# GRAMMAR AS MATHEMATICS

BY

## J. LAMBEK

*Dedicated to the memory of William Lloyd Garrison Williams.*

ABSTRACT. While there are a small number of reasonably deep theorems in mathematical linguistics, I wish to argue that grammar is mathematics at a very basic level, albeit "trivial" mathematics. Linguistic activities such as the production and recognition of sentences are quite analogous to the mathematical activities of proving theorems or making calculations, while learning a language involves something akin to the discovery or invention of postulates.

**Preamble.** I feel greatly honoured to have been chosen to give the Jeffery-Williams lecture, particularly as I was a friend of both Ralph Jeffery, who used to look over my shoulder at the summer research institute in Kingston, and Lloyd Williams, who had encouraged me to become a professional mathematician in the first place and often invited me to his home when I was but a humble undergraduate. I take this opportunity to talk on a subject which not everyone believes to be a part of mathematics (e.g. the NSERC Grant Selection Committee), but which Lloyd would have approved of, as he himself had gone to Oxford as a Rhode scholar after having majored in Classics.

0. **Intoduction.** Pythagoras, who coined the word "mathematics", was supposed to have said "everything is number". Most mathematicians would concede that at least all things in mathematics are numbers, although they might mumble under their breath: "or sets of numbers or sets of sets of numbers ...". Anyway, according to Pythagoras, "mathematics" (the word is related to the English word "mind") covered four subjects: arithmetic, geometry, astronomy and music. While we might have our doubts today about music, we could be persuaded that astronomy is "applied" mathematics. Moreover, ever since Descartes, we accept the reduction of geometry to "real" numbers, and these, according to Eudoxus and Dedekind, may be viewed as sets of rational numbers, which are the numbers Pythagoras had in mind.

Conspicuously absent from his classification is grammar. But wait, the four mathematical subjects of Pythagoras re-appear in the medieval university curriculum as the four liberal arts, the so-called "quadrivium". They are preceded by the "trivium" consisting of the three subjects: grammar, logic and rhetoric, hence labelled "trivial". No one would claim that rhetoric is mathematics, but logic has by now been

accepted as respectable mathematics almost everywhere and grammar is beginning to gain recognition at certain European universities as "mathematical linguistics" and, even in Montreal, I have been teaching an undergraduate course "computability and mathematical linguistics" for the last twenty years or so.

Actually, grammar deals with finite strings of symbols taken from a finite alphabet. According to Cantor, these strings are in one-to-one correspondence with the natural numbers, an idea which has been exploited by Gödel in the proof of his famous incompleteness theorem. However, my claim that grammar is mathematics is not based on Gödel's arithmetization, but rather on the observation that proof and computation are as fundamental in language as in mathematics.

If my position is controversial, at least my interest in language has been shared by many mathematicians, including Lloyd Williams, and I know of three who have made contributions to linguistics in the past: Eratosthenes, Wallis and Grassmann.

The kind of grammar I favour nowadays has been called by various names: semi-Thue systems, rewriting systems, productions grammars. In mathematics, semi-Thue systems were studied in connection with the word problem for semigroups (see e.g. Kleene 1952).

For our purposes, a *semi-Thue system* is a finitely generated free monoid together with a finitely generated pre-order on it. Thus, we are given a finite set $\mathcal{V}$ of symbols, usually called the *vocabulary* or *alphabet*, and a finite subset $\mathcal{P}$ of $\mathcal{V}^* \times \mathcal{V}^*$, called the set of *productions*. Here $\mathcal{V}^*$ is the free monoid generated by $\mathcal{V}$; it consists of finite strings of elements of $\mathcal{V}$, including the empty string 1, multiplication being just concatenation. From $\mathcal{P}$ we generate a pre-order relation on $\mathcal{V}^*$ with the help of the following axioms and rules of inference:

$$\Gamma \longrightarrow \Delta \quad \text{when } (\Gamma, \ \Delta) \in \mathcal{P} \quad \text{(productions)},$$

$$\Gamma \longrightarrow \Gamma \qquad \text{(reflexive law)},$$

$$\frac{\Gamma \longrightarrow \Delta \quad \Delta \longrightarrow \Lambda}{\Gamma \longrightarrow \Lambda} \text{(transitive rule)},$$

$$\frac{\Gamma \longrightarrow \Delta \quad \Gamma' \longrightarrow \Delta'}{\Gamma\Gamma' \longrightarrow \Delta\Delta'} \text{(substitution rule)}.$$

Thus, we may write $\Gamma \leqq \Delta$ to mean that $\Gamma \longrightarrow \Delta$ is provable in the above deductive system. We have used capital Greek letters to denote strings of symbols.

The substitution rule is easily seen to be equivalent to the following:

$$\frac{\Gamma \longrightarrow \Delta}{\Phi\Gamma\Psi \longrightarrow \Phi\Delta\Psi} \ .$$

As the free monoid also contains the empty string, special cases of this rule are:

$$\frac{\Gamma \longrightarrow \Delta}{\Phi\Gamma \longrightarrow \Phi\Delta} \ , \qquad \frac{\Gamma \longrightarrow \Delta}{\Gamma\Psi \longrightarrow \Delta\Psi} \ .$$

A *production grammar* $\mathcal{G} = (\mathcal{V}, \mathcal{P}, \mathcal{V}_i, \mathcal{V}_t)$ is a semi-Thue system in which two subsets of $\mathcal{V}$ have been specified: the *initial* vocabulary $\mathcal{V}_i$ and the *terminal* vocabulary $\mathcal{V}_t$. Note the symmetry of our definition: if $\mathcal{P}$ is replaced by its converse, hence

all arrows are reversed, and if $\mathcal{V}_i$ and $\mathcal{V}_t$ are interchanged, we get another production grammar, called the *dual* of $\mathcal{G}$.

For example, if we wish to generate English sentences, we might take $\mathcal{V}_t$ to be the set of all English words, say those listed in the Oxford dictionary, together with all their inflected forms, and $\mathcal{V}_i$ as the set $\{S,Q,C, \ldots\}$, where S = statement, Q = question, C = command et cetera are the different types of sentences we wish to generate. Assuming $\mathcal{V}$ and $\mathcal{P}$ to have been chosen judiciously, if $\Gamma$ is a string of English words, we expect to be able to prove $S \rightarrow \Gamma$ then and only then when $\Gamma$ is a grammatical declarative sentence.

The dual of the grammar just considered will serve to analyze certain strings of English words as different types of sentences. For example, in this dual grammar we might expect to prove

$$does\ he\ come \rightarrow Q.$$

If we are interested in translating from French to English, we might construct a production grammar in which $\mathcal{V}_i$ is the French vocabularly and $\mathcal{V}_t$ the English vocabulary. In such a grammar we would expect to be able to prove

$$il\ va \rightarrow he\ goes.$$

The translation grammar will of course be much more diffcult to construct than the separate grammars of the two languages.

In this brief exposition we shall discuss in detail three examples of languages and their grammars. The first is a formal language which has been specially constructed to illustrate a technical point, namely that the intersection of two contextfree languages need not be contextfree. The second is a fragment of English dealing with consanguineous kinship terminology, along lines which anthropologists have adopted to look at more exotic kinship systems. The third is a fragment of English which exhibits the structure of basic sentences such as *he likes her*. Our grammar should account not only for the syntax, but also for the morphology. In particular, it should not be able to generate the ungrammatical *\*he like she*. (Linguists usually put a star on incorrect forms.)

Let me emphasize that my main interest lies in the grammars of natural languages and the "trivial" theorems they contain, even though there exist many nontrivial mathematical theorems, or should I say "metatheorems", about formal languages. For the latter, the reader is referred to the beautiful exposition by Cohn (1975), who also discusses our first example below.

**1. A formal language which is not contextfree.** In examples of formal languages, one usually assumes that $\mathcal{V}_i = \{S\}$ consists of a single symbol S. The sentences of such a language are then all strings $\Gamma \in \mathcal{V}_t{}^*$ for which $S \rightarrow \Gamma$ is provable in the deductive system. For example, consider the terminal vocabulary $\mathcal{V}_t = \{a, b, c\}$, from which we want to generate the language

$$\mathcal{L} = \{a^k b^k c^k \mid k \geqq 1\}.$$

We take $\mathcal{V}_i = \{S\}$ and $\mathcal{V} = \{a,b,c,S,B\}$, with one auxiliary symbol B, and the following productions:

$$S \rightarrow abc,$$

$$S \rightarrow aBSc,$$

$$Ba \rightarrow aB,$$

$$Bb \rightarrow bb.$$

Indeed, here is a proof that $S \rightarrow$ aabbcc:

$$\underline{S \rightarrow abc} \qquad \underline{Ba \rightarrow aB}$$

$$\underline{BS \rightarrow Babc} \qquad \underline{Babc \rightarrow aBbc}$$

$$\underline{BS \rightarrow aBbc}$$

$$\underline{S \rightarrow aBSc} \qquad \underline{aBSc \rightarrow aaBbcc} \qquad \underline{Bb \rightarrow bb}$$

$$\underline{S \rightarrow aaBbcc} \qquad \qquad \underline{aaBbcc \rightarrow aabbcc}$$

$$S \rightarrow aabbcc$$

Such proofs in tree-form are a bit cumbersome and, in practice, they are replaced by so-called *derivations*, as follows:

$$S \rightarrow aBSc \rightarrow aBabcc \rightarrow aaBbcc \rightarrow aabbcc.$$

Here, the first step is the production $S \rightarrow aBSc$, the second step uses the production $S \rightarrow abc$ in the context aB–c, the third step uses the production $Ba \rightarrow aB$ in the context a-bcc ad the fourth step uses the production $Bb \rightarrow bb$ in the context aa-cc.

Alternatively, one may depict such a proof by a diagram:

$$
\begin{array}{ccc}
\multicolumn{3}{c}{\underline{\qquad\qquad S \qquad\qquad}} \\
aB & S & c \\
 & \overline{abc} & \\
\underline{aB} & & \\
\overline{bb} & &
\end{array}
$$

Linguists are fond of so-called *contextfree* productions, namely productions of the form $A \rightarrow \Gamma$ where the left side consists of a single symbol. (Sometimes one also insists that $\Gamma$ is not the empty string.) Of the productions listed above, only the first two are contextfree. Still, one might hope that another set of contextfree productions

will generate the same language. That this is not so follows from a famous result by Bar-Hillel, Perlis and Shamir, which we shall not prove here, as it is easily found in the literature (see e.g. Cohn 1975).

LEMMA (Bar-Hillel et al.) A necessary condition for an infinite language to have a contextfree grammar is that there should exist strings $\Gamma, \Delta, \Lambda, \Phi$ and $\Psi$, where $\Phi$ and $\Psi$ are not both empty, such that $\Gamma \Phi^n \Delta \Psi^n \Lambda$ is a sentence for each $n \geqq 1$.

To illustrate this condition, consider the fragment of English built from the terminal vocabulary $\{I, know, and\}$ and the productions

$$S \rightarrow I \ know, \qquad S \rightarrow S \ and \ S.$$

Now take $\Gamma = I \ know$ and $\Phi = \Delta = \Psi = \Lambda = and \ I \ know$. Let us apply the lemma to the language $\mathcal{L}$ under consideration in this Section and assume that it is contextfree. The cases $n = 1$ and $n = 2$ of the lemma require that there exist strings $\Gamma, \Phi, \Delta, \Psi$ and $\Lambda$ in $\{a,b,c\}^*$, with $\Phi$ and $\Psi$ not both empty, and integers $p, q \geqq 1$ such that

$$\Gamma \Phi \Delta \Psi \Lambda = a^p b^p c^p,$$

$$\Gamma \Phi \Phi \Delta \Psi \Psi \Lambda = a^q b^q c^q.$$

A straight-forward but tedious calculation shows that this leads to a contradiction.

It follows that the language $\mathcal{L}$ is not contextfree. On the other hand, it is easily seen to be the intersection of two contextfree languages:

$$\mathcal{L}_1 = \{a^m b^m c^n | m, n \geqq 1\},$$
$$\mathcal{L}_2 = \{a^m b^n c^n | m, n \geqq 1\}.$$

For example, $\mathcal{L}_1$ is generated by the productions

$$U_1 \rightarrow ab, U_1 \rightarrow aU_1 b,$$
$$V_1 \rightarrow c, V_1 \rightarrow V_1 c,$$
$$S_1 \rightarrow U_1 V_1,$$

where $S_1$ is the initial symbol.

2. **English consanguineous kinship terminology.** Kinship terminologies have been studied by anthropologists, who have been concentrating largely on languages spoken by isolated societies (see e.g. Buchler and Selby 1968). The usefulness of production grammars in this connection was first pointed out by Lounsbury (1965). The original emphasis was on reduction and equivalence rules, which may be illustrated in English by noting that a second cousin is viewed simply as a cousin and that the brother's wife is denoted by the same term as the husband's sister. These rules of primary interest to anthropologists were integrated with structural and morphological rules in

a series of articles written in collaboration with Michael Lambek (1981), who initiated the project, and Mira Bhargava (1983, 1985). Some fruitful ideas were contributed by George Bergman.

We shall here concentrate on the kinship terminology in the idiolect of Goodenough (1965); but, for brevity, we shall ignore any relations through marriage. The present treatment will differ slightly from that in my 1986 paper, by placing more emphasis on contextfree production, at the cost of increasing the number of productions.

We shall adopt the vocabulary:

$$\mathcal{V} = \mathcal{V}_i \smile \mathcal{V}_t \smile \{M, F, P, C, S, U, V, W, [, ], \}$$
$$\mathcal{V}_i = \{R\},$$
$$\mathcal{V}_t = \{father, mother, son, daughter, brother, sister, uncle, aunt, nephew, niece,$$
$$cousin, grand, great, \#\}$$

The reader may think of some of the auxiliary symbols as abbreviations:

$$R = \text{relation},$$
$$M = \text{male},$$
$$F = \text{female},$$
$$P = \text{parent (of)},$$
$$C = \text{child (of)},$$
$$S = \text{sibling (of)}.$$

Also we shall use U for descendant, V for ancestor and W for any relation which involves S. The terminal vocabulary includes the prefix *grand* and the symbol #, which is supposed to denote a space between words.

Before we postulate the productions, we take a look at the data which they should account for. The following table summarizes the consanguineous data discussed by Goodenough in 1965 for a certain dialect of American English, with one minor difference: we have spelled *grandnephew* in one word.

TABLE I

| | | | | | |
|---|---|---|---|---|---|
| $MSP^{n+1}$ | $\rightarrow$ | $(great\#)^n uncle$ | $FSP^{n+1}$ | $\rightarrow$ | $(great\#)^n aunt$ |
| $MP^{n+2}$ | $\rightarrow$ | $(great\#)^n grandfather$ | $FP^{n+2}$ | $\rightarrow$ | $(great\#)^n grandmother$ |
| $MP$ | $\rightarrow$ | $father$ | $FP$ | $\rightarrow$ | $mother$ |
| $MS$ | $\rightarrow$ | $brother$ | $FS$ | $\rightarrow$ | $sister$ |
| $MC^{m+1}SP^{n+1}$ | $\rightarrow$ | $cousin$ | $FC^{m+1}SP^{n+1}$ | $\rightarrow$ | $cousin$ |
| $MC$ | $\rightarrow$ | $son$ | $FC$ | $\rightarrow$ | $daughter$ |
| $MCS$ | $\rightarrow$ | $nephew$ | $FCS$ | $\rightarrow$ | $niece$ |
| $MC^{m+2}$ | $\rightarrow$ | $(great\#)^m grandson$ | $FC^{m+2}$ | $\rightarrow$ | $(great\#)^m granddaughter$ |
| $MC^{m+2}S$ | $\rightarrow$ | $(great\#)^m grandnephew$ | $FC^{m+2}S$ | $\rightarrow$ | $(great\#)^m grandniece$ |

The entries in Table I are self-explanatory. For example, $MSP^2$ is a kinship description to be read as 'male sibling of parent of parent" and the table tells us that it is rendered in English as *great # uncle*.

Concerning the kinship description of *cousin*, I have avoided the more precise rendering

$$GC^{m+1} SP^{n+1} \rightarrow \text{\textit{i-th cousin j times removed,}}$$

where $G = M$ or $F$ and

$$i = \min (m, n) + 1, \quad j = |m - n|,$$

as I find it difficult to calculate these two primitive recursive functions in my head.

We wish to postulate productions which will account for the data in Table I. For example, the productions should predict that

$$MSP^2 \rightarrow \text{\textit{great \# uncle}} \text{ in the context \# } - \text{ \#,}$$

assigning a kinship term to a kinship description. They should also predict that

$$R \rightarrow MSP^2 \text{ in the context \# } - \text{ \#,}$$

to account for the fact that $MSP^2$ is a kinship description in the first place and not, for example, *MPSP. Furthermore, we hope that they will explain the asymmetry between *grandnephew* and *great # uncle* and also account for the absence of *grandcousin*. Here then are our productions:

(1) *Structure rules.*
    $U \rightarrow C, U \rightarrow CU$;
    $V \rightarrow P, V \rightarrow VP$;
    $W \rightarrow S, W \rightarrow CW, W \rightarrow WP$;
    $R \rightarrow GU, GV, GW$ in the context $\# - \#$, where $G = M$ or $F$.

(2) *Reduction rules.*
    $C^2 SP \rightarrow CSP$,
    $CSP^2 \rightarrow CSP$.

(3) *Word assignments,* before # and after # or [*grand*].

| | | | |
|---|---|---|---|
| MSP | $\rightarrow$ *uncle* | FSP | $\rightarrow$ *aunt,* |
| MP | $\rightarrow$ *father,* | FP | $\rightarrow$ *mother,* |
| MS | $\rightarrow$ *brother,* | FS | $\rightarrow$ *sister,* |
| MCSP | $\rightarrow$ *cousin,* | FCSP | $\rightarrow$ *cousin,* |
| MC | $\rightarrow$ *son,* | FC | $\rightarrow$ *daughter,* |
| MCS | $\rightarrow$ *nephew,* | FCS | $\rightarrow$ *niece.* |

(4) *Prefix assignments*, where $G = $ M or F.

$$GP^2 \rightarrow [grand]GP;$$
$$GC^2 \rightarrow [grand]GC;$$
$$[grand] \rightarrow \begin{cases} great \ \# \ \text{before } grand \text{ or } great, \\ 1 \text{ before } cousin, \\ grand \text{ before other English words}; \end{cases}$$
$$GSP^2 \rightarrow \text{great} \ \# \ GSP.$$

We shall now discuss these productions.

(1) The structure rules account for:

$$U \rightarrow C^{m+1}, \quad V \rightarrow P^{n+1}, \quad W \rightarrow C^m SP^n;$$
$$R \rightarrow GC^{m+1}SP^{n+1}, GC^{m+1}(S), G(S)P^{n+1}, GS \text{ in the context } \# - \#;$$

where $m, n, \geqq 0$ and $G = $ M or F.

(2) The reduction rules account for

$$C^{m+1}SP^{n+1} \rightarrow CSP.$$

Rules of this kind had been first emphasized by Lounsbury (1965); they are more prominent in some exotic languages than in English.

(3) The word assignments allow us e.g. to rewrite MSP as *uncle* in the context $\# - \#$, but not, for example, $MSP^2$ as *uncle*P.

(4) Structural linguists might consider [*grand*] as a morpheme with three allomorphs: *great*#, *grand* and 1, the empty string.

We shall present some sample calculations, assuming everywhere the context $\# - \#$.

$$MC^4 \rightarrow [grand]MC^3 \rightarrow [grand]^2MC^2 \rightarrow [grand]^3MC$$
$$\rightarrow [grand]^3 son \rightarrow [grand]^2 grandson$$
$$\rightarrow [grand] \ great\# \ grandson \rightarrow \ great\# \ great\# \ grandson.$$
$$MC^3S \rightarrow [grand]MC^2S \rightarrow [grand]^2MCS$$
$$\rightarrow [grand]^2 \ nephew \rightarrow [grand] \ grandnephew$$
$$\rightarrow great\# \ grandnephew.$$
$$MC^2SP^2 \rightarrow [grand] \ MCSP^2 \rightarrow [grand]^2 \ MCSP$$
$$\rightarrow [grand]^2 \ cousin \rightarrow [grand] \ cousin \rightarrow cousin;$$

or more quickly:

$$MC^2SP^2 \rightarrow MCSP^2 \rightarrow MCSP \rightarrow \textit{cousin}.$$
$$MSP^3 \rightarrow \textit{great}\#MSP^2 \rightarrow \textit{great}\#\textit{great}\#MSP$$
$$\rightarrow \textit{great}\#\textit{great}\#\textit{uncle}.$$

Note that *grandnephew* is obtained from the same rule which already accounts for *grandson*, but that the kinship description of *great*# *uncle* cannot be evaluated by the rule which yields *grandfather*; it requires a new rule.

Webster's New Ideal Dictionary lists *great-nephew* as an alternative to *grand-nephew* and, in place of Goodenough's *great*#*uncle*, it offers a choice between *great-uncle* and *granduncle*. Webster's idiolect could be accounted for by altering the last two rules in (4) as follows:

$$[grand] \rightarrow grand \text{ or } \textit{great-} \text{ before other English words;}$$
$$GSP^2 \rightarrow [grand] \, GSP.$$

Why are there no grandcousins? We could try to introduce them by deleting the rule

$$[grand] \rightarrow 1 \qquad \text{before } \textit{cousin}$$

in (4) above. We could then compute

$$MC^2SP \rightarrow [grand]MCSP \rightarrow [grand]cousin \rightarrow \textit{grandcousin}.$$

However, there is no way of converting $MCSP^2$ to *grandcousin*, and this would lead to the undesirable situation where $A$ could be the grandcousin of $B$ whereas $B$ is not the grandcousin of $A$.

We should stress again that there is no unique set of postulated productions to account for the data in Table 1. In fact the productions offered here differ from those in my 1986 paper, for which Bill Anglin (1986) had given a soundness proof.

3. **Basic sentence structure in English**. Our third example deals with a small fragment of English, which we hope contains all basic declarative sentences subject to certain simplifying restrictions: noun phrases are taken to be personal pronouns, adverbs and adverbial phrases are avoided and only the simplest verb phrases are admitted, namely those built from transitive or intransitive verbs. Some such restrictions are necessary to keep this section within reasonable bounds, yet I hope that the sample is large enough to make a convincing case for our way of presenting English grammar.

The method advocated here is of course greatly influenced by the pioneering work of Chomsky. It differs from the latter by sticking to pure production grammars without additional machinery and by placing syntax and morphology on the same footing from the start, as I had advocated for French (1975), Latin (1979) and more recently

for English (1987). I expect to expand the fragment of English treated here in later publications.

Our terminal vocabulary will be a subset of the set of all English words, including inflected forms. The initial vocabulary will consist of a single symbol S. In addition, we shall make use of an auxiliary vocabulary borrowed from traditional grammar, which will be introduced gradually. To start with, we make the following abbreviations:

$$
\begin{aligned}
\text{Subj} &= \text{subject,} \\
\text{Pred} &= \text{predicate,} \\
\text{NP} &= \text{noun phrase,} \\
\text{P}_k &= k\text{-th person } (k = 1, 2, 3), \\
\text{Tens} &= \text{tense,} \\
\text{Inf} &= \text{infinitive,} \\
\text{T}_i &= i\text{-th simple tense } (i = 1, 2), \\
\text{Neg} &= \text{negation,} \\
\text{Asp} &= \text{aspect,} \\
\text{Pass} &= \text{passive of,} \\
\text{VP} &= \text{nuclear verb phrase,} \\
\text{Perf} &= \text{perfect participle of,} \\
\text{Part} &= \text{present participle of,} \\
\text{Obj} &= \text{object,} \\
\text{Acc} &= \text{accusative of.}
\end{aligned}
$$

We begin with a list of contextfree productions, which determine the structure of basic declarative sentences.

(1.1)  S $\longrightarrow$ Subj Pred;

(1.2)  Subj $\longrightarrow$ NP$_k$ P$_k$ ($k = 1, 2, 3$);

(1.3)  NP$_1$ $\longrightarrow$ *I*;
    NP$_2$ $\longrightarrow$ *we, you, they, ......*;
    NP$_3$ $\longrightarrow$ *he, she, it, one, ......*;

(1.4)  Pred $\longrightarrow$ Tens Inf;

(1.5)  Tens $\longrightarrow$ T$_i$ (Neg) (*V*) ($i = 1, 2$; *V* any modal verb);

(1.6)  Inf $\longrightarrow$ (Asp) (Pass) VP;

(1.7)  (Asp) $\longrightarrow$ (*have* Perf) (*be* Part):

(1.8)  VP $\longrightarrow$ $\begin{cases} V \text{ if } V \text{ is an intransitive verb} \\ \text{or a transitive verb with deletable object;} \\ V \text{ Obj if } V \text{ is a transitive verb;} \end{cases}$

(1.9)  Obj $\longrightarrow$ Acc NP$_k$.

We shall make some comments on these productions.

(1.1) A statement consists of a subject followed by a predicate.

(1.2) The subject consists of a noun phrase followed by a matching person marker: $P_1$, $P_2$, $P_3$ denote the first, second and third person respectively. These markers will act on the verb later. We have taken advantage of the fact that in modern English the three persons of the plural affect the verb in the same way as the second person of the singular. Other European languages still require six persons. In some languages, e.g. Latin, the subject may consist of a person marker alone.

(1.3) There is only one noun phrase which may precede $P_1$. $P_2$ may also be preceded by noun phrases built from plurals and $P_3$ by noun phrases built from count nouns or mass nouns, all with appropriate determiners. All these and some others will not be considered here.

(1.4) The predicate consists of a tense, rather widely conceived, followed by an infinitival verb phrase.

(1.5) The tense must contain $T_1$ = present or $T_2$ = past. It may contain a modal verb: *shall, will, can, may* or *must*, two of which traditionally also express the future tense. It has been found convenient to include an optional negation, as its appropriate place is just after $T_i$. We might just mention that two verbs, *need* and *dare*, become modal verbs when preceded by Neg.

(1.6) The infinitival verb phrase consists of an optional aspect, an optional passive transformation and a nuclear verb phrase.

(1.7) The aspect, if present at all, may be *have* Perf, *be* Part or *have*Perf*be*Part. Here Perf and Part are what I wish to call "inflectors", they act on the verb in a way to be specified later.

(1.8) The nuclear verb phrase consists of an intransitive verb, e.g. *come, go, work, sleep*, ......, or a transitive verb with deletable object, e.g. *eat, call, kill*, .....; it may also consist of any transitive verb followed by an object. The object cannot be deleted after e.g. *like, receive, resemble*, ..... .

(1.9) The object consists of the accusative inflector followed by a noun phrase. We have retained the subscript $k$, although it serves no purpose here.

With the help of the productions listed so far we can, for example, prove the following:

$$S \rightarrow \textit{he } P_3T_1 \textit{ call } \text{Acc } \textit{she,}$$

$$S \rightarrow \textit{he } P_3T_1 \text{ Neg } \textit{go,}$$

$$S \rightarrow \textit{he } P_3T_1 \text{ Neg } \textit{will } \text{Pass } \textit{call } \text{Obj,}$$

$$S \rightarrow \textit{he } P_3T_1 \textit{ have } \text{Perf } \textit{be } \text{Part } \textit{come.}$$

The expressions on the right of the arrows are not yet English sentences; they indicate the structure of certain declarative sentences, which may be computed with the help of further productions, however not contextfree ones.

Before stating these productions, we enlarge our auxiliary vocabulary by including the inflector

$$C_{ik} = \text{conjugation for } i\text{-th tense and } k\text{-th person}$$

and the morphemes

$$+not, +s, +ed, +ing.$$

We also assume that the dictionary lists the past tense $V^p$ and the perfect participle $V^q$ of any so-called irregular verb $V$. For example, $go^p = went, \quad go^q = gone.$

We are now ready to adopt the following new productions.

(2.1)            $P_k T_i \rightarrow C_{ik};$

$$C_{11} V \rightarrow \begin{cases} am & \text{if } V = be, \\ V & \text{otherwise;} \end{cases}$$

$$C_{12} V \rightarrow \begin{cases} are & \text{if } V = be, \\ V & \text{otherwise;} \end{cases}$$

$$C_{13} V \rightarrow \begin{cases} is & \text{if } V = be, \\ has & \text{if } V = have, \\ V & \text{if } V \text{ is a modal verb,} \\ V + s & \text{otherwise;} \end{cases}$$

$$C_{2k} \, be \rightarrow \begin{cases} was & \text{if } k = 1 \text{ or } 3, \\ were & \text{if } k = 2; \end{cases}$$

$$C_{2k} V \rightarrow \begin{cases} V^p & \text{if } V \text{ is irregular other than } be, \\ V + ed & \text{otherwise.} \end{cases}$$

(2.2)            $Neg \, V \rightarrow \begin{cases} V + not & \text{if } V \text{ is an auxiliary verb,} \\ do + not \, V & \text{otherwise.} \end{cases}$

(2.3)            $Acc \, I \rightarrow me,$

$Acc \, he \rightarrow him,$

$Acc \, she \rightarrow her,$

$Acc \, we \rightarrow us,$

$Acc \, they \rightarrow them,$

$Acc \, X \rightarrow X \text{ otherwise.}$

(2.4)            $Perf \, V \rightarrow \begin{cases} V^q & \text{if } V \text{ is an irregular verb,} \\ V + ed & \text{otherwise.} \end{cases}$

(2.5)            $Part \, V \rightarrow V + ing.$

(2.6)  $Pass \, V \, Obj \rightarrow be \, Perf \, V \, \text{-, where } V \text{ is any transitive verb which has a}$
        passive.

Here are some comments.

(2.1) $C_{ik}$ is an inflector which will act on any verb to produce its conjugation matrix, e.g. $C_{ik}go$ yields:

$$\begin{pmatrix} go & go & goes \\ went & went & went \end{pmatrix}.$$

In English, if we ignore the almost obsolete subjunctive, this matrix has $2 \times 3 = 6$ entries, of which usually only three are distinct. (The corresponding matrix in German has $4 \times 6 = 24$ entries, in French $7 \times 6 = 42$, in Latin $3 \times 5 \times 6 = 90$, in Hebrew $7 \times 2 \times 10 = 140$, in Arabic and Sanskrit many more.)

(2.2) Auxiliary verbs are the modal verbs, the verb *be* and the verb *have* when followed by the inflector Perf. (When followed by Obj, *have* counts as a transitive verb.) The negation discussed here is the grammatical negation, which may or may not coincide with the logical negation. For example, *can + not* is both the grammatical and the logical negation of *can*, whereas the grammatical negation *must + not* of *must* differs from its logical negation *need + not*.

(2.3) English has only few surviving accusative cases. In German, the inflector Acc acts visibly not only on pronouns, but also on nouns, adjectives and determiners.

(2.4) We assume the dictionary lists $be^q = been$, $have^q = had$, $go^q = gone$, etc. The inflector Perf will never appear before a modal verb.

(2.5) The inflector Part never appears before a modal verb either.

(2.6) Not all transitive verbs have passives, e.g. *resemble* does not. The reader may ignore the dash –, which denotes what Chomsky calls a "trace"; it plays no rôle in this brief exposition.

With the help of the new productions (2.1) to (2.6), we may almost compute our four sample sentences:

$$S \rightarrow he \ C_{11} \ call \ her \ \rightarrow \ he \ call + s \ her,$$

$$S \rightarrow he \ C_{13} \ do + not \ go \rightarrow he \ do + s + not \ go,$$

$$S \rightarrow he \ C_{13} \ will + not \ be \ \text{Perf} call-$$

$$\rightarrow he \ will + not \ be \ call + ed-,$$

$$S \rightarrow he \ C_{13} \ have \ be^q \ come + ing \rightarrow he \ has \ been \ come + ing.$$

All that remains is to compute the correct form of the morphemes *+not*, *+s*, *+ed*, *+ing*. We do this for the written forms; the spoken forms may be calculated in an analogous manner, but that would require a special notation for English phonemes.

Here are some sample spelling rules for *+not* and *+s*:

(3. 1)
$$can + not \rightarrow cannot, \ can't;$$
$$will + not \rightarrow will \ not, \ won't;$$
$$do + not \rightarrow do \ not, \ don't;$$
$$does + not \rightarrow does \ not, \ doesn't; \ etc.$$

(3. 2)    $X + s \rightarrow \begin{cases} Xes & \text{if } X \text{ ends in } z, s, x, sh \text{ or } ch, \text{ or in } o \text{ after a consonant,} \\ Yies & \text{if } X = Yy \text{ and } Y \text{ ends in a consonant or } qu, \\ Xs & \text{otherwise.} \end{cases}$

We shall skip the rather tedious spelling rules for +*ing* and +*ed*.

Here is a complete derivation of our first sample sentence:

$$S \rightarrow \text{Subj Pred}$$
$$\rightarrow NP_3\ P_3\ \text{Tens Inf}$$
$$\rightarrow \text{he } P_3\ T_1\ VP$$
$$\rightarrow \textit{he } C_{13}\ \textit{call } \text{Obj}$$
$$\rightarrow \textit{he call} + s\ \text{Acc } NP_3$$
$$\rightarrow \textit{he calls } \text{Acc } \textit{she}$$
$$\rightarrow \textit{he calls her}$$

Alternatively, we may respresent this derivation by a parsing diagram:

| | | S | | |
|---|---|---|---|---|
| Subj | | Pred | | |
| $NP_3$ | $P_3$ | Tens | | Inf |
| *he* | | $T_1$ | | VP |
| | $C_{13}$ | | *call* | Obj |
| | | *call* + *s* | Acc | $NP_3$ |
| | | *calls* | | *she* |
| | | | *her* | |

Such diagrams serve the convenience of the grammarian; I do not believe that they have any psychological reality. Anyway, it should be clear that the diagram is not a tree. To obtain a tree one would have to delete everything below the words *he* $P_3$ $T_1$ *call* Acc *she*.

## 4. Further discussion.

It is well-known to mathematicians that all and only recursively enumerable sets of symbols can be generated by production grammars from a single initial symbol. If we believe with Chomsky that the set of grammatical sentences of a natural language such as English is recursively enumerable, it becomes evident that the grammar of such a language can be presented in the form of productions.

What is not quite so obvious is how to write down the precise productions that are needed. This is like finding the postulates from which known theorems, as in geometry, can be deduced, an activity no mathematician need be ashamed of. In a sense, the speaker of a language "knows" such rules, though not consciously, and it may require a persistent Socratic process to uncover them.

Many authors feel uneasy about accepting unrestricted production grammars. For example, Cohn (1975) says: "... further restrictions are needed to impart a typical linguistic flavour. This is achieved by imposing restrictions on the rewriting rules." He then goes on to discuss context-sensitive, contextfree and finite-state productions. However, I hope to have convinced the reader that such restrictions, though interesting

mathematically, are not at all natural in linguistics if one wants to treat morphology with the same ease as syntax. For example, one quarter of the productions used in deriving our sample sentence *he calls her* were not even context-sensitive, inasmuch as the strings on the right side of the arrows were shorter than those on the left. This is surprising, as English comes close to being a so-called "analytic" language in having only very few surviving inflections.

As we have seen, to check the grammaticality of an English sentence is like finding the proof of a theorem. On the other hand, to produce a sentence involves at least two steps: creating its underlying structure and then converting this into a string of English words. While the first step is clearly non-deterministic, the second step is very much like computing a given numerical function for given arguments. Actually, as I have shown elsewhere (1987), there is a simple three-tape machine which carries out both of these steps in producing English sentences and which is also capable of analyzing given strings of English words. Of course, the machine is equivalent to a Turing machine, but it resembles a combination of two pushdown automata.

A few words should be said about ambiguity. For example, in our kinship grammar of Section 2 we could derive W → CSP in two different ways:

$$W \;\to\; CW \;\to\; CWP \;\to\; CSP,$$
$$W \;\to\; WP \;\to\; CWP \;\to\; CSP.$$

It seems reasonable to identify these two derivations. To give a convincing example of when two derivations should not be identified, let us consider a somewhat expanded form of our kinship grammar, which also includes the symbol

$$\Sigma = \text{spouse}.$$

In this grammar there would be two essentially inequivalent proofs that

$$R \to \textit{sister-in-law} \text{ (in the context } \# - \#);$$

one derivation goes via the kinship description F$\Sigma$S, the other via F$S\Sigma$. Or again, consider the English sentence "we shall pay \$1000 for your expenses and your airline ticket". There are two distinct ways of deriving this sentence, as the author has discovered to his cost.

To handle the questions raised by ambiguity or the lack of it one should introduce an equality relation between derivations in a production grammar. One then obtains a kind of two-category, to be precise, a strictly monoidal category (Hotz 1966, Benson 1975).

I have refrained from saying anything about the semantics of natural languages. This is a nontrivial problem, as the reader may gather from the brief discussion of grammatical versus logical negation in Section 3. Nonetheless, it has attracted the attention of a number of mathematicians and logicians. For example, a translation of

English into some form of type theory is implicit in the work of Curry (see e.g. Curry and Feys 1958), an idea which was considerably expanded by Montague (see e.g. Montague 1974). From a categorical point of view, this means that we interpret the given natural language in a Cartesian closed category or even a topos.

Benson (1970) had suggested that the interpretation should be viewed as a strictly monoidal functor from the production grammar, as categorized above, into the category of sets with canonical cartesian products. In place of the category of sets one might take any kind of category with a monoidal structure given by a Cartesian product, e.g. a topos, a Cartesian closed category or just a Cartesian category. One might in fact construct such a category freely from a production grammar by forcing the implicit tensor product given by concatenation to become a Cartesian product. Thus one might make mathematical sense out of the Sapir-Whorf thesis that the world as seen by a linguistic community is constructible from their language.

## REFERENCES

1. W. S. Anglin, *A mathematical analysis of a production grammar*, Theoretical Linguistics **13** (1986), 125–138.

2. Y. Bar-Hillel, M. Perlis and E. Shamir, *On formal properties of simple phrase structure grammars*, Z. Phonetik, Sprachwiss. Kommunikat. **14** (1961), 143–172.

3. D. B. Benson, *Syntax and semantics: a categorical view*, Information and Control **17** (1970), 145–160.

4. ——— *The basic algebraic structures in categories of derivations*, Information and Control **28** (1975), 1–29.

5. M. Bhargava and J. Lambek, *A production grammar for Hindi kinship terminology*, Theoretical Linguistics **10** (1983), 227–245.

6. ———, *A production grammar for Sanskrit kinship terminology*, McGill University Working papers in Cognitive Science **6** (1985).

7. B. Brainerd, *Introduction to the mathematics of language study*, Elsevier, New York 1971.

8. J. R. Buchler and H. A. Selby, *Kinship and social organization*, MacMillan, New York 1968.

9. N. Chomsky, *Syntactic structures*, Mouton, s'Gravenhage 1957.

10. ———, *Lectures on government and binding*, Foris Publications, Dordrecht 1982.

11. P. M. Cohn, *Algebra and language theory*, Bull. London Math. Soc. **7** (1975), 1–29.

12. H. B. Curry and R. Feys, *Combinatory Logic I*, North Holland, Amsterdam 1958.

13. W. H. Goodenough, *Yankee kinship terminology: a problem in componential analysis*. In: E. A. Hammel (ed.), *Formal semantics*, American Anthropologist **67** (5) Part 2 (1965), 259–287.

14. M. Gross, *Mathematical models in Linguistics*, Prentice Hall, Englewood Cliffs N.J. 1972.

15. G. Hotz, *Eindeutigkeit und Mehrdeutigkeit formaler Sprachen*, Elektronische Informationsverarbeitung und Kybernetik **2** (1966), 235–247.

16. S. C. Kleene, *Introduction to metamathematics*, Van Nostrand, New York 1952.

17. J. Lambek, *A mathematician looks at French conjugation*, Theoretical Linguistics **2** (1975), 203–214.

18. ——— *Computability and mathematical linguistics*, Notes for Mathematics 180–328, fourth edition, McGill University, Montreal 1976.

19. ——— *A mathematician looks at Latin conjugation*, Theoretical Linguistics **6** (1979), 221–224.

20. ——— *A production grammar for English kinship terminology*, Theoretical Linguistics, **13** (1986), 19–36.

21. ——— *Production grammars revisited*, Proceedings Sixth Amsterdam Colloquium, Amsterdam 1987, 175–195.

22. ——— and M. Lambek, *The kinship terminology of Malagasy speakers in Mayotte*, Anthropological Linguistics **22** (1981), 154–182.

23. F. G. Lounsbury, *Another view of Trobriand kinship categories*. In: E. H. Hammel (ed.) *Formal semantics*, American Anthropologist **67** (5) Part 2 (1965), 142–185.

24. R. Montague, *Formal philosphy, selected papers*, edited by R. H. Thomason, Yale University Press, New Haven 1974.

25. Webster's New Ideal Dictionary, Merrian Company, Springfield 1973.

*McGill University*