

ANNELY ROTHKEGEL

## IDIOMS IN AUTOMATIC LANGUAGE PROCESSING

The automatic identification of multi-word idioms in texts is part of the project "Automatic Lemmatization",<sup>1</sup> the main problem of which is to recognize lexical units of the stem dictionary and their grammatical characterization in a text which – as a rule – contains inflected forms, ambiguities, and idioms as multi-word expressions.

In this paper some examples shall be given to demonstrate the treatment of the latter ones by means of special dictionary entries and a special algorithm.

The problem is that groups of several words are to be identified as semantic and lexical units. The linguistic point of such a classification is based on sign theory. In this connection the individual word cannot generally be regarded as a linguistic unit.

In addition to the difficulties of identification there are problems of the dictionary. Most of the idioms are ambiguous; they have an idiomatic meaning as well as a literal one. These are expressions like *White House* ("government of the USA"), *green light* ("licence"), *draw red herrings* ("devert"), *kick the bucket* ("die") or in German *rechte Hand* ("assistant"), *blinder passagier* ("deadhead"), *grüne Fingerhaben* ("be a skillful gardener"), *in der Tat* ("indeed"), etc.<sup>2</sup>

It is an important point that there are semantic relations between the idiomatic and the literal meaning.<sup>3</sup> Therefore the use of idioms in sentences is often restricted, e.g. *kick the bucket* must be connected with a subject marked as "person" or at least as "animate", which is not the same for *die*. Relations of this kind – as far as they are known – should be made explicit in the lexicon.

---

<sup>1</sup> This is one of the projects at the Sonderforschungsbereich "Elektronische Sprachforschung", University of Saarbrücken, Germany.

<sup>2</sup> It should be mentioned that the problem of idioms is a general one and not dependent on a specific language. In this paper most of the examples are given in English. This doesn't mean that the system shown in the following is restricted to English idioms.

<sup>3</sup> Cp. R. W. LANGACKER (1968, p. 80) "standardized metaphors".

One way to do this is to integrate idioms into the normal dictionary by referring them to the one-word units which are likewise elements of the idioms. Thus it is possible to compare the semantic features of both and to mark the internal metaphorical relations still being under investigation.

The description of the automatic analysis can be divided into three parts:<sup>4</sup>

1) encoding of special morphologic, syntactic, and semantic markers, which are important for the identification;

2) integration of the multi-word units together with the grammatical information into the stem lexicon;

3) algorithm for the discovery of such units within the sentence analysis.

1. These expressions are divided into classes corresponding to their fixed syntactic patterns. Such a pattern can be described as part of a tree embedded in a phrase structure tree. The dominating node represents the part of speech of the whole expression. The categories verb (*V*), noun (*N*), adverb (*A*) are substituted by *V'*, *N'*, and *A'* each of them dominating an idiomatic subsystem. These systems are generated by small phrase structure grammars: for each pattern several productions are fixed in form and order to represent the whole derivation of the idiom.

Further we have to consider the fact, that the same structural components and their combinations can be found in various syntactic patterns which are dominated by different nodes, e.g. coordination connected with preposition:

<i>beer and skittles</i>	dominated by <i>N'</i>
<i>by leaps and bounds</i>	» by <i>A'</i>
<i>rain cats and dogs</i>	» by <i>V'</i>

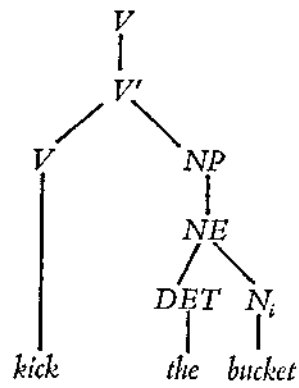
With regard to this it is necessary to use the auxiliary symbols needed to produce different patterns the structure of which is partly the same.

---

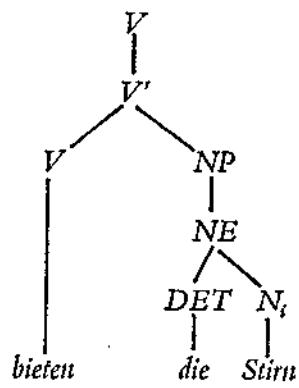
<sup>4</sup> An exhaustive description of German idioms is given in A. ROTHKEGEL (1973).

Examples:

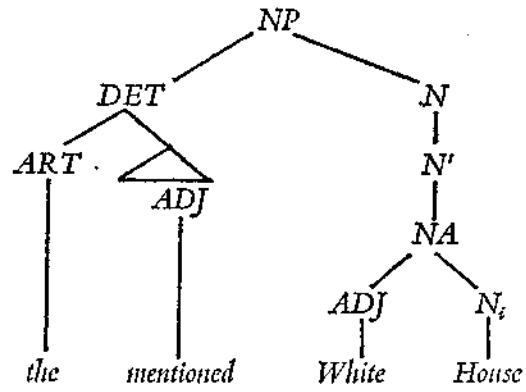
a) *kick the bucket*



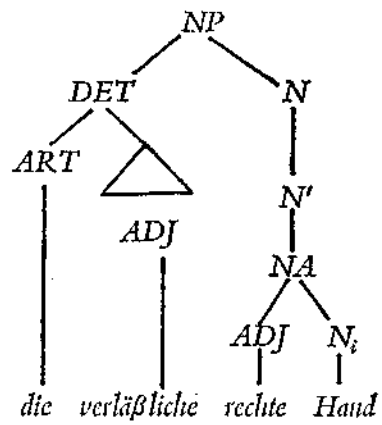
a') *die Stirn bieten*



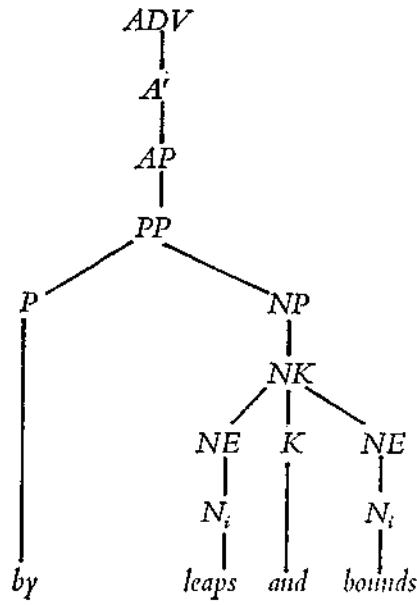
b) *White House in the mentioned White House*



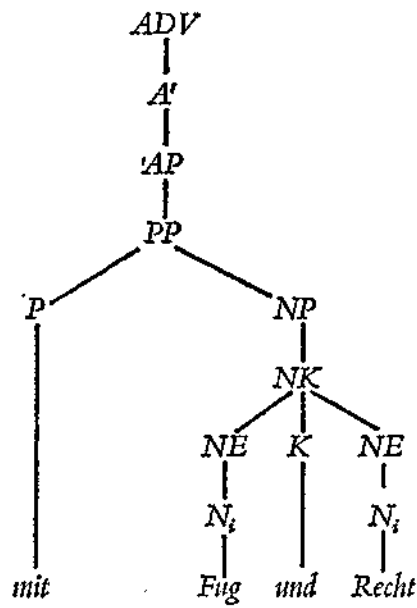
b') *rechte Hand in die verlässliche rechte Hand*



c) *by leaps and bounds*



c') *mit Fug und Recht*



The class an idiom belongs to is symbolized by an eight-figure number. It is formed systematically, so that the productions necessary for the syntactic structure can be derived from it. The number of the first place (on the left) determines the part of speech of the whole expression, the number of the last place refers to the part of speech of the nominal basic lexeme like noun, (noun in idiom), adjective, adverb, verb. The other numbers characterize the syntactical surface structure of the idiom (cp. the lists of the basic rules). By such a procedure new classes can be added without changing anything in the system.

## List of basic rules:

Number of the rule	rule	marked in class number as
(1)	$Z \rightarrow N_i$	1 (8. place)
(2)	$Z \rightarrow ADJ$	2 »
(3)	$Z \rightarrow ADV$	3 »
(100)	$N' \rightarrow NP$	4 (1. place)
(200)	$A' \rightarrow AP$	2 »
(210)	$AP \rightarrow NP$	if 1 → 9 at 2. place
(220)	$AP \rightarrow PP$	» A or B »
(230)	$AP \rightarrow PP PP$	» C »
(240)	$AP \rightarrow PP P$	» D »
(400)	$V' \rightarrow VP$	4 (1. place)
(410)	$VP \rightarrow V NP$	if 1 → 9 at 2. place
(420)	$VP \rightarrow V PP$	» A or B »
(430)	$VP \rightarrow V PP PP$	» C »
(10)	$NP \rightarrow NE$	1 (2. place)
(20)	$NP \rightarrow NA$	2 »
(30)	$NP \rightarrow NK$	3 »
(40)	$NP \rightarrow NE NE$	4 »
(50)	$NP \rightarrow NE NA$	5 »
(60)	$NP \rightarrow NE PP$	6 »
(60), (99)	$NP \rightarrow NE PP, P \rightarrow P_{em}$	7 »
(70)	$NP \rightarrow NA PP$	8 »
(70), (99)	$NP \rightarrow NA PP, P \rightarrow P_{em}$	9 »
(80)	$PP \rightarrow P NP$	A »
(80), (99)	$PP \rightarrow P NP, P \rightarrow P_{em}$	B »
(90)	$PP \rightarrow P NP P$	D »
(95)	$PP \rightarrow PP PP$	C »
(11)	$NE \rightarrow Z$	0 (from 3. place)
(12)	$NE \rightarrow DET Z$	1 »
(13)	$NE \rightarrow (DET) Z$	2 »

(21)	$NA \rightarrow ADJ Z$	0	»
(22)	$NA \rightarrow DET ADJ Z$	1	»
(23)	$NA \rightarrow (DET) ADJ Z$	2	»
(24)	$NA \rightarrow (ADJ) Z$	3	»
(25)	$NA \rightarrow DET (ADJ) Z$	4	»
(26)	$NA \rightarrow Z ADJ$	5	»
(31)	$NK \rightarrow NE K NE$	0	»
(32)	$NK \rightarrow NA K NA$	1	»
(33)	$NK \rightarrow NE K NA$	2	»
(34)	$NK \rightarrow NA K NE$	3	»
(41)	$PP PP \rightarrow P NP P NP$	0	»
(42)	$PP PP \rightarrow P_{em} NP P_{em} NP$	1	»
(43)	$PP PP \rightarrow P_{em} NP P NP$	2	»
(44)	$PP PP \rightarrow P NP P_{em} NP$	3	»
(99)	$P \rightarrow P_{em}$		

Example: *by leaps and bounds* or *mit Fug und Recht*

class number: 2A3000-1

1. place: 2 rule (200) :  $A' \rightarrow AP$
2. » : A dependent on the 1. place additional rule (220) :  $AP \rightarrow PP$   
for A (80) :  $PP \rightarrow P NP$
3. » : 3 rule (30) :  $NP \rightarrow NK$
4. » : 0 rule (31) :  $NK \rightarrow NE_1 K NE_2$
5. » : 0 rule (11) :  $NE_1 \rightarrow Z$
6. » : 0 rule (11) :  $NE_2 \rightarrow Z$
8. » : 1 rule (1) :  $Z \rightarrow N_i$

A generation of the productions is important for the analysis, because  
 1) the inventory and its order must be compared with the text and  
 2) the internal syntactic structure of the idiom can be represented as part of the sentence structure.

In addition to the productions there are morphological and syntactical restrictions. In a lot of cases they help to distinguish the idiomatic meaning from the literal one:

1) no discontinuity is permitted: *White House*, not *the white big house* or *the house is white*; *schwarzer Peter*, not *der schwarze kleine Peter* or *Peter ist schwarz*;

2) attribution is allowed: *the (mentioned) White House*, *eine (grosse) rolle Spielen*;

3) attribution is not allowed: *kick the bucket*, *die Stirn bieten*;

4) the noun must be in singular: *White House*, *rule the roost*, *kalte Küche*, *Fuss fassen*;

- 5) the noun must be in plural: *get cold feet, burn one's fingers kalte Füße bekommen, Wurzeln schlagen*;  
 6) noun with article: *have the face, die Stirn bieten*;  
 7) noun without article: *from pillar to post, play second fiddle, Fuss fassen, Hand anlegen*;  
 8) only in negation: *cut no ice, nicht von Pappe sein*;  
 9) no passive: *kick the bucket, den Kopf verlieren*.

Generally the encoding refers to the whole expression. The list of all marked idioms is the basis to get the lexical entries automatically by attaching the grammatical information to the corresponding basic lexeme of the idioms described in the following chapter.

2. The simplest way to connect the expression and the grammatical information would be to build up a separate lexicon of idioms as is practised in general. In this case the dictionary only has the function of a simple word-store and data carrier for the algorithm. From a linguistic view this is not satisfactory for reasons I mentioned above. Beyond that there is the fact that idioms belong to the field of word-formation. The basic lexemes of the idiom are already elements of the vocabulary represented in the dictionary. Therefore it is adequate to use the one-word entries of the normal dictionary as addresses of idiom entries.

With regard to the lexicon we distinguish two types of idioms:

i) the meaning of only one lexeme is determined by the partner lexeme as in *blinder Passagier, four-letter word*. Only *blind* respectively *four-letter* has an idiomatic meaning.

ii) The meaning of the whole expression strictly depends on the fixed connection of several lexemes. It is not possible to divide the meaning corresponding to the single constituents, e.g. *White House, kick the bucket, bring home the bacon, Staub aufwirbeln, den Mond anbellern*.

Concerning i there are no problems in the dictionary: the single lexemes (only basic lexemes) are marked as parts of an idiom by *I*, representing the class number, and by *R*, representing the information about the grammatical restrictions. The lexeme with the idiomatic meaning is characterized additionally by *S*, the component of the special semantic features. The corresponding lexical entries of the examples in i are as follows:

<i>blind</i>	: <i>I, R, S</i>	<i>four-letter</i>	: <i>I, R, S</i>
<i>Passagier</i>	: <i>I, R</i>	<i>word</i>	: <i>I, R</i>



With regard to  $\Pi$  there is no semantic congruence between the constituents of the idiom. Instead, the partner lexeme itself has the function of a marker. The component  $I_u$  represents the class number of this kind of connection,  $R$  and  $S$  referring to the whole expression are added only to one of the constituents. The following entries concern the idioms *White House*, *rechte Hand*, *kick the bucket*, *in die Hand nehmen*.

<i>White</i> : $I_{11}$	<i>House</i> : $I_{11}, R, S, \textit{white}$
<i>recht</i> : $I_{11}$	<i>Hand</i> : $I_{11}, R, S, \textit{recht}$
<i>bucket</i> : $I_{11}$	<i>kick</i> : $I_{11}, R, S, \textit{bucket}$
<i>Hand</i> : $I_{11}$	<i>nehmen</i> : $I_{11}, R, S, \textit{Hand}$

The component  $S$  also contains the information as to whether the idiom as a whole is ambiguous, i.e. whether there could be a literal meaning in the sentence. This point is important for the result of idiom identification.

3. The identification of idioms is relevant in that part of the sentence analysis in which the parts of speech are determined. The ambiguity of the idioms can only be solved on a higher level of the syntactic analysis. But here it is possible to decide whether there is an idiom in a sentence or not.

During the dictionary look-up the idiom marker  $I$  gets importance. As only basic lexemes are marked and not function words (e.g. only *Tat* of *in der Tat*, only *pink* of *in the pink*) there can be an idiom in a sentence even if only one  $I$ -marked word exists. On the other hand there can be several  $I$ -marked words which do not necessarily belong together and form an idiom.

The program of identification is divided into the following steps:

a) check whether there are the same or different class numbers of the  $I$ -marked words (elements of one idiom have the same class number.). Form groups of  $I$ -marked words which have the same class number.

b) Compare whether the quantity of basic lexemes of every group corresponds to the requirement of the class.

c) Generate the sequence of the parts of speech dependent on the class number and compare with the sequence of  $I$ -marked words and their context in the text.

d) Check the grammatical restrictions of the potential idiom and compare with the text.

e) Compare the partner lexemes of the text with the dates of the lexicon.

If one of the conditions from *b* to *e* is not fulfilled, the analysis stops with a negative result.

The following examples show the complexity of the analysis, when there are parts of idioms which potentially belong to different expressions. First a list of such idioms is given. It also contains a paraphrase of their idiomatic meaning, the class number, the quantity of required basic lexemes corresponding to this class, and the sequence of the parts of speech, which can be generated by means of the class number.

Examples:

idiom	class number	quantity of basic lexemes	sequence of parts of speech
<i>green light</i> (traffic light, licence)	120 ---- 1	2	ADJ, N <sub>i</sub>
<i>red light</i> (traffic light, danger)	120 ---- 1	2	•
<i>red tape</i> (bureaucratism)	120 ---- 1	2	•
<i>a red rag to a bull</i> (object which enrages s.o.)	18A211 - 1	3	DET, ADJ, N <sub>i</sub> , P, DET, N <sub>i</sub>
<i>in the light (of)</i> (in the view of)	2A11 --- 1	1	P, DET, N <sub>i</sub>
<i>on the rocks</i> (in financial difficulties)	2A11 --- 1	1	•
<i>in the red</i> (in debt)	2A11 --- 1	1	•
<i>see red</i> (become enraged)	410 ---- 2	2	V, ADJ
<i>see the light</i> (be born)	411 ---- 1	2	V, DET, N <sub>i</sub>
<i>see the red light</i> (see danger ahead)	421 ---- 1	3	V, DET, ADJ, N <sub>i</sub>
<i>draw red herrings</i> (devert)	420 ---- 1	3	V, ADJ, N <sub>i</sub>

The corresponding lexicon entries of the basic lexemes are as follows. Here only the idiom information is listed which is to be added to the usual grammatical markers.

<i>bull</i>	<i>I(18A211-1)</i>
<i>draw</i>	<i>I(420 ----1), R,S, red, herring</i>
<i>green</i>	<i>I(120 ----1)</i>
<i>herring</i>	<i>I(420 ----1)</i>
<i>light</i>	1. <i>I(120 ----1),R,S<sub>1</sub>, green</i> 2. <i>I(120 ----1),R,S<sub>2</sub>, green</i> 3. <i>I(120 ----1),R,S<sub>1</sub>, red</i> 4. <i>I(120 ----1),R,S<sub>2</sub>, red</i> 5. <i>I(2A11 ---1) R, S,</i> 6. <i>I(411 ----1)</i> 7. <i>I(421 ----1)</i>
<i>rag</i>	<i>I(18A211-1), R,S, red, bull</i>
<i>red</i>	1. <i>I(120 ----1)</i> 2. <i>I(18A211-1)</i> 3. <i>I(410 ----2)</i> 4. <i>I(421 ----1)</i> 5. <i>I(420 ----1)</i> 6. <i>I(2A11 ---1)</i>
<i>rock</i>	<i>I(2A11 ---1)</i>
<i>see</i>	1. <i>I(410 ----2),R,S, red</i> 2. <i>I(411 ----1),R,S, light</i> 3. <i>I(421 ----1),R,S, red, light</i>
<i>tape</i>	<i>I(120 ----1),R,S, red</i>

Analysis of idioms in sentences:

(1) *He cannot understand that the company is on the rocks.*

<i>I</i> -marked lexeme	class number
<i>rock</i>	2A11 --- 1

The program then follows the several steps of the identification scheme given above:

- a) There is only one group of *I*-marked words.
- b) Only one basic lexeme is required. There is agreement with the text.
- c) As sequence of parts of speech *P DET N* is generated, which corresponds to the context of *rocks* in the sentence.
- d) Grammatical restrictions: no discontinuity, definite article, noun in plural.

These conditions are fulfilled.

- e) As the class 2A11 --- 1 requires only one basic lexeme, a partner lexeme can not be compared.

Result: *on the rocks* can be an idiom. It is ambiguous.<sup>5</sup>

- (2) *But I saw the red light some time before.*

<i>I</i> -marked lexemes	class number	
<i>see</i>	410 ---- 2	
	411 ---- 1	
	421 ---- 1	
<i>red</i>	120 ---- 1	421 ---- 1
	18A211 - 1	420 ---- 1
	410 ---- 2	2A11 --- 1
<i>light</i>	120 ---- 1	411 ---- 1
	2A11 --- 1	421 ---- 1

Scheme of analysis:

- a) a group must have elements of the same class number or one element with a single class number. The following groups are formed:

		required lexemes
1) <i>see</i>	410 ---- 2	2
<i>red</i>	410 ---- 2	
2) <i>see</i>	411 ---- 1	2
<i>light</i>	411 ---- 1	
3) <i>see</i>	421 ---- 1	3

---

<sup>5</sup> This information is given by the *S*-component of the lexicon entry.

	<i>red</i>	421 ---- 1	
	<i>light</i>	421 ---- 1	
4)	<i>red</i>	120 ---- 1	2
	<i>light</i>	120 ---- 1	
5)	<i>red</i>	18A211 - 1	3
6)	<i>red</i>	420 ---- 1	3
7)	<i>red</i>	2A11 --- 1	1
	<i>light</i>	2A11 --- 1	

b) After comparing the required lexemes with the given ones groups 5 and 6 are omitted. Group 7 must be substituted by two different groups, because the class 2A11 --- 1 requires only one basic lexeme.

8)	<i>red</i>	2A11 --- 1	1
9)	<i>light</i>	2A11 --- 1	1

The idiom identification program continues with a group containing the most elements. In this case it is group 3.

c) Sequence of part of speech corresponding to class 421 ---- 1: *V DET ADJ N*. It is the same sequence as in the text. The text also corresponds to the

d) grammatical restrictions as "noun in singular", "definite article", "no further attribution".

e) The partner lexemes of *see* are *red* and *light*. There also is an agreement with the text.

Result: *see the red light* is a potential idiom. It is ambiguous.

In this case it is not necessary to check the other groups.

(3) *She saw a light on the green rocks.*

I-marked words	class number	
<i>see</i>	410 ---- 2	
	411 ---- 1	
	421 ---- 1	
<i>light</i>	120 ---- 1	411 ---- 1
	2A11 --- 1	421 ---- 1
<i>green</i>	120 ---- 1	
<i>rock</i>	2A11 --- 1	

Scheme of analysis:

a)	groups	class number	required lexemes
1)	<i>see</i>	410 ---- 1	2
2)	<i>see</i>	411 ---- 1	2
	<i>light</i>	411 ---- 1	
3)	<i>see</i>	421 ---- 1	3
	<i>light</i>	421 ---- 1	
4)	<i>light</i>	120 ---- 1	2
	<i>green</i>	120 ---- 1	
5)	<i>light</i>	2A11 --- 1	1
	<i>rock</i>	2A11 --- 1	

b) Group 1 and 3 are omitted because there is no agreement of the quantity of basic lexemes. Group 5 is to be substituted by two groups:

6)	<i>light</i>	2A11 --- 1	1
7)	<i>rock</i>	2A11 --- 1	1

c) The program continues with one of the biggest groups. The first is group 2. Sequence of parts of speech corresponding to class 411 ---- 1: *V DET N*. This is the same as in the text.

d) The grammatical restrictions are violated: the idiom must contain the definite article, but the sentence contains the indefinite article.

Result: group 2 does not represent any idiom. The next is group 4; the program restarts with

c) Sequence of parts of speech corresponding to class 120 ---- 1: *ADJ N*. The immediate context of *light* is not an adjective. Therefore the result is negative: group 4 also does not represent any idiom.

The next is group 6; restart at

c) Sequence of parts of speech for 2A11 --- 1; *P DET N*. There is no agreement with the sequence in the text which is *V DET N* (*saw a light*).

The last one is group 7:

c) Sequence of parts of speech: *P DET N*. The context of *rocks*

differs from the required sequence of the parts of speech. It is *DET ADJ N* (*the green rocks*).

Result: this sentence does not contain any idiom.

The next two examples show that a syntagm can be an idiom or a part of another syntagm.

- (4) *After having bought a car he is in the red*  
 (5) *She lives in the red house.*

In this case the part of speech of the basic lexeme given by the class number is relevant. If it is an adjective, the following context must be checked.

Finally, besides the practical purpose which is evident in computational language analysis automatic idiom identification is characterized by the general question, how to integrate linguistic aspects into computational linguistics. I think, there is an opportunity to get automatically new bases for language research. With regard to the lexicon there is another point: a computer lexicon does not need to be readable by a human user. Therefore it is possible to make explicit principles of productivity like word-formation and metaphor which would be a more adequate description of the vocabulary.

Symbols:

<i>A'</i>	subsystem of adverb
<i>ADJ</i>	adjective (part of speech)
<i>ADV</i>	adverb (part of speech)
<i>AP</i>	adverbial phrase
<i>DET</i>	determiner ( <i>ART</i> article)
<i>N'</i>	subsystem of noun
<i>N<sub>i</sub></i>	noun in idiom (part of speech)
<i>NA</i>	nominal phrase attributed by adjective
<i>NE</i>	simple nominal phrase
<i>NK</i>	nominal phrase with coordination
<i>NP</i>	nominal phrase (not specified)
<i>P</i>	preposition (part of speech)
<i>P<sub>en</sub></i>	preposition with article (e.g. <i>zum, im</i> )
<i>PP</i>	prepositional phrase
<i>V'</i>	subsystem of verb
<i>V</i>	verb (part of speech)
<i>VP</i>	verbal phrase
<i>Z</i>	dummy symbol for noun, adjective, adverb

## REFERENCES

- E. AGRICOLA (ed.), *Wörter und Wendungen. Wörterbuch zum deutschen Sprachgebrauch*, Leipzig 1963<sup>2</sup>.
- J. BAR-HILLEL, *Idioms*, in W. N. LOCKE, A. D. BOOTH (eds.), *Machine translation of languages*, New York 1957<sup>2</sup> (first ed. 1955), pp. 183-193.
- R. PH. BOTHA, *The function of the lexicon in transformational generative grammar*, The Hague 1968.
- W. L. CHAFE, *Idiomaticity as an anomaly in the Chomsky paradigm*, in «Foundations of Language», IV (1968) pp. 109-127.
- W. FLEISCHER, *Wortbildung der deutschen Gegenwartssprache*, Leipzig-Tübingen 1971<sup>2</sup> (first ed. 1969).
- B. FRASER, *Idioms within a transformational grammar*, in «Foundations of language», VI (1970), pp. 22-42.
- W. FRIEDERICH, *Moderne deutsche Idiomatik*, München 1966.
- J. J. KATZ, P. M. POSTAL, *Semantic interpretation of idioms and sentences containing them*, Quarterly Progress Report 70, Research Laboratory of Electronics, M.I.T., Cambridge (Mass.) 1963, pp. 275-282.
- R. KLAPPENBACH, *Feste Verbindungen in der deutschen Gegenwartssprache*, PBB(H) 82, 1961, pp. 443-457.
- R. W. LANGACKER, *Language and its structure: Some fundamental linguistic concepts*, New York 1968.
- I. A. MEL'ČUK, *O terminax 'ustojčivost' i 'idiomaticnost'*, in «Voprosy Jazykoznanija», IV (1960), pp. 73-80.
- A. ROTHKEGEL, *Feste Syntagmen. Grundlagen, Strukturbeschreibung und automatische Analyse*, Tübingen 1973.
- U. WEINREICH, *Problems in the analysis of idioms*, in J. PUHVEL, *Substance and structure of language* (ed.), Berkeley-Los Angeles 1969, pp. 23-81.
- H. WISSEMAN, *Das Wortgruppenlexem und seine lexikographische Erfassung*, in *Indogermanische Forschungen*, LXVI (1961), pp. 225-258.
- F. T. WOOD, *English Colloquial Idioms*, London 1969.