

# License to COLL

## How to bind bound words and readings to their contexts\*

Jan-Philipp Soehn · Dept. of German Linguistics, University of Jena · Email: jp.soehn@uni-jena.de

### 1 Motivation

Idioms are omnipresent in everyday language. Nonetheless, they have been widely neglected by linguists developing grammar fragments. And even where an account for idioms has been given, most approaches have their shortcomings (cf. Riehemann, 2001, ch. 4).

In this contribution we want to focus on decomposable and non-decomposable idioms<sup>1</sup> and idioms containing bound words. We concentrate on technical aspects of the analysis and refrain from presenting detailed linguistic corpus data due to space limitations. By “idiom” we mean idiomatic expressions that do not form complete sentences as would be the case for e. g. *His bark is worse than his bite*.

(1) *make waves* (“cause trouble”)

(2) *spill the beans* (“divulge a secret”)

The expressions in (1) and (2) are instances of decomposable idioms, i. e. their meaning can be derived from the idiom parts. Note that idiom parts are not necessarily to be understood literally. In (1), e. g., we can attribute the meaning “cause” to *make* and “trouble” to *waves*. The idiomatic meaning of the whole idiom consists of the idiomatic meanings of its parts.

Where this is not the case, an idiom is non-decomposable: the meaning of the whole phrase has nothing to do with the meaning of the words the idiom consists of. Consider (3) and (4):

(3) *saw logs* (“snore”)

(4) *shoot the breeze* (“chat”)

It is not clear how to assign the meaning “snore” to the words *saw* and *logs*, the same holds for “chat”.

Finally, we want to draw the attention to idioms comprising bound words or “cranberry words” (Aronoff, 1976). These are expressions which are highly collocationally restricted. Dobrovolskij (1988) compiled quite a lot of examples for German, Dutch and English.

---

\*The research to the paper was funded by the *Deutsche Forschungsgemeinschaft*. I am grateful to Stefan Müller, Christine Römer, Manfred Sailer, Adrian Simpson and the reviewers of HPSG'04 for comments and Michelle Wibraham for help with English.

<sup>1</sup>Cf. Nunberg et al. (1994), e. g., for this distinction.

(5) *to learn/do sth. by rote* (automatically, by heart)

(6) *to cock a snook* (to thumb the nose)

The underlined words are restricted to the given contexts. Sometimes there is some variation, as in *to lie/go/lay doggo* (Brit. slang; “to hide oneself”), but a free distribution is not possible. Such idioms can be either decomposable or non-decomposable.

### 2 Lexemes and Listemes

Before we present our analysis, we point out a way that enables us to select a specific word. This forms a prerequisite of our approach.

Idioms often consist of particular words which cannot be substituted by semantically equivalent terms. It seems in general that each word has a unique “identity” with an idiosyncratic behavior. The possibility to select a particular word would, thus, be a useful feature. Up to now, there has been a discussion about the necessity of having such kind of selection. One could argue that any data in question are to be handled as Constructions or collocations. But why impose such a “heavy thing” on an expression like *to furrow one's brow*? Would it not be plausible that the verb *furrow* simply selects a word of the form *brow*? For perfect tense in German a main verb has to be combined with the right auxiliary (*haben/sein*; in HPSG with the attribute AUXF, cf. Heinz and Matiasek, 1994, p. 222). Here one does nothing other than to select a particular lexeme.

Krenn and Erbach (1994) made an important contribution to idiom analysis within the HPSG framework. They suggested selecting particular lexemes via their feature LEXEME below CONTENT INDEX. This idea of having lexeme information in the CONTENT is questionable. A lexeme combines phonetic, morphological, syntactic and semantic properties all together, not only semantic information. Besides, their approach had several technical shortcomings (cf. Soehn and Sailer, 2003). We therefore propose that the LEXEME approach has to be discarded.

A different concept that helps to distinguish between individual words is that of a listeme<sup>2</sup>. As the

---

<sup>2</sup>This term has been introduced by Di Sciullo and Williams (1988) for a sign that is listed in the lexicon.

concept holds the characteristic of listedness in a lexicon, we use it in our grammar to identify a particular word or phrase. Thus, we insert LISTEME into the feature geometry below CATEGORY, emphasizing the morpho-syntactic character of information. More precisely, we put it below HEAD. This has two consequences: firstly, it is available for selection, as a HEAD value is below SYNSEM. Secondly, the LISTEME value of a projection is the same as the one of the head, as all HEAD features “percolate” according to the HEAD-FEATURE-PRINCIPLE. For our *furrow*-example that means that a modified direct object *his heavy brow* still has the same LISTEME value as *brow* alone.

A third question to address is the handling of pronominalization. It is necessary that pronouns have the same LISTEME value as their antecedent.<sup>3</sup> In Krenn and Erbach’s approach this was the major motivation of putting the LEXEME feature in the INDEX. To emulate this quality, we propose a constraint ensuring that each pronoun which is co-indexed with an antecedent takes over its LISTEME value. In the lexical entries of pronouns that value would be left underspecified in that way, that it consists of a disjunction of an identifying value (*she, her*, etc.) and a wildcard. In case of co-indexation the wildcard is identical to the LISTEME value of the antecedent and – by virtue of the constraint – becomes the actual and concrete LISTEME value of the pronoun. An informal description of such a pronoun constraint is illustrated in (7).

(7) PRONOUN-LISTEME-CONSTRAINT:

If a pronoun is co-indexed with an antecedent, it takes over the LISTEME value of that antecedent. Otherwise the LISTEME value of this pronoun is that of the other disjunct.

The value of LISTEME is an atomic sort as *brow, heavy, furrow, take, she* etc. In order to identify listemes for the same words having different meanings, we use numeric indices just as in a dictionary.

In summary, discarding the LEXEME approach, we propose a more adequate solution for the problem of selecting particular words, at least with respect to terminology, technical feasibility and the feature geometry. We introduce a feature LISTEME which is appropriate for the sort *head* taking atomic sorts as its value.

### 3 Licensing Contexts

Getting to the analysis, we have to define a second attribute in the feature geometry. We declare objects of sort *sign* to bear a list-valued feature COLL (Context Of Lexical Licensing), first introduced by Richter and Sailer (1999). The COLL list may contain objects of sort *barrier*. These *barriers* are particular nodes in the

<sup>3</sup>E. g. in the phrase *He furrowed it*. the pronoun has the same LISTEME value as its antecedent, satisfying the subcategorizational requirement of the verb.

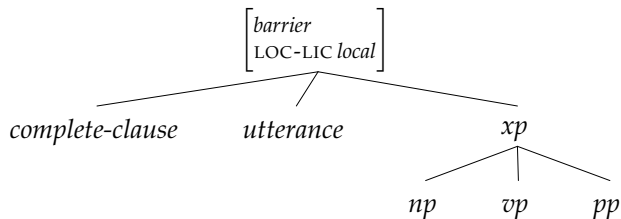


Figure 1: Sort hierarchy for *barrier*

syntactic configuration, like XPs, complete clauses or utterances (a complete clause with an illocutionary force). The concept of barriers is borrowed from the tradition of generative grammar, where these form boundaries for government and binding principles. We avail ourselves of this concept and use similar barriers for the restriction of distributional phenomena.

*barrier* objects have an attribute LOCAL-LICENSER (LOC-LIC) which has a value of sort *local*. In the lexical entry of an idiomatic word one can thus specify a *barrier* on its COLL list with a specific *local* configuration. Subsorts of *barrier* are illustrated in figure 1. The subsorts of *barrier* correspond to nodes in the syntactic tree with particular properties. The following relations identify the nodes which relate to the barriers *complete-clause* and *vp*, respectively.<sup>4</sup>

$$\forall \square \left( \left( \text{is\_complete-clause}(\square) \leftrightarrow \left[ \begin{array}{l} \text{phrase} \\ \text{SS} \left[ \begin{array}{l} \text{STATUS complete} \\ \text{LOC CAT} \left[ \begin{array}{l} \text{HEAD verb} \\ \text{SUBCAT elist} \end{array} \right] \end{array} \right] \end{array} \right] \right) \right) \right)$$

$$\forall \square \left( \left( \text{is\_vp}(\square) \leftrightarrow \left[ \begin{array}{l} \text{phrase} \\ \text{SS} \left[ \begin{array}{l} \text{STATUS incomplete} \\ \text{LOC CAT} \left[ \begin{array}{l} \text{HEAD verb} \\ \text{SUBCAT netist} \end{array} \right] \end{array} \right] \end{array} \right] \right) \right) \right)$$

The LICENSING-PRINCIPLE (informally in 8) makes sure that if there is a barrier specified on a word’s COLL list, there is an actual barrier in the phrase our word occurs in. This barrier must fulfill the *local* requirements and it has to be minimal, i. e., there is no other potential barrier of the same kind between the word and the actual barrier.

(8) LICENSING-PRINCIPLE (LIP):

For each *barrier* object on the COLL list of a sign *x* and for each phrase *z*:

- the LOCAL value of *z* is identical with the LOC-LIC value,
- iff *z* dominates *x*, *z* can be identified as the barrier specified<sup>5</sup> and *z* dominates no sign *y* which in turn dominates *x* and forms an equivalent barrier.

Hence, a word for which a barrier is defined cannot occur elsewhere; its distribution is already specified in the lexical entry.

This concludes the description of technical requirements for our approach to idioms. Note that we have

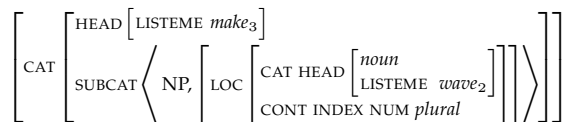
<sup>4</sup>Cf. (Richter, 1997, pp. 68f) for the STATUS feature.

defined a very small number of new sorts and attributes to be included in the signature. All idiosyncratic information comes from the lexicon, as we will see in the next section.

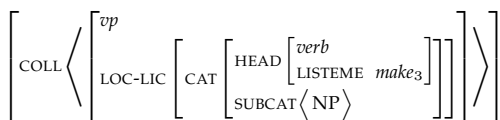
### 3.1 Decomposable Idioms

Let us show how a decomposable idiom can be analysed with our proposal. Take for instance the idiom in (1) *make waves*<sup>6</sup>. We can assign the meanings “cause” and “trouble” to *make* and *waves* and assume that there are two lexical entries for the idiomatic usage of these words.<sup>7</sup>

The idiomatic *make* subcategorizes for a plural noun with the word form *wave* (the idiomatic version) creating a VP with the meaning “cause trouble”.



*wave*<sub>2</sub> for its part bears a non-empty COLL list which looks as follows:



The distribution of the idiomatic noun *waves* is restricted in that it must be the complement of idiomatic *make*. The LIP makes sure that the specified *vp* on the COLL list is identical to the actual VP containing *make* and *waves*. That would have the following semantics:  $\lambda x. [\text{waves}''(y)](\text{make}''(x, y))$ <sup>8</sup>. Defining the barrier as a VP correctly implies that passivization of this idiom is not possible.<sup>9</sup>

Our example *spill the beans*<sup>10</sup> can be analysed analogously. As we assume regular syntactic composition to be in force, we predict that different specifiers (*some beans*) or modifications (as *some very compromising beans*) are grammatical.

A special case of the idiom not occurring in its canonical form is that of pronominal reference. In fact, pronominalization is quite hard to handle in idiom analysis. Cf. the following example:

<sup>6</sup>as in “Italian film makes waves” from <http://news.bbc.co.uk/1/hi/entertainment/film/3171907.stm> (All weblinks were found by Google on 01-27-2004)

<sup>7</sup>Another meaning of the idiom is “call attention” or “attract interest”.

<sup>8</sup>In this contribution we have not enough space to go into details of semantics. Under CONTENT LF we give the logical form of the expression, using a double apostrophe to indicate an idiomatic meaning.

<sup>9</sup>Riehemann found 5 examples out of 243 (2%) where the idiom parts do not occur within the same VP. If one wants to account for those (including passivization and a relative clause) the barrier is simply to be set accordingly.

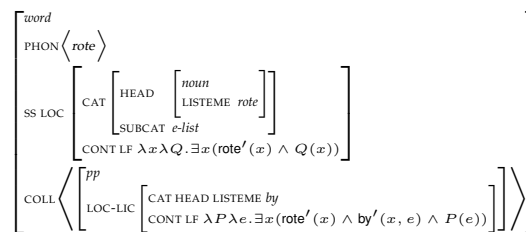
<sup>10</sup>as in “Tom Cruise has spilled the beans on Nicole Kidman’s relationship with US musician Lenny Kravitz.” from <http://www.smh.com.au/articles/2003/11/29/1070081589377.html?from=storyrhs>

(9) *Eventually she spilled all the beans. But it took her a few days to spill them all.*<sup>11</sup>

Here the pronoun *them* refers back to the idiomatic *beans*. As described in section 2 a pronoun has the same LISTEME value as its antecedent, so *them* gets its correct meaning. This being the case, the subcategorization requirements of idiomatic *spill* in both clauses are satisfied. The antecedent of *them* in turn is licensed by its own COLL value stating that the idiomatic *beans* can only occur together with the verb *spill* in its idiomatic use. The barrier is a *complete-clause* which allows e. g. passive or relative constructions. Thus, our proposal can handle pronominalization data, too.

### 3.2 Bound Words

With our approach we can handle bound words as well. The idiom in (5) *to learn sth. by rote*<sup>12</sup> contains a word that never occurs in other contexts than as a complement of a by-PP. The idiom is decomposable: *rote* means something like “routine”. The relevant parts of its lexical entry can thus be stated as follows:



By defining the CONTENT value of the barrier *pp* we prevent a modification of *rote*, which would be ungrammatical. The PP can modify any verb, allowing the occurrence of (*know, learn, sing, do,...*) *sth. by rote*.

To account for the example in (6), the lexical entry of *snook* requires a *vp* barrier with an appropriate LISTEME value of the head, as seen for the idiom *make waves*. We can restrict the distribution of these bound words in the same way as we handle idiomatic words contained in a decomposable idiom.

### 3.3 Non-decomposable Idioms

For idioms that have a non-decomposable meaning we define phrasal lexical entries (PLE), according to Sailer (2003) and following the idea of Gazdar et al. (1985). PLEs are lexical entries for syntactically complex expressions. Thus, they have properties of both words and phrases. As words, they are licensed by their lexical entry. As phrases, lexical rules cannot apply to them and syntactic operations like topicalization can be excluded by defining structural requirements in their DTRS attribute. According to standard HPSG assumptions we adopt Immediate Dominance

<sup>11</sup>Riehemann (2001), p. 207

<sup>12</sup>as in “Students forced to learn history by rote” from [http://www.shanland.org/Political/News\\_2002/students\\_forced\\_to\\_learn\\_history.htm](http://www.shanland.org/Political/News_2002/students_forced_to_learn_history.htm)

Schemas that license ordinary phrasal signs. In order to exclude the application of ID-Schemas to a phrase licensed by a PLE we can redefine the ID-PRINCIPLE in the following way:

$$[\textit{phrase}]_{\text{COLL } e\text{-list}} \rightarrow \left( \begin{array}{l} \text{HEAD-COMPLEMENT-SCHEMA } \vee \\ \text{HEAD-ADJUNCT-SCHEMA } \vee \\ \text{HEAD-MARKER-SCHEMA } \vee \\ \text{HEAD-FILLER-SCHEMA} \end{array} \right)$$

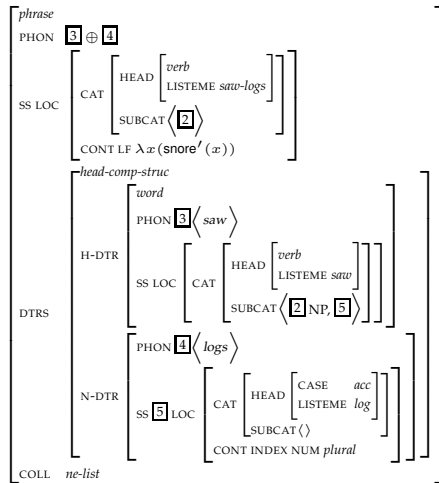
Accordingly, we have to change all principles of grammar that are concerned with regular combination of signs in such a way that they only apply to phrases bearing an empty COLL list. This can simply be done by adding a line in the antecedent (remember that the principles consist of an implication) stating [COLL *e-list*].

In order to specify which lexical entries must have an empty COLL list, we introduce subsorts of *listeme*, namely *coll\_listeme* and *no\_coll\_listeme*, and make the following constraint:

$$[\textit{sign}]_{\text{SS LOC CAT HEAD LISTEME } no\_coll\_listeme} \rightarrow [\textit{coll\_elist}]$$

Note that all lexical entries have different values of LISTEME and, conversely, the set of all LISTEME values covers the entirety of lexical entries.

We have now made a distinction between regular phrasal signs which have an empty COLL list and non-regular or idiomatic phrases having a non-empty COLL list.<sup>13</sup> Thus, in a PLE of an idiom like (3) *saw logs*<sup>14</sup> we define its COLL list as non-empty. Besides, this idiom cannot be passivized without losing its idiomatic reading. Passivization is already excluded by the nature of the PLE itself: an object in accusative case is required and thus, *logs* can never occur as the subject.



In defining a non-empty COLL value, we provide a unified way to treat decomposable and non-decomposable idioms, marking their quality of being idiomatic. Parts of decomposable idioms bear a non-empty COLL list, which restricts their occurrence to

<sup>13</sup>The distribution of COLL values could be easily constrained by another principle which we omit here for reasons of space.

<sup>14</sup>as in "Two young boys stand by their mother's bed while she saws logs in her sleep." from [http://www.collegestories.com/filmfrat/igby\\_goes\\_down.html](http://www.collegestories.com/filmfrat/igby_goes_down.html)

certain contexts. Nondecomposable idioms also have a non-empty COLL list, exempting them from regular syntactic and semantic principles.

In addition, the occurrence of nondecomposable idioms can be restricted to certain contexts via the same feature. This is important for idiomatic intensifiers, among others, like *as a sandboy* in *to be happy as a sandboy* or *as a kite* in *to be high as a kite*.

## 4 Alternative Analyses

### 4.1 A Different COLL Mechanism

The analysis we suggest here is an enhancement of a proposal by Richter and Sailer (1999). However, in Sailer (2003) the author described a variant of the COLL mechanism: In this thesis, the value of COLL is a singleton list that may contain a sign. That sign is the overall expression in which the idiomatic word occurs. Take for example the idiom *spill the beans*: in the lexical entry of the idiomatic word *beans* its COLL value is specified as a sign containing the semantic contributions of a definite article, the idiomatic word *spill* and *beans* itself in the right scopal relations. Sailer defines the so-called COLL-PRINCIPLE ensuring that the sign specified in a COLL list dominates the sign bearing that list. As a consequence, information of the overall utterance is available at lexical level and, conversely, local information is available on each node in the structure.

Thus, even though Sailer introduces only one new attribute, this approach is very unrestrictive and if one taps its full potential, nearly all grammatical phenomena can be described, even if they have nothing to do with collocations. Selection, e. g., would only be a special case of a collocation. Because of this power and unrestrictedness, that version of COLL is to be met with criticism.

### 4.2 A Constructional Approach

Riehemann (2001) makes another concrete proposal for the analysis of idioms. She adopts many ideas of Construction Grammar and carries them forward to the HPSG framework. Her approach requires a complex machinery of new sorts and attributes to cover not only the amount of existing idioms but also their occurrences in different syntactic configurations. She has to assume, e. g., distinct subsorts of a *spill\_beans\_idiom\_phrase* for the idiom occurring in different constructions (e. g. a *head-subject-phrase* or a *head-filler-structure*). Even if the existence of sorts for different constructions themselves is well established in Construction Grammar, it is questionable to assume different subclasses of linguistic signs, only because they contain idiomatic items in different syntactic structures. In other words, why assume different sorts for one single idiom only because it occurs in different constructions?

Moreover, Riehemann herself has to admit that her approach cannot handle cases of pronominal reference like (9), because idiomatic *spill* is not licensed as it seems to appear by itself and not within a *spill\_beans\_idiom\_phrase*. In addition, Riehemann is unable to account for bound words, as she cannot constrain their distribution once she assumes lexical entries for them.

In summary, it seems to us that a lexical approach is to be preferred over a structural one. Nevertheless, her arguments in favor of a constructional analysis of non-decomposable idioms are convincing. Our counterpart to that are phrasal lexical entries which we assume for this kind of idiomatic expressions.

## 5 Prospects for a Modular Approach

We have proposed one way of analyzing idioms and similar phenomena of distributional idiosyncrasies. It can handle distributional characteristics of idiomatic words and even difficult cases like pronominalization.

We decided to take a word-level collocation-based account using the COLL feature. This approach is modular in two ways. Firstly, the barriers can be adjusted “vertically” according to the range (XP, complete clause or utterance) needed for a particular idiomatic expression. Secondly, by the LOC-LIC feature we can specify any characteristics within the local information. We could now go on and define other attributes of *barrier* like PHON-LIC to define any requirements of the phonetic string of that barrier. In that way our approach is also horizontally modular.

An application of such a PHON-LIC feature would be the modelling of occurrence restrictions of the English indefinite article *an*. This phenomenon together with other cases of sandhi is discussed by Asudeh and Klein (2002). With our approach, we define the lexical entry of *an* as follows: the PHON-LIC value of the barrier *np* on the COLL list is the phonetic string  $\langle an \rangle +$  a phonetically realized vowel.

Thus, with a quite general approach to idioms using the COLL feature, we can handle very particular phenomena, too. We leave it to further research to explore the possibilities that our approach holds.

## References

- Aronoff, Mark (1976). *Word Formation in Generative Grammar* (3rd print 1985 ed.). MIT Press, Cambridge MA. Linguistic Inquiry Monographs.
- Asudeh, Ash and Klein, Ewan (2002). Shape Conditions and Phonological Context. In F. van Eynde, L. Hellan, and D. Beermann (Eds.), *Proceedings of the 8th International HPSG Conference*, pp. 20–30. CSLI Publications. <http://csli-publications.stanford.edu/HPSG/2/>.
- Di Sciullo, Anna-Maria and Williams, Edwin (1988). *On the Definition of Word* (Second ed.). Linguistic Inquiry Monographs. MIT Press, Cambridge, Mass.
- Dobrovolskij, Dmitrij (1988). *Phraseologie als Objekt der Universalienlinguistik*. VEB Verlag Enzyklopädie Leipzig. Linguistische Studien.
- Gazdar, Gerald, Klein, Ewan, Pullum, Geoffrey, and Sag, Ivan (1985). *Generalized Phrase Structure Grammar*. Cambridge, Mass.: Harvard University Press.
- Heinz, W. and Matiassek, J. (1994). Argument Structure and Case Assignment in German. In J. Nerbonne, K. Netter, and C. Pollard (Eds.), *German in Head-Driven Phrase Structure Grammar*, pp. 199–236. CSLI Publications. Lecture Notes 46.
- Krenn, Brigitte and Erbach, Gregor (1994). Idioms and Support Verb Constructions. In J. Nerbonne, K. Netter, and C. Pollard (Eds.), *German in Head-Driven Phrase Structure Grammar*, pp. 365–396. CSLI Publications. Lecture Notes 46.
- Nunberg, Geoffrey, Sag, Ivan A., and Wasow, Thomas (1994). Idioms. *Language* 70, 491–538.
- Richter, Frank (1997). Die Satzstruktur des Deutschen und die Behandlung langer Abhängigkeiten in einer Linearisierungsgrammatik. Formale Grundlagen und Implementierung in einem HPSG-Fragment. In E. Hinrichs, D. Meurers, F. Richter, M. Sailer, and H. Winhart (Eds.), *Ein HPSG-Fragment des Deutschen, Teil 1: Theorie*, Number 95 in Arbeitspapiere des SFB 340, pp. 13–187. Universität Tübingen.
- Richter, Frank and Sailer, Manfred (1999). LF Conditions on Expressions of Ty2: An HPSG Analysis of Negative Concord in Polish. In R. D. Borsley and A. Przepiórkowski (Eds.), *Slavic in Head-Driven Phrase Structure Grammar*, pp. 247–282. Stanford: CSLI Publications.
- Riehemann, Susanne Z. (2001). *A Constructional Approach to Idioms and Word Formation*. Ph. D. thesis, Stanford University, Stanford, CA.
- Sailer, Manfred (2003). Combinatorial Semantics and Idiomatic Expressions in Head-Driven Phrase Structure Grammar. Phil. Dissertation (2000). Arbeitspapiere des SFB 340. 161, Eberhard-Karls-Universität Tübingen.
- Soehn, Jan-Philipp and Sailer, Manfred (2003). At First Blush on Tenterhooks. About Selectional Restrictions Imposed by Nonheads. In G. Jäger, P. Monachesi, G. Penn, and S. Wintner (Eds.), *Proceedings of Formal Grammar 2003*, pp. 149–161.