, the next part of the message is planned: we think while we speak, and, while speaking, we think (Kempen & Hoenkamp 1986).

position held by many *psychologists* (Paivio 1971; Kempen 1977; Osgood, 1980; Anderson 1983; Miller & Johnson-Laird 1985), *linguists* (Fillmore 1977; Langacker 1983; Habel 1988) and *computer scientists* (Hill 1984; Sowa 1984; Arbib 1986). As a matter of fact, there are many similarities between language and perception, both from a structural and a procedural point of view.

Both are compositional, and both have well-formedness and completeness conditions (though in the sense of Gestalt psychology rather than in the mathematical sense). Furthermore, natural language and images are produced and perceived in a similar way, that is, globally. Both of them are to some extent holistic entities. We start to produce or to recognise a global structure (pattern) which we then fill in with details, that is, we tend to go from the general to the specific. The way how this might be done at the conceptual level is discussed in Zock (1996).

What does this mean for sentence generation? It simply means, that rather than processing word-by-word, or concept-by-concept (local strategy), humans process larger chunks, trying to match entire conceptual structures on linguistic structures (global strategy). In other words, processing is done via pattern matching.<sup>17</sup> This is probably true on all levels: conceptual, syntactic, lexical, and even phonological.

Before turning to the problem of how people recognise linguistic structures on the basis of conceptual structures, we would like to comment on the relationship between conceptual structures and their linguistic counterparts, that is, words and syntactic structures.

## 5 Conceptual structures and syntactic structures are to a great extent parallel

Linguistic structures (order of words) and conceptual structures (order in which thoughts become available, i.e., spring into our mind), while not entirely parallel, correlate to a large extent,<sup>18</sup> that is, items belonging together conceptually tend to appear side by side in the surface structure.<sup>19</sup> This is

<sup>17</sup> It is probably for this very same reason that people are able to recognise misspelled words, despite the speed of reading. We certainly don't look for every character, yet we are able to see the mistakes, especially if they occur at specific points. Strangely enough, it is for similar reasons that people overlook mistakes. We don't perceive what is, but we perceive what ought to be (Bruner 1973).

<sup>18</sup> The regularities concerning their mappings are discussed in Section 7.

<sup>19</sup> See Behagel's first law (cited in Vennemann 1975), Anderson's graph deformation principle (Anderson 1983), or Levelt's pioneering work on linearisation (Levelt 1981, 1982).

reasonable with regard to economy (memory). If conceptual and linguistic structures were not parallel to some extent, we would constantly be faced with a storage problem.<sup>20</sup> For, whenever the word expressing a given conceptual fragment cannot be inserted into the surface string, it needs to be held in working memory until it can be attached to the string. If translating a conceptual structure into linguistic form consists above all in finding words for parts of the graph (semantic network)<sup>21</sup> and in ordering them,<sup>22</sup> it seems reasonable to try to maintain as much as possible the conceptual connectivity in the surface form (syntax). That is, the words should be uttered as much as possible in the same order in which the concepts for which they stand, have been generated.

The question that arises now is, how do we put it all together, or, how do we compute syntactic structures? In order to answer this question let's take an example. Suppose you wanted to express the following message, message which could have been planned and expressed incrementally:

[MAN] <-(AGT) - [CATCH] - (OBJ) -> [FISH]

[MAN] <-(AGT) - [MOVE] - (LOC) -> [GROUND]

[MOVE] <-(MAN) - [FAST]

One way to start the process is by lexicalising, that is, by trying to find words that cover (express) the message planned.<sup>23</sup> This kind of mapping is prob-

<sup>20</sup> Of course, there are limits and exceptions (long distance phenomena, split infinitives, etc.) which make strict parallelism impossible. Language is linear, thought is relational. The representation of the latter being a graph, it is in principle possible to add information at any moment. For exceptions to this principle of adjacency, see Stockwell's discussion on the 'heavier element principle' and the 'topicalisation principle' (Stockwell 1977: 68-69 and 75-76).

<sup>21</sup> Words hardly ever stand for single concepts, generally they stand for definitions. Yet definitions require a graph. This being so, it makes sense to express (or code) these definitions in terms of conceptual structures (graphs). In that respect there is no fundamental difference between the underlying meaning of words and sentences. According to our view, the process of lexicalisation consists in matching word definitions (conceptual graphs defining the meaning of the word) on an utterance graph, i.e., a structure containing the message to be conveyed. How such a message might be built has been discussed in (Zock 1996).

<sup>22</sup> Obviously, there is more to determining surface form: computation of part of speech, insertion of function words, morphological operations (inflections, agreement), etc.

<sup>23</sup> According to our view, lexicalisation is performed in two steps. During the first step only those words are selected that pertain to a given semantic field (for example, movement verbs). At the next step the lexical expert selects from this pool the term that best expresses the intended meaning, i.e., the most specific term (maximal coverage). For more details see (Nogier & Zock 1992; Zock 1996).

ably done stepwise, as it is hard to imagine that a speaker is able to find simultaneously all the words of a very big conceptual chunk. This means that, having found a word for the chunk [MAN] <- (AGT) - [CATCH] - (OBJ) -> [FISH],<sup>24</sup> the speaker tries to find a candidate for the remaining part of the message.<sup>25</sup>

The result of this stepwise consumption of the message graph is a lexicalised conceptual graph (LCG) whose links will then be replaced by functional information (subject, direct object, etc.) and lexical categories (part of speech, step 3). On the basis of this Preliminary Syntactic Structure (PSS) a tree is built. In order to perform the final operations (inflection, etc.) morphological information is added. For more details see (Nogier 1991; Nogier & Zock 1992; Zock 1996). This entire mapping process is depicted in Figure 2.

### 5.1 Discussion

The careful reader has certainly noticed the following facts: (i) the same mechanism was used for lexicalisation and precomputation of syntactic structure: pattern matching;<sup>26</sup> (ii) words were chosen prior to syntactic-structure;<sup>27</sup> (iii) the initial message graph may be considerably simplified

<sup>24</sup> For the sake of simplicity we will ignore here the fact that a “man catching fish”, and a “fisherman” are not necessarily the same. While the former may be an *amateur*, the latter is a *professional*.

<sup>25</sup> Of course, other strategies are possible. Rather than trying to find the next word (breadth first), the speaker could work in depth, trying to finalize the surface form of the first lexical element, that is, determine its syntactic category, i.e., part of speech. In that case he would pursue lexicalisation only once the final form of the preceding element has been computed. Another alternative would be to start the process by determining syntactic structure, inserting then lexical items into the computed syntactic slots. Which strategy is chosen under what condition remains an empirical question.

<sup>26</sup> The assignment of *part of speech* to words (or, more precisely, to word stems) is performed at the PSS level. This could be done via the mapping rules described in Figure 4. However, these rules might be insufficient, in particular if there are several candidates. Ultimately we do need a grammar in order to check at the different points of the chain the possibility of a given category. While the mapping rules specify how a given concept or conceptual chunk may be mapped onto a syntactic category or syntactic structure, formal grammars specify what categories can, or should occur at a specific point in time, that is, at a specific point in the chain.

<sup>27</sup> Please note, that at this stage words are not inflected. What we get are the base forms of words. Note also that, unlike Kempen & Hoenkamp (1987), or Nogier (1991), who conflate the last two steps into one, we do not assume that syntactic information (syntactic functions, syntactic categories) is computed simultaneously with the rootforms. *Part of speech* and *syntactic functions* are determined later. There is one empirical finding though which is troublesome for our approach: when people fail to find the right word, they tend to come up with an alternative (synonym) which belongs to the same syntactic category as the one they were looking for.



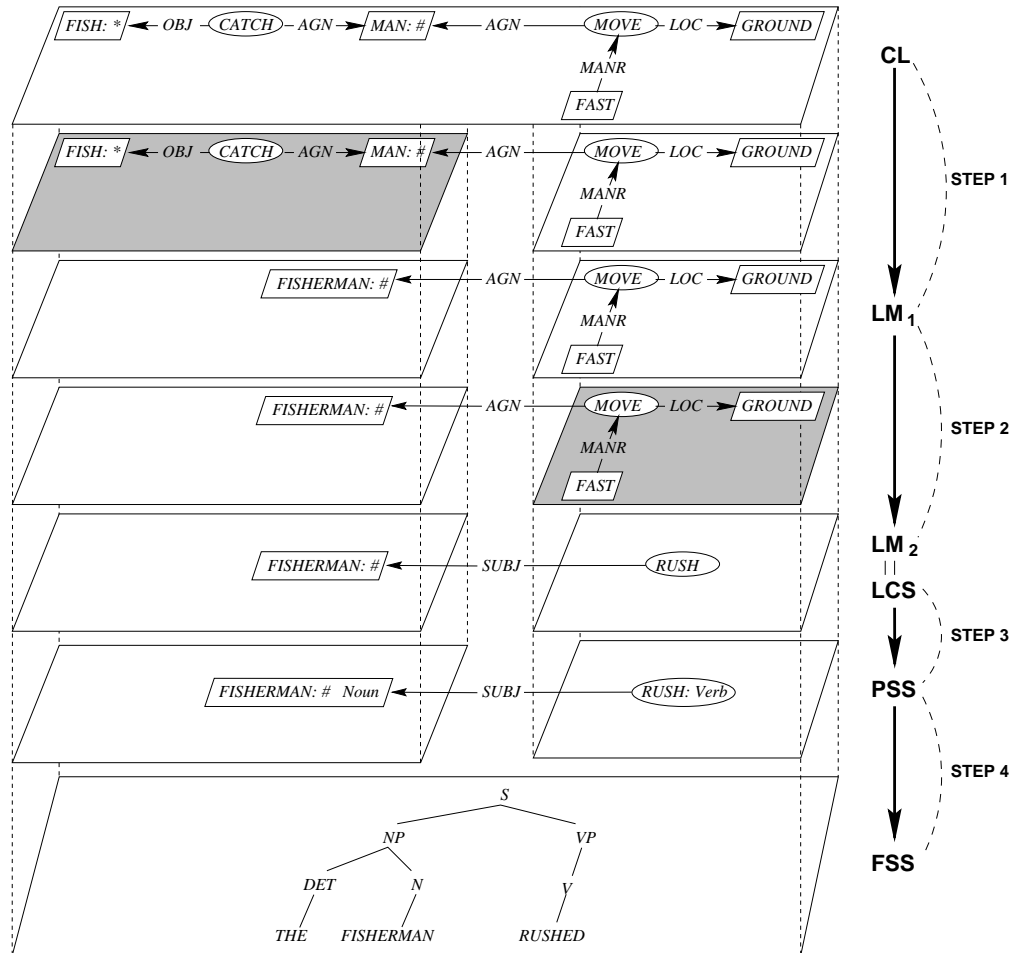


Fig. 2: *The lexicon as a mediator between the conceptual and linguistic levels*<sup>28</sup>

CL:	conceptual level	LCS:	lexicalised conceptual structure
LM <sub>1</sub> :	lexical mapping for word <sub>1</sub>	PSS:	preliminary syntactic structure
LM <sub>2</sub> :	lexical mapping for word <sub>2</sub>	FSS:	final syntactic structure

(contraction) if the syntactic structure is computed via the lexicon (contraction of the message graph to a lexicalized conceptual graph). Actually, this is one of the principal reasons for computing words before syntactic structure: large parts of the message graph become reduced to a relatively small lexical graph (compare the graphs at the conceptual level and the LCS

<sup>28</sup> For illustrative purposes we assume here serial processing. Yet there seems to be evidence that people compute words (and perhaps even syntax) in parallel. For pointers to the relevant literature, see (Levelt 1989).

level in Figure 2).<sup>29</sup> If the syntactic structure is computed before words are chosen, syntactic constraints may be used to choose among lexical alternatives (synonyms). This may happen in the case of parallel structures, where one part constrains the other. Compare:

- (i) *We were expecting the worst, but hoping for the best.*
- (ii) \**We were expecting the worst, but hoped for the best.*

In such a case one uses the same tense in both clauses.

Obviously, one has to justify the fact that lexicalisation precedes the determination of syntactic structure. Basically there are three possible strategies:

1. Syntactic structure is determined prior to word choice. This strategy is implied in traditional generative grammars, where the syntactic tree is built top down. The words are inserted fairly late during the derivational process into syntactically specified slots.
2. Word choice precedes syntax (syntactic trees are built bottom up).
3. Words and their syntactic structure are computed in parallel.<sup>30</sup>

Of course, either of these strategies could be used, though, we don't believe very much in the first option for the following reasons.

First and above all, the speaker wants to convey meanings. Yet, syntax conveys little meaning compared to words. Second, if syntactic structures are to be computed on the basis of conceptual graphs whose nodes are more elementary than words, it is hard, if not impossible, to compute the syntactic structures first: the message graphs are simply too big to make such a strategy very feasible. In addition, it doesn't really make sense to compute the syntactic structure at this point of the process, since large parts of the graph will be reduced at the next step (lexicalisation) anyhow. Last, but not least, syntactic structure depends to a large extent on the subcategorisation features of the words used (see Table 6). Hence, it cannot be computed entirely without knowing the words that are being used. Obviously, it would be nice to decide on this issue on the basis of empirical work. Unfortunately, for the time being we lack conclusive psycholinguistic evidence. For a good discussion though see Aitchison (1983, Chapter 11).

---

<sup>29</sup> One may wonder though whether this is not an artificially induced problem, resulting from our way of modelling the process.

<sup>30</sup> One could also think of a hybrid solution: depending on the situation, priority is given to syntactic or to lexical choice.

## 6 What do syntactic structures depend upon?

Syntactic structures depend basically on three sets of variables, or choices: conceptual, pragmatic and linguistic.

**Conceptual choices.** Different conceptual structures map generally onto different syntactic structures:

CONCEPTUAL STRUCTURE	LINGUISTIC STRUCTURE	STRING
[PERSON:#]<-attr-[SPEED:HIGH]	Pron+Copula+Adj	<i>She is fast.</i>
[PERSON:#]<-agnt-[MOVEMENT]	Pron+V+Adv	<i>She runs fast.</i>
	↑ attr	
	[speed:high]	

As we have shown already (Tables 3 and 5), there is no one-to-one mapping between conceptual structures and their linguistic correlates (parts of speech, syntactic structures). Different conceptual structures or relationships may be expressed by the same kind of linguistic structure.<sup>31</sup> This is particularly obvious for genitives:

Peter's car ...	[possession]
Peter's brother ...	[family relationship]
Peter's leg ...	[inalienable possession, part of]

Conversely, the same conceptual structure or relationship, for example the notion of possession, may map onto different linguistic structures or forms (paraphrase):

This car <b>belongs</b> to the president.	[verb]
This is the car <b>of</b> the president.	[preposition]
This is the president's car.	[case: genitive]
This is <b>his</b> car.	[possessive pronoun]

The conceptual structure is by no means sufficient in order to decide on the linguistic structure. Pragmatic and linguistic information is also necessary. What can be assumed to be known? Is there a word for a given concept? Can this word be used in this particular way?

Suppose we were to express the following idea: ⟨PERSON⟩ be ⟨PROFESSION⟩ ⟨PLACE⟩. In that case one can use in English a **verb**, if the person referred to is a *teacher*, but not if s/he is a *professor* or *doctor*:

<sup>31</sup> This is probably for reasons of economy. The number of possible conceptual structures is enormous. If there were only *one-to-one correspondences*, we would have to learn a great many different syntactic structures. Furthermore, we would have to create a new syntactic structure every time we invented a new conceptual structure.



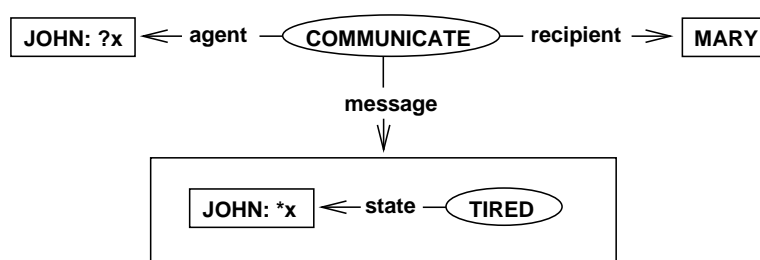


Fig. 3: *John communicated to Mary the fact that he was tired*

The above example illustrates the interaction between lexical choices and syntactic structure. In other words, a syntactic structure cannot be chosen in isolation, or solely on conceptual grounds.<sup>34</sup>

Another example that nicely shows the far reaching consequences of lexical choice is the following. Suppose we were to express in French the following idea:  $\text{HELP}_{\text{past perfect}}(\text{JOHN}, \text{MARY})$ . Suppose furthermore that *John* and *Mary* were known. Depending on the chosen *verb* (*aider*, *venir en aide*, *rendre service*) various aspects of the surface form would change (see Table 6): the *pronoun* (*la* vs. *lui*), the *auxiliary* (*être* vs. *avoir*), the object's *agreement* marking on the verb (*aidée* vs. *aidé*). For more details see Zock (1994).

Table 6: “He has helped her”  
(consequences of the verb choice on *clitics* and *auxiliaries*)

SENTENCE	PRONOUN	AUXILIARY	AGREEMENT
Il l'a aidée.	<i>la</i>	<i>avoir</i>	verb agrees
Il lui est venu en aide.	<i>lui</i>	<i>être</i>	no agreement
Il lui a rendu service.	<i>lui</i>	<i>avoir</i>	no agreement

In the next section we will show how one can recognise certain fundamental syntactic structures on the basis of the formal characteristics of the conceptual structure.

<sup>33</sup> By convention we shall use a *question mark* for odd sounding forms, and an *asterisk* for ungrammatical sentences. The non native speaker of English may be puzzled by this subtlety of English, yet it exists and is explainable. The reason why, *I am tired*, following, *John told Mary*, sounds odd, is due to the fact that the verb *tell* means something like *to report*. A report being something that follows the event it describes or reports, it looks strange if the speaker switches from a reported event to a present event, i.e., direct speech.

<sup>34</sup> For more details on how linguistic structures may vary as a function of pragmatic, conceptual and linguistic choices, see Zock (1988; 1994).

## 7 Prototypical patterns

If we look at syntactic structures, we discover after a while strong correlations with their conceptual counterparts: conceptual structures (see Table 3 and Figure 4). For example, *nouns* usually represent **entities**, *verbs* express actions, *adverbs* stand for manner, etc.

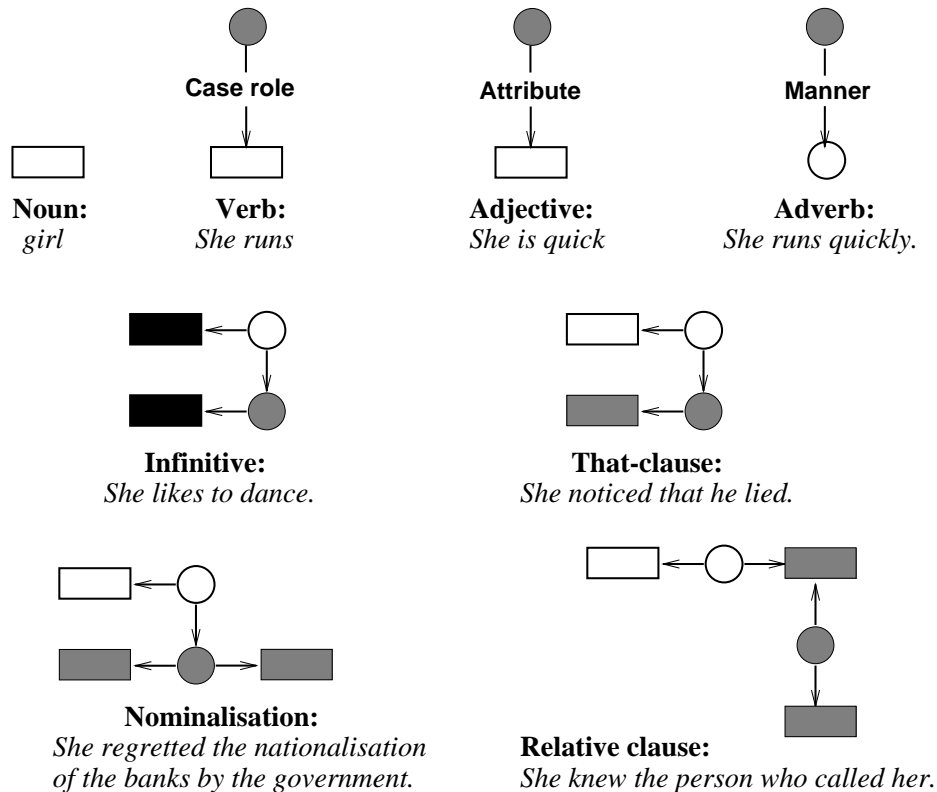


Fig. 4: *Some prototypical structures*<sup>35</sup>

Once we have discovered that, we can use this knowledge the other way round, that is, for generation. Hence, **entities** may map onto *nouns*, **actions/event/processes** may be expressed by *verbs*, etc.<sup>36</sup> Obviously,

<sup>35</sup> The shaded nodes signal particularly relevant information for a given syntactic structure.

<sup>36</sup> As we have pointed out already, this need not always be the case, **actions** may well map onto *nouns* (nominalisation), **directions** may be expressed by *prepositions* or a *verb*, etc. Also, our examples hold for English. While the ontological categories

what holds for concepts holds also for larger conceptual chunks. Hence, looking at the lexicalised conceptual structure one can predict to some extent not only the part of speech, but also the potential syntactic structure (see Table 5).

Figure 4 encodes some typical patterns. By convention, we will use circles for the predicates (verbs, adjectives, adverbs)<sup>37</sup> and boxes for the arguments (nouns, propositions). It should be noted however, that the mapping approach, in order to become feasible, imposes special constraints on the knowledge representation, consistency: elements being morphologically different, yet playing semantically a similar role, should be coded the same way. This is the case with verbs, adjectives and adverbs (see Figure 4). Though syntactically different, they are conceptually alike, hence coded the same way (by a circle). All of them are predicates, the difference hinging in the nature and function of their arguments: **entities** for verbs or adjectives, **actions**, **events**, or **processes** for adverbs. In addition, information of the links are relevant. In order to decide whether predicates pointing to an entity should map onto a verb or an adjective, one must check whether the link is a *case role* (in the case of nouns), or an *attribute* (adjective). One might also appreciate the similarity between *infinitives* and *that-clauses*. The difference hinging on coreference and on the kind of verb: *that-clauses* require specific kinds of verbs.

Figure 5 shows a simple conceptual structure (5A) and the way how one may get progressively (5B, 5C) to its linguistic counterpart (5D). Obviously, there is more than one way to build the corresponding tree (top down, bottom up, bidirectionally).<sup>38</sup> Anyhow, the language user having performed this kind of process a number of times, will gradually become able to go directly from the conceptual structure to the (prefinal) linguistic form.

In the next section we will show how one can recognise one particular syntactic structure, relative clauses.

---

(**entities**, **actions**, **attributes**, etc.) are to a large extent universal, the mappings are language dependent.

<sup>37</sup> That's why we will not use them for indicating the links. Another deviation with regard to John Sowa's notation is the direction of the links for adverbs.

<sup>38</sup> Please note that we have here a mixed strategy of pattern matching and incremental processing. It might also be worthwhile mentioning that 3B is a hybrid form: it contains conceptual and syntactic information.

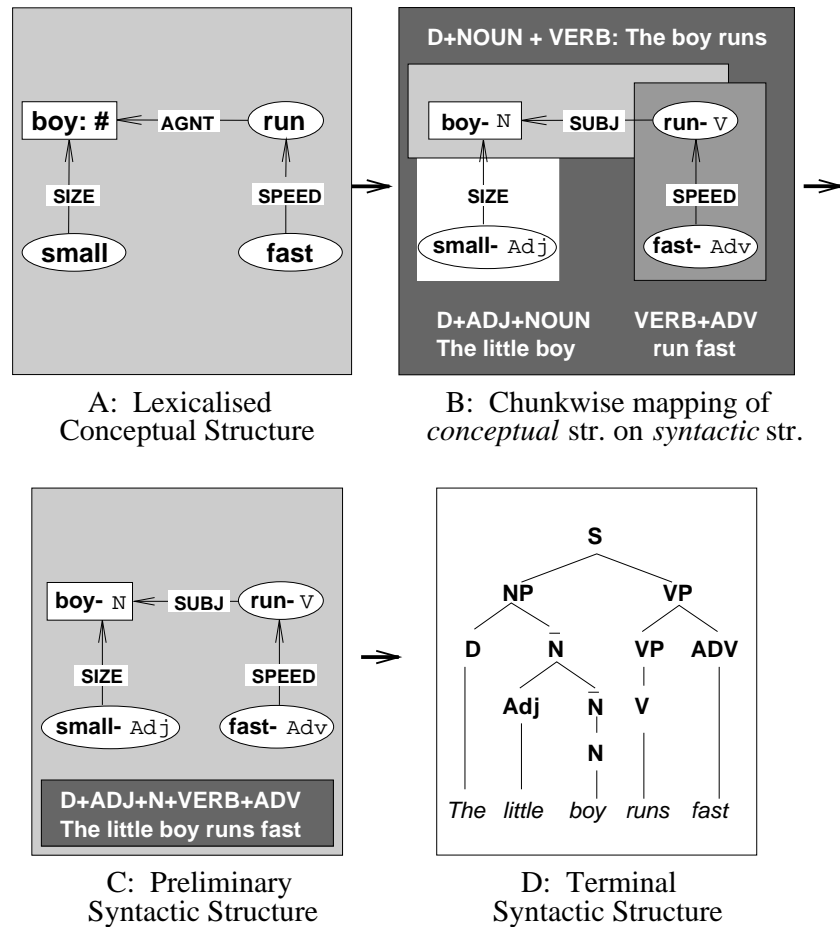


Fig. 5: *From conceptual to linguistic structures*

## 8 Where do relative clauses come from, how can they be recognised, and what do they depend upon?

These are actually three questions, for which, due to space constraints, we will provide only sketchy answers.

Relative clauses add information. In that respect they are similar to adjectives. The information given may be crucial for the identification of the referent (restrictive relative clause) or not (nonrestrictive relative clause). This latter case is generally marked by a comma or a pause.

A typical situation for a relative clause arises when some entity participates in more than one event: `little(boy:#5) & fast(run(boy:#5))`.

This is a conceptual condition which can be captured by a mapping rule: an entity being pointed at by two opposing arcs (see Figure 5). This formal



characteristic might be used by a speaker, recognising this structure as a potential *relative-clause candidate*.<sup>39</sup> We use the word “potential”, because the condition mentioned here above, though necessary, is by no means sufficient. Actually, the conceptual structures could be expressed in either of three ways: (1) by a simple sentence, (2) two independent clauses, or (3) a relative clause.<sup>40</sup>

- (1) *The little boy runs fast.*
- (2) *The boy is little. He runs fast.*
- (3) *The boy who is little runs fast.*

Besides signalling the fact that an entity participates in more than one event, relative clauses signal relative prominence. Put differently, relative clauses factorise and highlight information. In addition, they are devices for increasing processing time and for allowing for spontaneity (expression of afterthoughts, i.e., thoughts that were not planned at the onset of articulation). The extra time they allow for may be needed for encoding the rest of the main clause.

Yet, as we shall see, there is more to the generation of relative clauses than just coreference. Take for example the following propositions (see Figure 6):

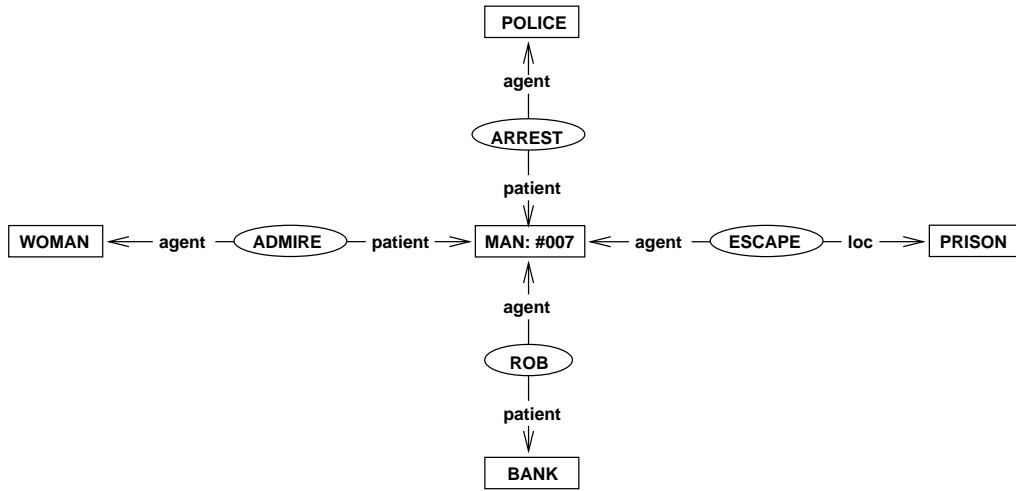
- |   |   |
|---|---|
| (1) <code>rob(man007, bank)</code>      | (3) <code>escape(man007, prison)</code> |
| (2) <code>arrest(police, man007)</code> | (4) <code>admire(woman, man007)</code>  |

When communicating these events one has to consider several factors. *Chunking*: shall all these events be expressed in a single sentence, that is, a series of independent clauses, a coordinated sentence, or a relative clause? In the latter case one has to pay attention to the role played by the coreferential element (`man`) in each clause. What *role* does the to-be-embedded element play (`agent`, `patient`)? Does it play the same role in both events (that is, in the future matrix clause and subordinate clause)? The clauses (1,3) and (2,4) are symmetrical in the sense that in both cases the *man* plays the same role. In the first case (1,3) he is the **agent**, whereas in the second case (2,4) he is the **patient**. This symmetry can, of course, have an effect on realisation.

---

<sup>39</sup> As one can see, mapping rules, or the cues a speaker may be sensitive to, are not only language specific, but also dependent on the knowledge representation formalism. If messages are coded in terms of semantic networks or conceptual graphs, then nodes and links become of focal interest, if one uses first order logic (propositions), then identity of reference might be a crucial element to be looked for.

<sup>40</sup> Please note, again there are subtle, though important differences between these forms at the pragmatic level, differences which we will not deal with here.

Fig. 6: *Conceptual input*


---

Yet, there are still quite a few other factors playing a role: *order of mention* (linear order of clauses, topicalisation), *communicative status* (definite vs. indefinite), *presupposition* (known, unknown), *tense*. Some of these decisions are prior to syntactic processing, and their consequences (communicative status) should be part of the input. Suppose we were to express the following events (see Figure 7):<sup>41</sup>

Event-1 (E1)    `rob(man, bank)`  
 Event-2 (E2)    `escape(man, prison)`

By varying systematically certain parameters such as *chunk size* (independent vs. complex clause, coordination vs. subordination), *order of events*, *topicalisation*, *communicative status* of the participants (definite vs. indefinite),<sup>42</sup> *tense*, etc. we will notice, that certain structures are not possible, while others, though grammatically correct, sound simply odd. It is by analyzing this data that we could get an answer to our question “what factors codetermine relative clauses”.

---

<sup>41</sup> For the sake of simplicity, these representations will not include information concerning time, tense, aspect, mood, etc. It should be noted though that this kind of information may play an important role, constraining the use of a particular structure.

<sup>42</sup> For example, the communicative status of the common object has to be identical in the main clause and subordinate clause, as it would be incorrect to produce: *A man who had robbed the bank escaped from prison*. Put differently, one can’t relativise input structures like `rob(man:*, bank) & escape(man:#, prison)` or `rob(man:#, bank) & escape(man:*, prison)` because the man, though being coreferential, does not have the same communicative status.

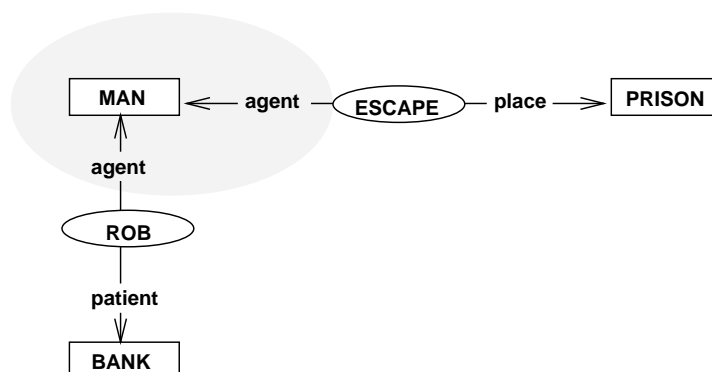


Fig. 7: *A man robbing a bank and escaping from prison*

Unfortunately, for reasons of space, we cannot perform such a systematic analysis here, though it would be worth the effort.

### 8.1 Discussion

Table 7 highlights several interesting problems. The conceptual input encoded in Figure 7 is clearly underspecified. Hence, the precise syntactic form is undecidable. This shows up in Table 7 by the large variety of possible forms, forms which express subtle differences though. It also shows up in the changes of tense which are not trivial at all. Let's have a closer look at some of the sentences.

Sentence 1a, though grammatically correct, sounds odd as it seems to lack information, namely *having robbed a bank, the man was put into jail*. This effect would be less striking in a coordinated structure, where the listener would get the feeling that the speaker expressed a list of actions performed by some man: *A man has robbed a bank and escaped from prison*.

By comparing the sentences (1b, 1d) and (2b, 2c), one can see to what extent 'communicative status' (definitive vs. indefinite) and 'tense' (see 1a and 1b – 1e) may affect the interpretation, hence, acceptability of the sentence.

Table 7: *Stylistic effects due to the variation of parameters like chunk size, order of mention, topicalisation, etc.*

INDEPENDENT CLAUSES	SUBORDINATION
<hr/>	
(1) ORDER: E1 > E2;      TOPIC: man	
(1a) <i>A man has robbed a bank. He escaped from prison.</i>	(1b) <i>A man who had robbed a bank escaped from prison.</i>
	(1c) ? <i>A man who had robbed a bank had escaped from prison.</i>
	(1d) <i>The man who had robbed a bank escaped from prison.</i>
	(1e) ? <i>The man who had robbed a bank had escaped from prison.</i>
<hr/>	
(2) ORDER: E1 < E2;      TOPIC: man	
(2a) <i>A man (had) escaped from prison. He had robbed the bank.</i>	(2b) ? <i>A man who (has) escaped from prison. had robbed a bank.</i>
	(2c) <i>The man who (has) escaped from prison had robbed a bank.</i>
<hr/>	
(3) ORDER: E1 > E2;      TOPIC: bank	
(3a) <i>A bank was robbed by a man. He had escaped from prison.</i>	(3b) <i>The bank was robbed by the man who had escaped from prison.</i>
<hr/>	
(4) ORDER: E1 < E2;      TOPIC: prison	
(4a) <i>From prison escaped a man. He had robbed a bank.</i>	(4b) * <i>From prison escaped the man who had robbed the bank.</i>
	(4c) ? <i>The prison from which the man who robbed the bank escaped . . .</i>
<hr/>	

E1 > E2 means: event-1 precedes event-2.

Topicalisation is another factor. While we can describe the scene from the man's point of view, we can't start linearisation from *the prison*: (4b) being an incorrect sentence, while (4c) is incomplete. Syntactic constraints like passivisability may play a role too. For example, the verb "to escape" can't be passivised.

Presuppositions are another factor. The following sentences, presuppose specific information concerning temporal order.

*The man who had robbed the bank escaped from prison.*

*The man who escaped from prison had robbed the bank.*

*The bank was robbed by the man who had escaped from prison.*

Put differently, order of events may impose special constraints on syntactic structure. The choice between *coordination* or *subordination* may de-

pend on conceptual information given with the input. Also, implicatures may vary depending on linear order. Compare Levelt's well known example (Levelt 1989):

*She became pregnant. They got married.*  
*They got married. She became pregnant.*

Likewise, compare the following two sentences, which express basically the same events, yet with different emphasis.

1. *Hitler has often been compared with Napoleon, although there are many differences between the two men.*
2. *Although there are many differences between Hitler and Napoleon, the two men have often been compared.*

While the first version focuses on the differences, the second stresses the similarities between the two men.

Figure 8 is similar to the preceding one; again the two clauses are symmetrical. However, this time the entity to be relativised plays the role of the *patient*.

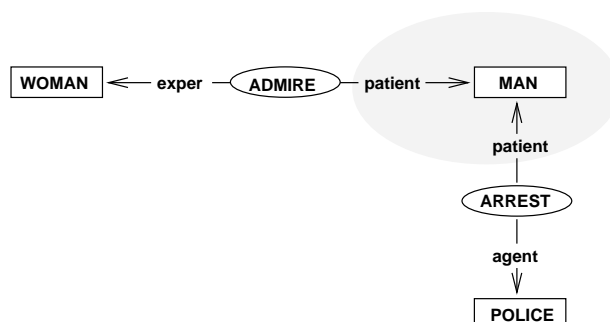


Fig. 8: *A man arrested by the police and admired by a woman*

---

Depending on certain temporal or attentional givens (focus), the structure in Figure 8 can be expressed in the following ways:

- 1a) *The man (whom) the woman admired was arrested by the police.*
- 1b) ? *The man (whom) the police arrested was admired by the woman.*
- 1c) ? *The man (who was) arrested by the police was admired by the woman.*
- 2a) *The woman admired the man whom the police arrested.*
- 2b) *The woman admired the man (who was) arrested by the police.*
- 3a) *The police arrested the man whom the woman admired.*

It should be noted, that this time, linearisation can be started from any node, that is, from a linguistic point of view there are no focus constraints. This is probably due to the fact that all the verbs used can be passivised. It may be worth mentioning however, that the passive voice occurs only in the main clause. Sentences (1b) and (1c), while not incorrect, sound distinctly odd.

The next structure (Figure 9) is different from the former in that the “man” plays a different role in each event: the structure is asymmetrical. In one case the “man” plays the role of the *agent*, whereas in the other he is the *patient*.

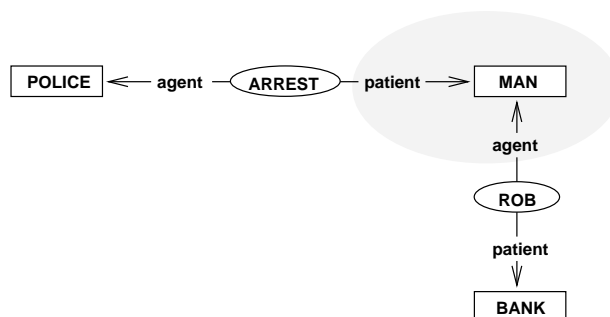


Fig. 9: *A man robbed a bank; He was arrested by the police*

---

Of course, this fact may reflect in the surface structure. According to the role played, the “man” will surface as *subject* or as *object* of the main verb.

Again there are constraints on the topicalisation, but for different reasons. The last sentence (3) sounds odd, as it gives the impression that the police had arrested the man before his robbing the bank. In general, it is not very good to subordinate a clause that expresses a consequence.

- 1) *The man who had robbed the bank was arrested by the police.*
- 2) *The police arrested the man who had robbed the bank.*
- 3) ? *The bank was robbed by the man whom the police arrested.*

## 9 Discussion

What can be learnt from looking at these networks? In order to recognise a candidate for a given syntactic structure one must consider several factors: the *type of concept* (predicate/argument), its *communicative status* (definite/indefinite), its *role* with regard to the whole (predicate dominating an argument or dominating another predicate), the *nature* (case role) and

*direction of the arcs* (incoming vs. outgoing). Yet, several other points are worth mentioning:

1. It is not enough to look just at one predicate or argument (local strategy), one has to look at larger chunks. Typically, the formal characteristics of the surrounding predicates or arguments also play a role. The relevant information being spread all over, one has to look at entire conceptual configurations.<sup>43</sup>
2. The formal characteristics (conceptual conditions) of the underlying conceptual structures are by no means sufficient for determining the syntactic form. They only suggest potential candidates. Other factors need to be taken into account, most prominently: the **size of conceptual chunks** to be verbalised,<sup>44</sup> **shared knowledge** (definiteness), **topicalisation** (active vs. passive voice, type of embedding), the **subcategorisation features** or **syntactic requirements** of a particular word,<sup>45</sup> and, last but not least, the **relative prominence** of each clause (saliency, focus), i.e., what shall be put into perspective, that is, be expressed by a main or subordinate clause? Syntactic structures are generally the result of an interaction between conceptual, linguistic and discourse choices.
3. The conceptual structure taken as input needs to contain a lot more information than the graphs shown here. Otherwise it is not possible to decide on a communicatively adequate syntactic structure.<sup>46</sup>

In conclusion, despite obvious correlations, the parallelism between conceptual structure and linguistic form is relative. There is no one-to-one mapping. While an *Agent-Action-Patient* structure is likely to be expressed in English by an S-V-O pattern, we cannot tell solely on these grounds that the speaker will render this idea in an active form. Similarly, the choice between a *that-clause* and an *infinitive* cannot always be made on purely conceptual grounds. The structure building properties (syntactic characteristics) of the verb must also be taken into account, all the more as the

---

<sup>43</sup> We believe that the idea of central vs. peripheral view (i.e., the idea that there is more to come) has some psycholinguistic reality.

<sup>44</sup> According to the amount of information the speaker tries to integrate into a sentence frame he may end up with several independent clauses, or one complex, heavily embedded clause (subordinate clauses).

<sup>45</sup> Not all verbs can be nominalised or passivised. Some words constrain other words (collocations), etc.

<sup>46</sup> A question that arises in that context is how complex a pattern may be, that is, how much information it can contain, without losing one of its most fundamental characteristics: recognisability.

choice of a particular verb may turn out to be incompatible with the chosen syntactic structure. In sum, one cannot strictly separate the syntax from the lexicon.

## 10 Conclusion

We have argued in this paper that human performance, that is, verbal fluency as observed in spontaneous discourse production, is only explicable if one hypothesises global strategies combined with pattern-matching: the speaker operates on larger chunks. We have also pointed out that for reasons of economy (storage), conceptual structures and linguistic structures ought to be parallel, at least to some extent (principle of structure preservation). Finally, we have outlined a strategy for fast computation of syntactic structures. We have shown that the same mechanism, — pattern matching, — could be used simultaneously for choosing words and (pre)computing syntactic structures. We have suggested that this be done in the following way. A given message is reduced to a lexicalised conceptual graph (lexicalisation). This graph serves as input for the (pre) computation of syntactic structure: ontological categories (entity, state, deep case roles, etc.) are replaced by syntactic categories (part of speech, syntactic functions). This preliminary syntactic structure is then checked against lexical subcategorisation features (for example passivisability of verbs) and the result is handed to the morphological component for final operations (agreement, insertion of function words).

We have claimed furthermore, that in order to precompute the syntactic structure, the user looks at the formal characteristics of the conceptual input.<sup>47</sup> It is by looking at cues like the *type of concept* (predicate vs. argument), its relative *position* with respect to the whole, *type and direction of the links* (ingoing vs. outgoing), etc., that s/he decides, i.e., precomputes, whether a given concept or conceptual chunk should map on, let us say, a noun, a verb, an adjective, an infinitival-, relative-, or that-clause, etc.

We would like to take this opportunity to clarify here our position with regard to formal grammars and incremental processing. By suggesting to map lexicalised conceptual structures directly onto (final) linguistic forms, i.e., syntactic structures, we may have given the reader the wrong impression that one could do without a formal grammar. Obviously, this is not quite so. There is still assembly, and the legal combinations, i.e., pieces that

---

<sup>47</sup> This makes generally sense only once lexicalisation has taken place.



can go together, are still specified by a formal grammar. For example, if we used unification, than we would simply try to unify larger chunks than in most other approaches. As one can see, we do not mean to bypass the grammar, we simply mean that, depending on the situation and the speaker's proficiency, grammaticalisation is performed on larger chunks.

This last statement should also clarify our position with regard to incremental processing. We do not mean to criticise the basic idea behind it, quite to the contrary. We simply challenge the view of word-to-word processing. Indeed, there may be a whole spectrum of units, going from very small items (words or even below) to fairly large chunks (fixed expressions, patterns). The size of the unit, and the way to process them may vary depending on the user and the situation (cognitive states).

One major criticism though towards formal grammars. According to our knowledge they do not make explicit the correspondances, or mappings between conceptual structures and linguistic forms (see Figure 1).

Hence, formal grammars miss a link. Yet, this link is fundamental and lies at the heart of our approach. Obviously, many details are still lacking, especially concerning the size of the units, the recovery strategies (what to do in case of failure), and last, but not least the mapping rules. Actually, so far we have shown only the tip of the iceberg. We need to make a list of the possible conceptual structures and their linguistic counterparts, and we have to specify the mapping rules and their constraints. Despite all these shortcomings, and despite the fact that our approach still lacks formal treatment and the test of implementation, — though the PROTECTOR generation system (Nicolov, Mellish & Ritchie 1997) is a serious step in that direction,<sup>48</sup> — it embodies in principle at least two interesting facts: procedural knowledge can be stated explicitly via the mapping rules; processing can be speeded up by operating on larger chunks rather than atomic units.

**Acknowledgements.** The author would like to express his gratitude to all those who were so kind to comment on the initial draft: Dominique Estival, Aravind Joshi, Guy Lapalme, Yves Lepage, William Levelt, Terry Patten, Alain Polguère, Ehud Reiter and Dan Tufis. Special thanks go to Nicolas Nicolov, who has devoted a considerable amount of time for the long discussions through which a lot of important theoretical points were clarified and who helped with editing the manuscript.

---

<sup>48</sup> The use of D-Tree Grammars in PROTECTOR allow the system to operate on larger conceptual units.

## REFERENCES

- Aitchison, Jean. 1983. *The Articulate Mammal: an Introduction to Psycholinguistics*. London: Hutchinson.
- Anderson, John. 1983. *The Architecture of Cognition*. Cambridge, Mass.: Harvard University Press.
- Arbib, Michael, E. Conklin & J. Hill. 1986. *From Schema Theory to Language*. New York: Oxford University Press.
- Ausubel, David. 1980. "Schemata, Cognitive Structure and Advance Organisers: A Reply to Anderson, Spiro and Anderson". *American Educational Research Journal* 17:3.400-404.
- Bartlett, Frederik. 1932. *Remembering*. Cambridge, U.K.: Cambridge University Press.
- Becker, Joseph. 1975. "The Phrasal Lexicon". BBN Report No. 3081. Cambridge, Mass.: Bolt Beranek & Newman.
- Bobrow, Daniel. 1968. "Natural Language Input for a Computer Problem Solving System" *Semantic Information Processing* ed. by Marvin Minsky, 33-145. Cambridge, Mass.: MIT Press.
- & G. Norman. 1975. "Some Principles of Memory Schemata". *Representation and understanding: studies in cognitive science* ed. by Daniel G. Bobrow & A. M. Collins, 31-149. New York, U.S.A.: Academic Press.
- Bock, Kathryn, H. Loebell & R. Morey. 1992. "From Conceptual Roles to Structural Relations: Bridging the Syntactic Cleft". *Psychological Review* 99:1.150-171.
- Brown, Roger 1958. "Linguistic Determinism and the Part of Speech". *Psychological Review* 65:1.14-21.
- Bruner, Jerome. 1973. *Beyond the Information Given: Studies in the Psychology of Knowing* ed. by J. Anglin. New York, U.S.A.: W.Norton.
- Burton, Richard. 1976. "Semantic grammar: an Engineering Technique for Constructing Natural Language Understanding Systems". BBN Report No. 3453. Cambridge, Mass.: Bolt Beranek & Newman.
- Chomsky, Noam. 1959. "Review of Verbal Behavior, by B.F. Skinner". *Language* 35:26-58.
- Clark, Herbert & E. Clark. 1977. *Psychology and Language: An Introduction to Psycholinguistics*. New York, U.S.A.: Harcourt, Brace Jovanovich.
- & S. Haviland. 1977. Comprehension and the Given-New Contract. *Discourse Production and Comprehension* ed. by R.O. Freedle, 1-40. N.J., U.S.A.: Norwood.

- de Smedt, Koenrad. 1990. *Incremental Sentence Generation: A Computer Model of Grammatical Encoding*, Ph.D. dissertation, (also NICI TR 90-01). Nijmegen, The Netherlands: Nijmegen Institute for Cognition Research and Information Technology.
- van Dijk, Teun. 1977. "Semantic Macro-Structures and Knowledge frames in Discourse Comprehension". *Cognitive Processes in Comprehension* ed. by M. A. Just & P. A. Carpenter, 3-31. N.J., U.S.A.: Hillsdale.
- Fillmore, Charles. 1977. "Scenes-and-frames Semantics". *Linguistic Structures Processing* ed. by Antonio Zampolli, 55-82. Amsterdam: North Holland.
- Friedman, Daniel & M. Feldeisen. 1987. *The little LISP*. Cambridge Mass.: MIT Press.
- Fries, Charles. 1952. *The Structure of English. An Introduction to the Construction of English Sentences*. New York: Harcourt, Brace.
- Goffman, Erving. 1974. *Frame analysis. An Essay on the Organisation of Experience*. Cambridge, Mass.: Harper & Row.
- Habel, Christopher. 1988. "Cognitive Linguistics: The Processing of Spatial Concepts". LILOG Report 45. Stuttgart, Germany: IBM.
- Harris, Zellig. 1951. *Methods in Structural Linguistics*. Chicago: University of Chicago Press.
- Hendrix, Gary. 1977. "The LIFER Manual: A Guide to Building Practical Natural Language Interfaces". Technical Note 138, Menlo Park: SRI.
- Hill, Jane & M. Arbib. 1984. "Schemas, Computation and Language Acquisition". *Human Development* 27:282-296.
- Hovy, Edward. 1990. "Unresolved Issues in Paragraph Planning". *Current Research in Natural Language Generation* ed. by Robert Dale, Chris Mellish & Michael Zock, 17-41. London: Academic Press.
- Kant, Emmanuel. 1781. *Critique of Pure Reason*, translated by Max Müller. Garden City, N.Y.: Anchor Books.
- Kempen, Gerard. 1977. "Conceptualising and Formulating in Sentence Production". *Sentence Production: Developments in Research and Theory* ed. by S. Rosenberg, 259-274. Hillsdale, N.J.: Erlbaum.
- & E. Hoenkamp. 1987. "An Incremental Procedural Grammar for Sentence Formulation". *Cognitive Science* 11:201-258.
- Koffka, Kurt. 1935. *The principles of Gestalt psychology*. New York: Harcourt Brace & World.
- Lado, Robert. 1964. *Language Teaching: A Scientific Approach*. New York: McGraw Hill.
- Langacker, Ron. 1983. *Foundations of Cognitive Grammar I & II*. Bloomington: Indiana University Linguistics Club.

- Levelt, William. 1981. "The Speaker's Linearisation Problem". *Phil. Transactions Royal. Soc.* London B295,305-315.
- . 1982. "Linearization in Describing Spatial Networks". *Processes, Beliefs and Questions* ed. by Stanley Peter & E. Saarinen. Dordrecht, Holland: Reidel.
- . 1989. *Speaking*. Cambridge, Mass., U.S.A.: MIT Press.
- Mandler, Jean. 1979. "Categorial and Schematic Organisation in Memory". *Memory Organisation and Structure* ed. by C. Puff, 259-299. New York: Academic Press.
- McKeown, Kathleen. 1984. *Text Generation*. Cambridge, U.K.: Cambridge University Press.
- Mel'cuk, Igor. & A. Zholkovskij. 1970. "Towards a Functioning 'Meaning-Text' Model of Language". *Linguistics* 57:10-47.
- Miller, George. 1956. "The Magical Number Seven, Plus or Minus Two: Limits on our Capacity for Processing Information". *Psychological Review* 63:2.81-97.
- & Philip Johnson-Laird. 1985. *Language and Perception*. Cambridge, U.K.: Cambridge University Press.
- Minsky, Marvin. 1975. "A Framework for Representing Knowledge". *The Psychology of Computer Vision* ed. by Patrick Winston, 211-277. New York: McGraw Hill.
- . 1985. *The Society of Mind*. New York, U.S.A.: Simon & Schuster.
- Nicolov, Nicolas, Chris Mellish & Graeme Ritchie. 1997. "Approximate Generation from Non-Hierarchical Representations". *Recent Advances in Natural Language Processing* ed. by Ruslan Mitkov & Nicolas Nicolov, Amsterdam: John Benjamins.
- Nogier, Jean François. 1991. *Génération Automatique de Langage et Graphs Conceptuels*. Paris: Hermès.
- & Michael Zock. 1992. "Lexical choice by pattern matching". *Knowledge Based Systems* 5:3.200-212 (also in *Current Directions in Conceptual Structures Research* ed. by T. Nagle, J. Nagle, L. Gerholz & P. Eklund, 413-435, Berlin & New York: Springer Verlag).
- Olson, David. 1970. "Language and Thought: Aspects of a Cognitive Theory of Semantics". *Psychological Review* 77:257-273.
- Osgood, Charles. 1971. "Where do sentences come from?". *Semantics: An Interdisciplinary Reader in Philosophy, Linguistics, and Psychology*, ed. by D. Steinberg & L. Jakobovits, 497-529. Cambridge: Cambridge University Press.
- . 1980. *Lectures on Language Performance*. New York: Springer Verlag.

- Paivio, Allan. 1971. *Imagery and Verbal Processes*. New York: Holt, Rinehart & Winston.
- Patten, Terry, Michael Geis & Barbara Becker. 1992. "Toward a Theory of Compilation for Natural Language Generation". *Computational Intelligence* 8:1.77-101.
- Piaget, Jean. 1970. "Piaget's Theory". *Carmichael's Manual of Child Psychology* ed. by P. Mussen, vol.I, 318-323, New York: Wiley.
- Raphael, Bertram. 1968. "SIR: A computer Program for Semantic Information Retrieval". *Semantic Information Processing* ed. by Marvin Minsky, 146-226. Cambridge, Mass.: MIT Press.
- Rivers, Wilga. 1972. *Speaking in Many Tongues: Essays in Foreign Language Teaching*. Rowley: Newbury House.
- Roberts, Paul. 1962. *English Sentences*. New York: Harcourt Brace & World.
- Rösner, Dietmar. 1987. "The Automated News Agency: SEMTEX - A Text Generator for German". *Natural Language Generation: New Results in Artificial Intelligence, Psychology and Linguistics* ed. by Gerard Kempen, 133-148. Dordrecht: Martinus Nijhoff Publishers.
- Rumelhart, David. 1975. "Notes on a Schema for Stories". *Representation and Understanding: Studies in Cognitive Science* ed. by Daniel G. Bobrow & A. M. Collins, 211-236. New York: Academic Press.
- & A. Ortony. 1976. "The Representation of Knowledge in Memory". *Schooling and the Acquisition of Knowledge* ed. by R. C. Anderson, R. J. Spiro & W. E. Montague, 99-133. Hillsdale, N.J.: Erlbaum.
- Schank, Roger & R. Abelson. 1977. *Scripts, Plans, Goals and Understanding: an Inquiry into Human Knowledge Structures*. N.J., U.S.A.: Hillsdale.
- Simmons, Robert & J. Slocum. 1972. "Generating English Discourse from Semantic Networks". *Communications of the Association for Computing Machinery (CACM)* 15:10.891-905.
- Skinner, Burrhus. 1957. *Verbal Behavior*. New York, U.S.A.: Appleton-Century-Crofts.
- Sowa, John. 1984. *Conceptual Structures: Information Processing in Mind and Machine*. Reading, Mass.: Addison Wesley.
- Stockwell, Robert. 1977. *Foundations of Syntactic Theory*. N.J., U.S.A.: Prentice Hall.
- Swartout, William. 1983. "XPLAIN : A System for Creating and Explaining Expert Consulting Systems". *Artificial Intelligence* 21:3.285-325.
- Tannenbaum, Percy & F. Williams. 1968. "Generation of Active and Passive Sentences as a Function of Subject or Object focus." *Journal of Verbal Learning and Verbal Behavior* 7:246-250.

- Vennemann, Theo. 1975. "An Explanation of Drift". *Word Order and Word Order Change* ed. by C. Li, 269-305. Austin, Texas: University of Texas Press.
- Weizenbaum, Joseph. 1966. "ELIZA. A Computer Program for the Study of Natural Language Communication between Man and Machine". *Communications of the Association for Computing Machinery (CACM)* 9:36-45.
- Wilks, Yorick. 1975. "A Preferential Pattern-Seeking Semantics for Natural Language Inference". *Artificial Intelligence* 6:1.53-74.
- Winograd, Terry. 1972. *Understanding Natural Language*. New York: Academic Press.
- . 1975. "Frame Representation and the Declarative-Procedural Controversy". *Representation and Understanding: Studies in Cognitive Science* ed. by Daniel G. Bobrow & A. M. Collins, 185-210. New York: Academic Press.
- Zock, Michael, Gerard Sabah & C. Alviset. 1986. "From Structure to Process: Computer-assisted Teaching of Various Strategies for Generating Pronoun-Constructions in French". *Proceedings of the 11th International Conference on Computational Linguistics (COLING-86)*, 566-570. Bonn, Germany.
- Zock, Michael. 1988. "Natural Languages are Flexible Tools, That's What Makes Them Hard to Explain, to Learn and to Use". *Advances in Natural Language Generation: an Interdisciplinary Perspective* ed. by Michael Zock & Gerard Sabah, 181-196. London: Pinter.
- . 1990. "If you Can't Open the Black Box, Open a Window! A Psycholinguistically-Motivated Architecture of a Natural Language Generation Component". *Proceedings of COGNITIVA-90*, 143-152. Madrid, Spain.
- . 1994. "Language in Action, or, Learning a Language by Watching It Work". *Proceedings of the 7th Twente Workshop on Language Technology: Computer-Assisted Language Learning*, 101-111. Twente, The Netherlands.
- . 1996. "The Power of Words in Message Planning". *Proceedings of the 16th International Conference on Computational Linguistics (COLING'96)*. Copenhagen, Denmark.

## Index

- correspondences, 4
  - many-to-many, 18
  - one-to-one, 18
- d-tree grammar (DTG), 32
- focus constraints, 29
- French language, 3
- generation strategies, 14
- global strategy, 12
- hybrid knowledge representation,
  - 2
- linearisation, 27
- linguistics
  - cognitive, 9
  - structural, 3, 9
- mapping rules, 1, 3, 7, 9, 23
- natural language generation (NLG),
  - 1
- parallelism-correlations, 12, 13, 30
- pattern matching, 1, 2, 8, 12
- patterns
  - prototypical, 3, 21
- presupposition, 30
- PROTECTOR, 32
- sentence generation, 12
- structures
  - conceptual, 3, 4, 6, 12, 30
  - linguistic, 3, 4, 11, 12
  - syntactic, 6, 18, 30
- syntactic patterns
  - recognition of, 1, 4, 9, 12, 20,
    - 22
- thematic roles, 6
- topicalisation constraints, 27, 30