# Basic Relations and Stereotype Relations in the Semantics of Compound Nouns[*]

Melanie J. Bell

*Anglia Ruskin University, UK*
*melanie.bell@anglia.ac.uk*

This paper tests the hypothesis of Fanselow (1981) that the semantic relations in compound nouns are of two types: 'basic' and 'stereotype'. It is shown that the probability of a compound falling into either of Fanselow's proposed categories can be largely predicted using semantic and distributional properties of the constituent nouns, as well as the degree of lexicalisation of the compound as a whole. The so-called 'basic' relations, namely constitution, location, identity, resemblance and meronymy, are more likely in compounds that are not lexicalised, that have productive modifiers and/or semantically-specific heads, and whose constituents are perceived as representing concrete rather than abstract concepts. It is argued that such relations might be regarded as basic in several ways: they relate to states of physical entities, have a high level of generality and may be associated with semantic and phonological transparency.

Keywords: *Compound Noun, Semantic Relation, Semantic Class, Basic, Stereotype, Lexicalisation, Productivity, Specificity, Concreteness, Transparency, Generality, Modification, Morphological Family*

## 1. Introduction

The semantic interpretation of a compound word involves not only the senses of the constituent words that make up the compound, but also the semantic relation or relations between them. Whereas the former are overtly represented in the form of the compound, the latter are unexpressed, and because of this, interpretation depends to a greater or lesser extent on extralinguistic knowledge. For example, the semantic relation in *tea cup* might be labelled FOR, since a tea cup is a cup intended for drinking tea. The semantic relation in *gold cup* might be either MADE OF (a cup made of gold) or RESEMBLES (a cup with the colour of gold). Similar analyses can be applied to longer compounds. In *desk-top computer*, for example, the relation between DESK and TOP might be classified as either HAS (the desk has a top) or WHOLE-PART (the top is part of the desk), while the relation between DESK-TOP and COMPUTER might be classified as either FOR or ON, since a desk-top computer is a computer intended for use on a desktop. Any additional difficulty in interpreting compounds with more than two constituents lies in the initial semantic parsing, rather than in the possible semantic relations, so for ease of exposition, since the focus here is on semantic relations, the examples used in the rest of this paper will be compounds consisting of two nouns, N1 (the left-hand constituent) and N2 (the right-hand constituent).

Many attempts have been made to define and categorise the possible relations between compound constituents, both in the mainstream and computational linguistic literature, (e.g. Jespersen 1942, Hatcher 1960, Lees 1960, Marchand 1969, Levi 1978, Warren 1978, Lauer 1995, Nastase & Szpakowicz 2003, Girju, Moldovan, Tatu & Antohe 2005, Ó Séaghdha 2007). The relations identified are sometimes defined in terms of descriptive phrases, e.g. 'N2 is made of N1', and sometimes in terms of underlying predicates, e.g. BE. However, such classifications are problematic for a number of reasons and no fully satisfactory or inclusive system has been devised, with the result that most taxonomies include a category 'other'. Even with that option included, naïve human raters have very poor levels of agreement in classification tasks, and even trained lexicographers achieve at most 65-75% agreement, depending on the dataset used (Ó Séaghdha

2007: 45ff.). When semantic classification of compounds is required for experimental purposes, one solution is to use only those items about which raters agree. In applications where it is necessary to assign a value to every item, for example in computational natural language processing, the usual practice is either to take a majority vote amongst a group of raters or to thrash out the difficult items until some consensus is achieved. Ó Séaghdha (ibid.: 51-53) summarises the difficulties of achieving annotator agreement for compound semantic classification.

The first problem in trying to list the semantic relations that can occur in compounds is that it is difficult to achieve exhaustive coverage of the relations found, let alone define the boundaries of what is possible. Downing (1977), in a seminal experimental study, attempted to discover whether there are any constraints on the kinds of semantic relationships that can be expressed by compounding. She concluded that the list of possible relationships is virtually limitless. Furthermore, she found that using a limited number of classes involved the loss of 'much of the semantic material considered by the subjects to be necessary or essential to the definitions' (ibid.: 826). These problems have led some authors, e.g. Jespersen (1942: 137), Bauer (1979) and Lieber (2004: 53), to suggest that only the values of N1 and N2 in a compound are semantically specified, with the relation being interpreted pragmatically from context and encyclopaedic knowledge.

The second problem with compound classification systems is that they are hard to apply consistently. This is partly because, in some cases, more than one relation might be taken to apply to the same compound. For example, *knife handle* might be assigned to any of the classes ON (a handle on a knife), FOR (a handle for a knife), WHOLE-PART (a handle that is part of a knife) or HAS (a handle that a knife has). In other cases, it might be hard to assign any category at all without a context. Bell & Plag (2012: 493) give the example *peach thing*, which might denote 'something peach coloured, made of peaches, or designed for holding peaches, amongst other possibilities'.

A third problem is that it is hard to separate the senses of constituents themselves from the relation between them. Tarasova (2013: 107) shows that in a stratified random sample of compounds from the New Zealand

printed media, 76% of constituent nouns show a significant preference for one semantic relation over others. Furthermore, in a study based on the British National Corpus (BNC), Maguire, Wisniewski & Storms (2010) show that this tendency applies not only to individual nouns, but also to semantic classes of nouns. For example, when N1 represents a time period and N2 represents an event, the semantic relation is 'N2 during N1' in 89% of cases (Maguire, Wisniewski & Storms 2010). In keeping with these findings, some approaches to compound semantics regard the semantic relations as emerging from the elaborated lexical entries of the constituents. Such approaches include the cognitive tradition of Langacker (1987), the generative lexicon framework of Pustejovsky (1991) and the lexical semantic approach of Lieber (2009, in contrast to Lieber 2004).

Despite all the difficulties involved in classifying compound semantic relations, Ó Séaghdha (2008: 19) concludes that 'the inventories that have been proposed are more notable for their commonalities than their differences'. One possible reason for this core similarity might be that some compound relations are more basic than others (cf. Fanselow 1981), so that all taxonomies of semantic relations tend to converge in recognition of the basic relations while varying in their analysis of others. The aim of this paper is to investigate the hypothesis that some compound semantic relations can be regarded as 'basic' and that the probability of a compound involving such a relation can be predicted from the distribution and perceived concreteness of its constituents, as well as their semantic class and degree of specificity. The rest of the paper is organised as follows: Section 2 gives the background to the study, including the rationale for the hypotheses to be tested, Section 3 outlines the methodology used in the empirical investigation, Section 4 describes the statistical models generated and Section 5 concludes with a discussion of the results.

## 2. Background

### 2.1. Basic Relations and Stereotype Relations

Fanselow (1981) argues that compound semantic relations can be divided into two types: a limited number of basic relations *(Grundrelationen)*

and a much larger, possibly unlimited, number of stereotype relations *(Stereotyprelationen)*. He identifies five basic relations, which are AND *(und)*, MADE OF *(gemacht aus)*, SIMILAR TO *(ähneln)*, PART OF *(ist Teil von)*, and LOCATED RELATIVE TO *(ist lokalisiert bezüglich)*, which can refer to location in space or time (ibid.: 156). Some of his examples, which are all from German, are shown in Table 1. In this paper, for ease of exposition, I will use the terms 'basic relation' to refer to the five types identified by Fanselow as basic, and 'stereotype relation' to refer to other compound relations. However, this use of the terminology should not be taken to indicate any commitment to Fanselow's position.

**Table 1.** Basic relations in compounding according to Fanselow (1981)

| Relation | Examples | |
| --- | --- | --- |
| *und*<br>AND | *Hausboot*<br>houseboat | *Radio-Uhr*<br>radio-clock |
| *gemacht aus*<br>MADE OF | *Steinblock*<br>stone block | *Roggenbrot*<br>rye bread |
| *ähneln*<br>SIMILAR TO | *Flammenschwer*<br>flame-bladed sword<br>(lit. flame sword) | *Blutbuche*<br>copper beech<br>(lit. blood beech) |
| *ist Teil von*<br>PART OF | *Autokotflügel*<br>car wing | *Kammzinke*<br>comb tooth |
| *ist lokalisiert bezüglich*<br>LOCATED RELATIVE TO | *Dezembertagung*<br>December meeting | *Küstenstraße*<br>coast road |

According to Fanselow, (ibid.: 158), the defining characteristic of the basic relations is that they depend on organisational principles of perception or of semantic classification that are independent of the meanings of individual words. In contrast, the stereotype relations arise from the conceptual structures of the compound constituents, and compounds can also develop their own stereotypes. To illustrate stereotype relations, Fanselow (ibid.) gives the examples of *newspaper woman (Zeitungsfrau)* and *book shop*

*(Buchgeschäft)*; he argues that in *newspaper woman* the inferred relation DELIVERS is attributable to the conceptual structure of *newspaper* and that in *book shop* the conceptual structure of *shop* is responsible for the inferred relation SELLS. Fanselow (ibid.) hypothesises that general principles of rationality and co-operation mean that speakers who coin compounds will only do so where the unexpressed relation is likely to be apparent to the listener or reader, and that the most obvious source of such relevance is the meanings of the compound constituents. However, he further asserts that not all compounds conform to this generalisation and that in such cases, where no relation can be inferred from the conceptual structures of the constituents, one of the basic relations applies. Fanselow's (ibid.) operational definition of 'basic relation' is that a compound AB has a basic relation if the most explicit paraphrase of AB contains nothing that has do to with the meaning of either A or B. For example, in *politician-composer (Politiker-Komponist)* and *coastal town (Küstenstadt)*, Fanselow (ibid.) argues that the respective relations, AND and LOCATED RELATIVE TO, have nothing to do with the meanings of the constituents, and that these are therefore basic relations. In contrast, he cites *nursery (Kinderzimmer)* as an example where the definition of basic relation does not apply; he argues that *nursery* is not explicitly paraphrased as a 'room for children' but rather as a 'room where children usually live or are cared for', and that this stereotype relation is attributable to the conceptual structure of *room.*

Some of Fanselow's (1981) examples are more convincing than others, and the assertion that basic relations have nothing to do with the meanings of the constituents is probably too strong. For example, in the case of *coastal town*, it could plausibly be argued that the relation LOCATED RELATIVE TO arises from the meaning of *coast*, since a coast is an area in which things can be located, namely 'the land beside or near to the sea' (Oxford Advanced Learner's Dictionary 1995). Nevertheless, the basic relations do seem have some things in common which distinguish them from other relations. In particular, they relate to properties shared by all physical entities, namely that they consist of matter (MADE OF), occupy a time and space (LOCATED RELATIVE TO), and have a size, shape, colour and chemical constitution (all of which can be conceptualised using the relation SIMILAR TO). Furthermore, an object may belong to more

than one class of objects simultaneously (AND), and may be composed of components (PART OF). In other words, Fanselow's (ibid.) basic relations of constitution, spatial and temporal location, identity and meronymy arise from the laws of physics, rather than any more specific property of particular concepts. As Fanselow (ibid.) puts it:

> If one learns the meaning of e.g. hammer, then one has to learn that it holds of things with a specific form and function, if one learns the meaning of nail, that these are things to be hammered into walls etc. But one does not need to learn that their denotations are located somewhere, that they can belong to other denotations, or that they are made out of something etc. (Fanselow 1981: 158)

Fanselow's (1981) assertion that one does not need to learn that the denotations of *hammer* and *nail* are located somewhere or that they are made out of something is of course only true provided one knows that these words denote physical objects. His argument seems to be that, assuming we know the head of a compound represents a physical entity, it will always be possible to form compounds with that head using any of the basic relations, without needing to know any further detail about the concept represented. In comparison, the stereotype relations are less universally applicable and sometimes relate more to human activity and experience than to the laws of nature. Although taxonomies vary, relations that do not fall into Fanselow's (ibid.) basic type include, for example, USE, FOR, ABOUT and CAUSE (e.g. Levi 1978: 76-77). The first three of these, relating to use, purpose and topic, only make sense in terms of human, and perhaps other animal, behaviour. The causal relation is distinguished from the basic relations by the fact that not every object has an easily identifiable cause or result and, furthermore, causation suggests an event whereas the basic relations describe states.

In his choice of the term 'basic' for relations associated with physical states, Fanselow (1981) may or may not be assuming that the physical universe is more basic than the psychological or social. The recognition that certain compound semantic relations can potentially apply whenever the head element represents a physical entity, hinges neither on terminology nor on philosophical stance. However, it does lead to a hypothesis about

how the constituents of compounds with basic relations might differ from those of compounds with stereotype relations. Because the basic relations describe physical states, it is to be expected that the likelihood of a compound using a basic relation will be greater when the constituent nouns represent physical entities. We might therefore predict that the probability of a compound involving a basic relation will be positively correlated with the perceived concreteness of N2 and perhaps also of N1.

## 2.2. Semantic Relations and English Compound Stress

Circumstantial evidence for the distinction between basic and stereotype relations comes from the phonological stress patterns of English compound nouns. Although many English compounds have the left stress pattern characteristic of Germanic compounds, e.g. `tea cup` and `help desk,` there are also many that are normally pronounced with main stress on the second noun, e.g. *silk `shirt* and *afternoon `nap* (where ` indicates main stress).

   It has long been asserted in the Anglicist literature that certain semantic relations between compound constituents, as well as particular semantic classes of N1 or N2, are associated with stress on N2, i.e. with so-called right stress or right prominence. Table 2 shows the correspondence between the categories identified by various authors in the nineteenth and twentieth centuries. It can be seen that these categories largely overlap those described by Fanselow (1981) as having basic relations. The first two rows correspond roughly to Fanselow's (ibid.) AND: a gentleman farmer is both a gentleman and a farmer, and a man cook is both a man and a cook. The third row corresponds to Fanselow's (ibid.) MADE OF (silk thread is thread made of silk and rubber boots are boots made of rubber), while the fourth row corresponds to his SIMILAR TO (a bow window is a window whose shape resembles that of a bow). The next two rows correspond to the temporal and spatial aspects of Fanselow's (ibid.) LOCATED RELATIVE TO: a morning paper is a newspaper published in the morning and a kitchen sink might be classified as a sink in a kitchen. Although the PART OF relation has not usually been associated with right prominence, there is some overlap between this and the LOCATED RELATIVE TO relation, e.g. a kitchen sink could also be considered part of a kitchen and a garage door might be

**Table 2.** Semantic categories claimed to produce right prominent compound nouns in English

| Sweet 1891 | Marchand 1969 | Fudge 1984 | Liberman & Sproat 1992 |
|---|---|---|---|
| | Copulative combinations: *gentleman farmer* | N1 is semantically central and N2 specifies it further: *knight bachelor* | N2-is-a-N1: *rogue elephant* |
| N1 defines the sex or age of the N2: *man cook* | Combinations with N1 denoting sex or age: *baby__* | | |
| N1 denotes material of which N2 is made: *silk thread* | All combinations of the *stone wall* type | N1 is a material and N2 is made of N1: *cotton dress* | N2 is thing-made-out-of-N1: *rubber boots* |
| N1 expresses something that resembles N2: *bow window* | | | |
| | | N1 is a time or season: *morning paper* | N1 is time-when-N2-occurs: *fall weather* |
| | Combinations with N1 denoting relational position: *bottom__* | N1 is a location: *kitchen sink* | N1 is place-where-N2-is-found: *garage door* |
| | | N1 specifies the value of N2: *pound note* | N1 is measure and N2 is thing-measured: *mile run* |
| | | N1 and N2 form a proper name: *William Smith* | N1 is classifier and N2 is name: *Bayou Goula* |
| N2 is a general place-word and N1 is a proper name: *Oxford Road* | | N2 is geographical or a thoroughfare and N1 is the name applied to N2: *Thames valley Shaftesbury Avenue* | N1 is proper-name applied to N2: *Connecticut Yankee* |

considered part of a garage. The final three rows in Table 2 show relations that Fanselow (ibid.) explicitly excludes from consideration, namely those in which the first noun is a measure term, and those that are, or include, proper names. These will therefore not be considered further in the present paper.

   Plag, Kunter, Lappe & Braun (2008) conducted a large-scale empirical test of the hypothesis that semantic relations and the semantic classes of constituents can predict English compound prominence. They used compounds taken from the Boston University Radio Speech Corpus (Ostendorf, Price & Shattuck-Hufnagel 1996), a corpus of spoken news recordings, from which it was possible to rate the prominence of compound tokens using acoustic criteria. They categorised these compounds semantically using five classes of constituents and eighteen semantic relations. Each compound was categorised by two independent judges, and only those items for which the judgements concurred were included in the subsequent analysis. The semantic variables were then used, along with other potential predictors, as predictors of prominence in a regression analysis. Plag, Kunter, Lappe & Braun (ibid.) found that certain semantic classes and relations were the best predictors of prominence, including the relations shown in Table 3.

**Table 3.** Semantic relations shown to produce right prominent compound nouns by Plag, Kunter, Lappe & Braun (2008)

| Semantic relation |
| --- |
| N1 is N2 |
| N2 is made of N1 |
| N2 located at N1 |
| N2 during N1 |
| N1 has N2 |
| N2 is named after N1 |

It can be seen that the rightward-leaning relations in the Boston data largely correspond to those identified in the earlier descriptive literature, and therefore to Fanselow's (1981) basic relations: 'N1 is N2' corresponds to AND, 'N2 is made of N1' corresponds to MADE OF, 'N2 located at N1' and

'N2 during N1' correspond to LOCATED RELATIVE TO. The relation 'N1 has N2' also has some degree of overlap both with LOCATED RELATIVE TO and with PART OF; for example, as mentioned in the introduction, a knife handle might be classified not only as a handle that a knife has but also as a handle on a knife or a handle that is part of a knife. The correlation of this sub-set of relations with a particular stress pattern suggests that those relations hypothesised by Fanselow (1981) to form a basic set in German might also form some sort of natural class in English.

Although Plag, Kunter, Lappe & Braun (2008) confirmed that certain semantic relations are predictive of stress pattern in English noun-noun compounds, subsequent work (e.g. Plag, Kunter, Lappe & Braun 2007, Plag 2010) has shown that the strongest predictors of a compound's stress are the identities of its constituent nouns. In other words, any given noun in either N1 or N2 position will tend to be associated with a particular stress pattern. Since it is known that particular constituents are also associated with particular semantic relations (Tarasova 2013: 101-107), this means that there is a three-way correlation in English compounds between constituent identity, stress pattern and semantic relation. Furthermore, since rightward stress is associated with the basic relations, it is likely that those constituents predictive of rightward stress will also favour a basic relation. If such constituents can be shown to have properties in common, these properties might explain not only the tendency of these nouns to produce rightward stress in compounds but also their tendency to occur with one of the basic relations. If so, then what the basic relations have in common, and perhaps even the sense in which they can be regarded as 'basic', might be best understood in terms of the properties of their constituents.

Bell & Plag (2013) investigated what properties of nouns determine their tendency to produce a particular stress pattern in compounds. They found that for nouns in N1 position, the tendency to produce right stress increases as the number of syllables increases and as the positional family size increases. Positional family size is the number of different nouns with which a given constituent combines in a given corpus. For example, the positional family for *house* in N1 position might include such combinations as *houseboat, house party* and *house cat*. Since the first noun in English NN compounds usually functions to modify the second, the positional family

size of N1 can be taken to reflect the productivity of that noun as a modifier. Provided the calculation of family size includes novel as well as established compounds, then the greater the family size, the more nouns N1 can modify and the more productive it therefore is as a modifier (cf. Plag 1999: 22ff.). Since the probability of right stress increases with the positional family size of N1, we can infer that right stress is most likely when N1 is a productive modifier. For nouns in N2 position, Bell & Plag (ibid.) found that the probability of right stress again increases as the number of syllables increases, and also as the number of senses of N2 decreases. The more senses a noun has, the less semantically specific it might be taken to be, so that rightward stress is correlated with more semantically specific nouns in N2 position, which in English is usually the position of the semantic and syntactic head. In summary, productive modifiers and semantically specific heads, as well as longer constituents in either position, predispose English compounds to right stress.

Bell & Plag (2012, 2013) interpret the effect of constituent length on compound stress as a reflection of the tendency for speakers of all languages to avoid long strings of unaccented syllables (cf. Ladd 1996: 244). This seems likely to be a purely phonological effect and it is hard to see how it would correspond to particular semantic relations. On the other hand, the effects of modifier productivity and head specificity on compound stress are more semantic in nature, and it therefore seems more plausible that they might extend to predicting semantic relations. Since basic relations are associated with right stress, we might expect that constituent properties that increase the probability of a compound having right stress will also increase the probability of it having a basic relation, in other words, that basic relations will be associated with productive modifiers and semantically specific heads.

A final prediction arising from work on English compound stress concerns the effect of lexicalisation. All recent studies on the subject (e.g. Plag, Kunter, Lappe & Braun 2007, 2008; Bell & Plag 2012, 2013) have established that the probability of a compound having right stress is inversely correlated with its degree of semantic lexicalisation, whatever measure of lexicalisation is chosen. It follows that, since basic relations are correlated with right stress, they will also be inversely correlated with

lexicalisation; in other words, basic relations are unlikely to be found in highly lexicalised compounds.

## 2.3. Basic Relations and Semantic Classes

In addition to the association of individual constituents with particular semantic relations, certain semantic classes of compound constituents are also associated with particular relations. This intuitively plausible correlation was confirmed empirically by Maguire, Wisniewski & Storms (2010). In a study based on the BNC they found, for example, that the MADE OF relation occurred in 68% of compounds where N1 represented a substance and N2 represented an artefact, the DURING relation occurred in 89% of compounds where N1 represented a time period and N2 represented an event, and the LOCATED relation occurred in 91% of compounds where N1 represented an area and N2 represented an animal. It is therefore to be expected that particular classes of constituents will be associated with Fanselow's (1981) basic relations. Perhaps most obviously, modifiers that represent materials or substances are likely to occur with the basic relation MADE OF, while modifiers that represent locations in time or space are likely to occur with the basic relation LOCATED RELATIVE TO.

## 2.4. Hypotheses

To summarise the discussion in Section 2, I have argued that if the semantic relations identified as basic by Fanselow (1981) can be distinguished from other compound semantic relations, then the distinction is likely to involve properties of the constituents of the compounds in question. The aim of the paper is to predict the probability of any compound having a basic relation on the basis of the properties of its constituents. It was argued in Section 2.1 that the basic relations correspond to states of physical entities, as opposed to events or human activity; from this follows the hypothesis that basic relations will be more likely in compounds whose constituents have a high level of perceived concreteness. In Section 2.2, it was argued that the basic relations are likely to be associated with similar properties to those that predict right stress in English compounds; from this follow

the hypotheses that basic relations will be more likely in compounds with productive modifiers and semantically specific heads, and less likely when the compound as a whole is semantically lexicalised. Finally, in Section 2.3, it was argued that basic relations will be more likely in compounds whose constituents belong to particular semantic classes. In addition, it is likely that these various effects will interact with one another. The paper will test the following predictions:

1. The probability of a compound having a basic relation will be positively correlated with the perceived concreteness of its constituents, i.e. the greater the perceived concreteness of the constituents, the more likely is the compound to involve a basic relation.

2. The effect of perceived concreteness will interact with semantic specificity. For example, where a constituent has a high concreteness rating but is highly polysemous, a particular compound might use a sense that is less concrete than the core sense.

3. The probability of a compound having a basic relation will be positively correlated with the productivity of N1 as a modifier, i.e. the more productive the modifier, the more likely is the compound to have a basic relation.

4. The probability of a compound having a basic relation will be greater in compounds whose modifier represents a material, a substance, a place, or a point or period in time.

5. All of the above effects will be modified by the extent to which the compound is lexicalised, such that lexicalisation will reduce the overall probability of a basic relation.

## 3. Method

### 3.1 Dataset

The dataset used was the set of 1000 compounds originally selected for the analyses described in Bell (2013). Because the data was originally intended for phonological analysis, the compounds were selected from the demographically-sampled spoken conversation section of the British National Corpus. This was done by searching for strings of two nouns,

ordering them randomly, then checking them manually in context until 1000 compounds were found (since not all noun-noun strings are compounds).

## 3.2 Coding the Dependent Variable

In order to produce statistical models to test the hypotheses outlined in Section 2.4, it was first necessary to code each of the compounds in the dataset for the presence or absence of a basic relation. But as described in the introduction, coding compound semantic relations is notoriously difficult because of the inherent ambiguity or vagueness of the constructions involved. Compounds can show at least two types of ambiguity. Firstly, the same noun-noun string can have more than one possible denotation, e.g. *steel warehouse* can represent either the set of warehouses made of steel or the set of warehouses in which steel is stored. In the first case, *steel warehouse* involves the basic relation MADE OF, whereas in the second case, it involves a stereotype relation based on the stereotype of a warehouse as a place where things are stored. Secondly, over and above this possible ambiguity of denotation, many compounds can be interpreted in a variety of ways, even when the denotation remains the same. For example, peanut butter can be conceptualised as a kind of butter made from peanuts. With this interpretation, PEANUT BUTTER is the transparent combination of the two concepts PEANUT and BUTTER with the basic relation MADE OF. On the other hand, exactly the same peanut butter can also be conceptualised as a kind of food stuff unrelated to butter (in the dairy sense). With this interpretation PEANUT BUTTER is a single concept, not simply the combination of PEANUT and BUTTER; in Fanselow's (1981) terms, the compound has developed its own stereotype. The fact that *peanut butter* is well known to be produced with both left and right prominence patterns supports the hypothesis that it has two possible semantic analyses.

Fanselow (1981: 184) mentions that ambiguity is very hard to avoid in compounding, and discusses the example of *car wing (Autokotflugel)*. A car wing is not necessarily a wing that is part of a car, or a wing that is located on a car, since the wing can be removed from the car and the corresponding car does not need to continue to exist, nor does there need to be a car at all in order for a thing to be a car wing. This kind of relation is particularly

difficult to code consistently since, in addition to being coded as a wing on a car or a wing that is part of a car, *car wing* could also be coded as a wing for a car or a wing that a car has. Fanselow (ibid.) argues that the location relation actually has less usage than commonly assumed; for example, he suggests that it would be strange to interpret *night worker (Nachtarbeiter)* as a 'worker at night time', presumably because a person who is a night worker remains so even during the day.

   Compounds in which N1 could possibly be interpreted as denoting the spatial or temporal location of N2 may or may not entail this relation. For example, *X is a bedroom wall* usually entails *X is in a bedroom*, since an identical wall in a different type of room would be unlikely to be classified as a bedroom wall. On the other hand, although some judges might code *door handle* as a handle located on a door, *X is a door handle* does not entail *X is on a door*: a door handle is still perceptibly a door handle when it is, for example, on a shelf in a shop. The first type, where the location is entailed, is very similar to the PART OF relation, e.g. a bedroom wall might be regarded as part of a bedroom. The second type, where a locative relation is not entailed, might alternatively be classified as having the relation FOR: a door handle is a handle for a door, i.e. designed to go on a door, rather than necessarily on one. This purposive relation has what Levi (1978: 99) calls a 'characteristic vagueness'. She gives the examples of *fertility pills*, where the pills are for fertility in the sense of increasing it, contrasted with *headache pills,* where the pills are for headaches in the sense of decreasing them. Clearly the most explicit paraphrases of such purposive compounds involve the stereotypes of the constituents, here *fertility* and *headache* respectively. In other words, compounds where the relation can be coded as FOR tend to involve stereotype relations, and this is true even in cases where some judges might assign a locative relation; for example, a door handle is a handle for opening a door, and this explicit paraphrase arises from the stereotype of a door as something to be opened. On the other hand, compounds in which a locative relation is entailed can be regarded as having a basic relation. In order to reduce ambiguity, the compounds used in this study were therefore coded using a criterion of entailment: compounds were only classed as having the relation LOCATED RELATIVE TO if the compound's interpretation was judged to entail the

location. For consistency, a similar criterion was applied to the coding of other relations.

A professional lexicographer was employed at approximately the commercial rate to code the semantic relations. For each relation coded, she was asked to make a semantic judgement about whether it applied to each of the compounds in the dataset. Because the items were presented out of context, there were inevitably some that were ambiguous or whose meaning was otherwise unclear. In such cases, judgements were based on what the rater considered to be the likely or possible meaning or meanings, rather than making any attempt to look them up. The same tasks were also completed by the author. The procedure was that the tasks were presented to the lexicographer with a few examples done by the author. The lexicographer then completed the first hundred items and returned the results for comparison. Discrepancies were discussed and misunderstandings ironed out before the remainder of the task was completed. At the end, the results were again compared and any discrepancies discussed. Any disagreements which could not be resolved resulted in the relevant item being excluded from further analysis. The details of the task for each of Fanselow's (1981) basic relations are outlined below:

AND
The compounds were classified according to whether or not the following statement was true:
   *X is (an) NN* entails *X is (an) N2* and *X is (an) N1*
   e.g. *X is a singer songwriter* entails *X is singer* and *X is a songwriter*

MADE OF
The compounds were classified according to whether or not the following statement was true:
   *X is (an) NN* entails *X is (an) N2* and *X is made of/with N1*
   e.g. *X is a silk shirt* entails *X is a shirt* and *X is made of silk*
      *X is tuna risotto* entails *X is risotto* and *X is made with tuna*

SIMILAR TO
This relation was treated as being equivalent to AND, except that the

meaning of the first constituent is shifted before the relation is applied. The compounds were classified according to whether or not the following statement was true:

   *X is (an) NN* entails *X is (an) N2* and *X is metaphorically (an) N1*

   e.g. *X is a queen ant* entails *X is an ant* and *X is metaphorically a queen*

      *X is a giant step* entails *X is a step* and *X is metaphorically a giant*

      *X is a stiletto heel* entails *X is a heel* and *X is metaphorically a stiletto*

In the above examples, the ant is similar to a queen in terms of social position, the step is similar to a giant in terms of size, and the heel is similar to a stiletto in terms of shape.

LOCATED RELATIVE TO

The compounds were classified according to whether or not the following statement was true:

   *X is (an) NN* entails *X is (an) N2* and *X is at/in/on (an) N1*

   e.g. *X is a London school* entails *X is a school* and *X is in London*

      *X is a Monday morning* entails *X is a morning* and *X is on a Monday*

In these examples, the condition of entailment is satisfied because an identical school in a different location would not be a London school and a morning on any other day of the week is not a Monday morning.

In some cases, a compound was felt to be too ambiguous for the entailment condition to be satisfied. For example, an 'olive thing' could be a dish made with olives, in which case the statement *X is made with N1* would be true. However, an 'olive thing' could also be an olive-coloured cloth or an implement for removing olive stones and, in either of these cases, the statement would be false. So *X is an olive thing* does not entail *X is made of olives*. In other cases, although the compound was less vague, it was still felt to have more than one distinct meaning, only one of which satisfied the condition of entailment. For example, if *steel warehouse* denotes warehouses made of steel, the condition is satisfied, but if it denotes warehouses for storing steel, the condition is not satisfied. In such cases, the compound was marked as ambiguous and excluded from the analysis.

PART OF

This relation was not included in the original coding exercise for the analyses in Bell (2013), but was added later for the purposes of the present paper. Four adult native speakers of standard British English identified those compounds in the dataset where the first noun represented a whole thing and the second noun represented part of that thing, e.g. *car wing* (the wing is part of the car) or *comb tooth* (the tooth is part of the comb). The compounds were coded as positive for this relation only in cases where at least three of the four judges agreed that the relation applied.

### 3.3 Coding the Predictors

Each of the compounds in the data was coded for the variables hypothesised to predict the likelihood of a compound having a basic relation. These are summarised in Table 4.

**Table 4.** Predictors initially present in the analysis

| N1 semantic classes | Semantic specificity | Distribution of constituents | Concreteness of constituents | Lexicalisation |
|---|---|---|---|---|
| material | number of synsets N1 | family size ratio of N1 | concreteness of N1 | spelling ratio |
| time or place | number of synsets N2 | family size ratio of N2 | concreteness of N2 | |

*3.3.1 N1 Semantic Classes*

In order to get the most objective possible rating of semantic categories, the WordNet lexical database was used (Miller 1995). In this database, words are grouped together into sets of cognitive synonyms (synsets), each of which represents a distinct concept. A word can belong to a number of different synsets, each representing a different sense of the word, and synsets are linked in the database to related synsets which might represent, for example, hypernyms, hyponyms, synonyms or antonyms of the sense in question. Given a list of words, it is therefore possible to use WordNet to find which of those words are hyponyms of a particular concept, for

example MATERIAL. The compounds in the dataset were classified according to whether N1 was listed as a hyponym of any of the synsets shown in Table 5.

**Table 5.** WordNet synsets for N1 classes associated with basic relations

| N1 semantic class | WordNet synset | WordNet definition |
|---|---|---|
| material | substance (sense 1) | that which has mass and occupies space |
| | fabric (sense 1) | cloth, material, space |
| | building material | material used for constucting buildings |
| time or place | location (sense 1) | a point or extent in space |
| | time period | period of time |

For each value of N1 in the data, the WordNet 2.1 browser was manually searched to find out whether any sense of N1 was listed as a hyponym of the synsets in Table 5. Since not all senses of such nouns were hyponyms of the relevant categories, it was then necessary to determine which sense applied to the compound or compounds in the dataset. To do this, two native speakers of English, the author and a volunteer with a degree in modern languages, manually checked WordNet and decided which sense of N1 applied to each of the compounds in question. For example, in *paper plate, paper* has the sense 'a material made of cellulose pulp derived mainly from wood or rags or certain grasses', whereas in *paper title*, it has the sense 'a scholarly article describing the results of observations or stating hypotheses'. In some cases, more than one sense seemed simultaneously applicable, in which case both senses were recorded as applying. In other cases, the compound was ambiguous, and more than one sense of N1 was possibly applicable depending on the interpretation of compound. These senses were recorded as 'possibly applicable'. The results of these two rating exercises were compared. If both raters thought that a sense was either applicable or possibly applicable, then it was accepted as a possible sense in the context of the compound. If either rater felt that a sense was not applicable where the other thought it was, or possibly was, the item was referred to a professional lexicographer who made a definitive judgement. Having decided on the relevant sense or senses of N1 for each compound, a

check was then made to see whether any of these senses was a hyponym of any of the synsets in Table 5. Compounds for which N1 was a hyponym of 'substance', 'fabric' or 'building material' were combined to give a category 'N1 is a material'. Compounds where N1 was a hyponym of 'location' or 'time period' were combined in a category 'N1 is time or place'.
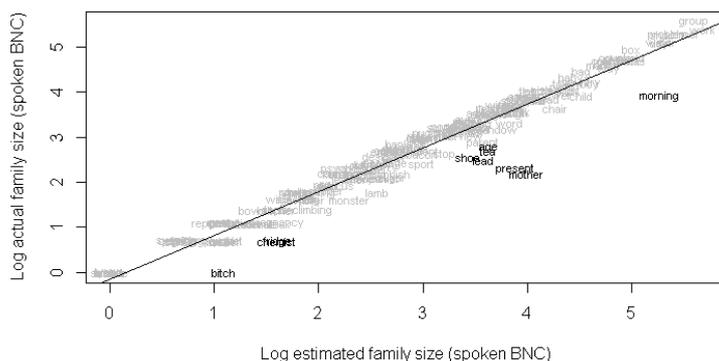
### 3.3.2 Semantic Specificity

WordNet was also used to estimate the semantic specificity of the compound constituents. The synset count for a word is the number of synsets in which it appears in the WordNet database. Since each synset corresponds to a different concept, the number of synsets to which a word belongs is an estimate of the number of senses it has. In other words, the greater a word's synset count, the more polysemous and hence less semantically specific it can be taken to be. Synset count was manually extracted from the WordNet index file for nouns, for all values of N1 and N2 in the data. In cases where a constituent did not appear in WordNet, the compound was excluded from further analysis.

### 3.3.3 Distribution of Constituents

The positional constituent family sizes of N1 and N2 were found to be highly correlated with their respective synset counts, so that meaningful statistical models could not include both these types of predictor. In order to reduce the collinearity in the data to an acceptable level (c-number = 20.46), the so-called family size ratios were used instead of the positional constituent family sizes. The family size ratio of a compound constituent is the logarithm of its family size in the target position divided by its family size in the reverse position (Bell 2013). To take *house boat* as an example, the N1 positional family would include e.g. *house party, house wine, house cat* etc., while the reverse N1 family would include e.g. *acid house, town house, fashion house* etc.; for N2, the positional family would include e.g. *speed boat, river boat, ferry boat* etc., and the reverse family would consist of e.g. *boat house, boat hook, boat people* etc. The greater the family size ratio for a noun in a particular position, the stronger is the tendency of that

noun to occur in that position rather than the other position. Tarasova (2013: 103) finds that many nouns do show a preference for either the N1 (modifier) or N2 (head) position in compounds, and that the family size in either position is significantly inversely correlated with the family size in the other position. In other words, the more productive a noun is as a compound head, the less productive it is likely to be as a modifier, and vice versa. This is confirmed by the findings of Baayen (2010), in whose data about 32% of compound constituents occurred exclusively as N2 while about 41% occurred exclusively as N1.

   Positional and reverse family sizes were estimated from the whole BNC, using the BYU-BNC interface (Davies 2004-), by searching for spaced noun-noun strings in which the relevant values of N1 or N2 occurred in first or second position as appropriate. For those constituents whose singular and plural forms were tagged in the corpus as separate lemmas, the positional family of each form was extracted from the corpus. The families of the singular and plural lemmas were then combined, and any duplicated types were removed, to give the total number of combinations for the constituent in question. However, since not all noun-noun strings are compounds, care was taken to ensure that the raw counts accurately reflected actual family sizes.



**Figure 1.** Correlation between estimated and actual constituent positional family sizes in the spoken BNC, showing outliers

To do this, accurate family sizes were first calculated for a representative sample of the constituents in the data. For each of these constituents, every noun-noun string in which it occurred was checked in context in the spoken BNC, to ascertain whether it actually occurred as a compound. The logarithms of the accurate family sizes thus obtained were then plotted against the logarithms of the estimated family sizes, based on string counts, as shown in Figure 1. The constituents represented in bold in Figure 1 are outliers that are more than two standard deviations below the line of best fit. These are *bitch, chemist, fridge, mother, present, lead, shoe, tea, age,* and *morning*. These outliers were found to fall into a small number of categories, as described in Bell & Plag 2012:

> … the constituents were vocatives in N2 position (as in the example *tea mother*…); had homonymic forms (e.g. *lead*); were part of high-frequency formulae such as *morning meaning good morning*; had a high probability of being mistagged, e.g. the post-nominal adjective *present*, as in *highlights any special feature present* ; or had very small family sizes. (Bell & Plag 2012: 18)

For all constituents in the whole dataset that potentially fell into any of these five classes, accurate family sizes were then manually extracted from the corpus as described above. Again, the actual family sizes were plotted against the estimated family sizes, revealing a strong and highly significant positive correlation between the two measurements ($r=0.97$, $p<0.0001$). Compounds for which either constituent fell more than 1.5 standard deviations below the line of best fit were excluded from subsequent analysis. For the remaining compounds, the evidence indicates that the estimated family sizes based on raw counts of noun-noun strings are highly correlated with actual family sizes, and can therefore safely be used to represent family sizes in statistical models.

### 3.3.4 Concreteness of Constituents

For perceived concreteness, I used the publically available ratings obtained by Brysbaert, Warriner & Kuperman (2014). These authors obtained

concreteness ratings for 40 thousand English lemmas using online crowdsourcing. The raters indicated level of concreteness on a five-point scale, where 'concrete' meant that the meaning of the word could be learnt through direct sensory experience, as opposed to abstract meanings, which have to be learnt by explanation using other words.

### 3.3.5 Lexicalisation

The degree of compound lexicalisation was estimated using the spelling ratio of the compounds in the BNC (cf. Bell & Plag 2012, 2013). Spelling ratio is the logarithm of the number of times the compound occurs with unspaced or hyphenated spelling divided by the number of times it occurs with spaced spelling. The measure is based on the assumption that writers are more likely to produce a compound with hyphenated or unspaced spelling when they perceive it as representing a single concept, and that such singularity is most likely when the compound is lexicalised.
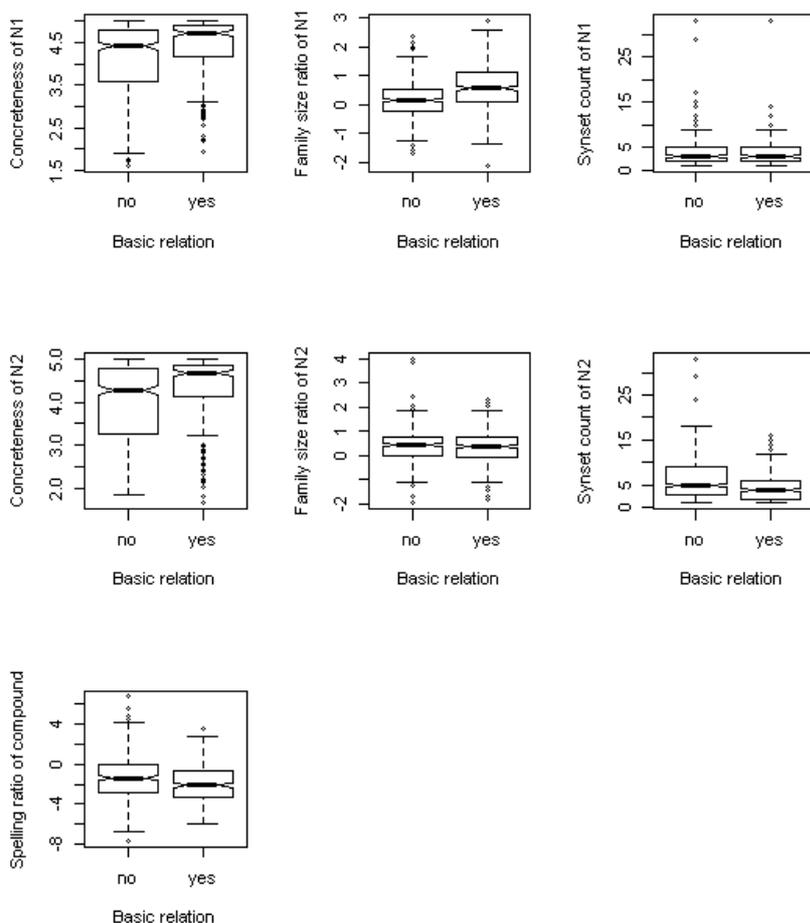
## 3.4 Statistical Analysis

The data was first explored graphically to examine the distributions of the predictor variables in compounds with and without basic relations. Following this, a logistic regression model was constructed in which the dependent variable was the probability of a compound having a basic relation and the predictors were the variables listed in Table 4. The model was constructed using the lrm function which is part of the rms package in the statistical software environment R. The initial model included all the predictors shown in Table 4, as well as a number of two- and three-way interaction terms. To allow for the possibility that different properties of a constituent might interact with one another, I initially included interactions between the three numerical predictors for each constituent, that is to say, between the concreteness rating, family size ratio and synset count. Since it is known, e.g. from the work of Gagné & Spalding (2014), that properties of the two constituents interact in the relational interpretation of compounds, I also included interaction terms between N1 and N2 variables, that is to say between N1 and N2 concreteness, N1 and N2 family size ratio

and between N1 and N2 synset count. Finally, because I hypothesised that the degree of lexicalisation of a compound might interact with any of the other properties, I included interaction terms between spelling ratio and all other predictors. Non-statistically-significant predictors were removed from the model step-wise in the standard process of model simplification. The final model resulting from this simplification process was then subjected to model criticism, using a combination of bootstrap validation (with 200 iterations) and penalised maximum likelihood estimation. These processes reduce the risk that the model will over-fit the data and reduce any undue influence on the model of any extreme values in the data.

## 4. Results

### 4.1 Distribution of Variables

Figure 2 shows the distributions of the numerical predictors for compounds that involve a basic relation compared with those that do not. The thick line in each box represents the median value and the extent of the box itself represents the interquartile range. The notches correspond to the confidence interval around the median; this means that if the notches on two boxes do not overlap, the medians can be regarded as significantly different. It can be seen that there are significant differences between compounds with and without basic relations in terms of the concreteness of both N1 and N2, as well as the spelling ratio of the compound, the family size ratio of N1 and the synset count of N2. All these differences go in the expected directions. The constituents of compounds with basic relations are perceived as significantly more concrete than the constituents of compounds without basic relations; compounds with basic relations have lower spelling ratios, i.e. are less lexicalised, than compounds without basic relations; the modifiers (N1) in compounds with basic relations have higher family size ratios, i.e. are more typically modifiers, than the modifiers in other compounds; and the heads (N2) of compounds with basic relations have lower synset counts, i.e. are more semantically specific, than the heads of other compounds.

**Figure 2.** Distribution of numerical predictors in compounds with and without basic relations.
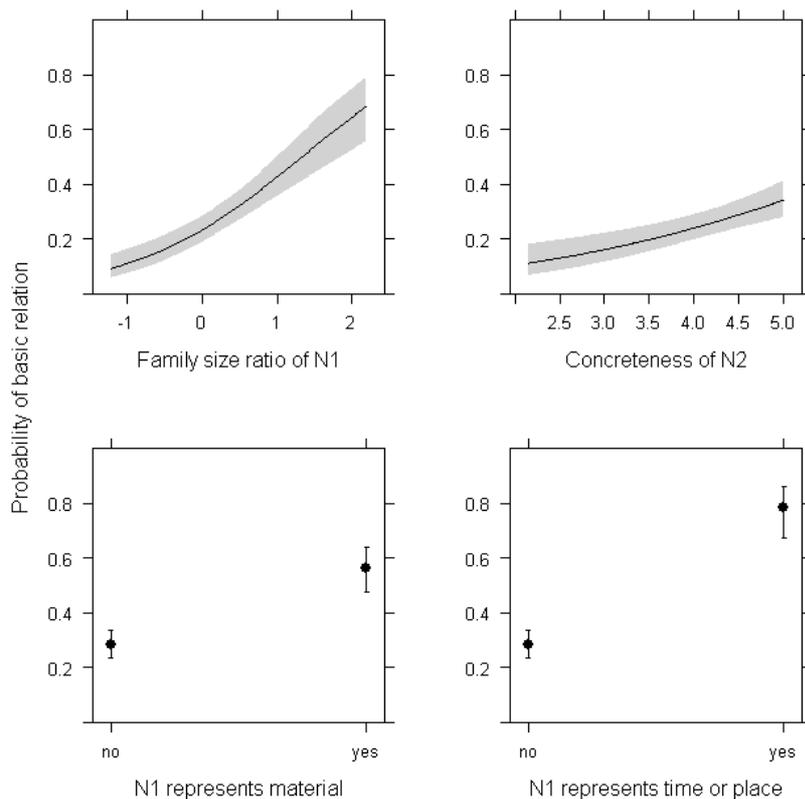
## 4.2 Logistic Regression Model

Table 6 shows the final logistic regression model, after step-wise removal of non-significant predictors, bootstrap validation and penalised maximum likelihood estimation. A positive coefficient indicates that an increase in the

predictor is associated with an increase in the probability of a compound having a basic relation, while a negative co-efficient indicates that an increased value of the predictor reduces the probability of the compound having a basic relation. The model provides support for all the predictions made in Section 2.4, with the exception of Prediction 2.

**Table 6.** Final model for basic relations, N = 788, C = 0.822

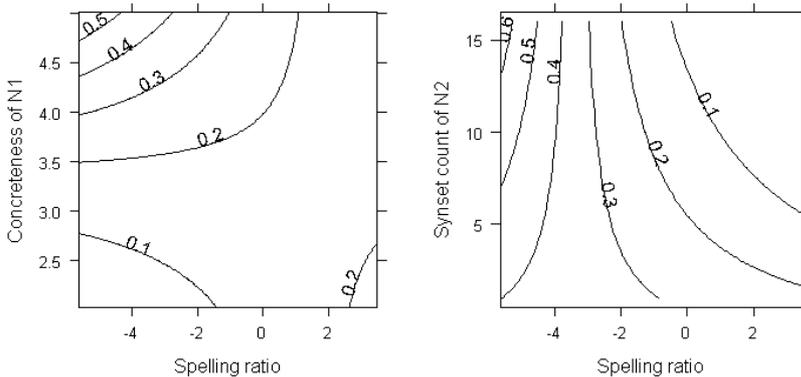|  | Coef | S.E. | Wald Z | Pr(>\|Z\|) | Penalty Scale |
|---|---|---|---|---|---|
| Intercept | -4.5184 | 0.9924 | -4.55 | <0.0001 | 0 |
| FamilySizeRatioN1 | 0.8917 | 0.1357 | 6.57 | <0.0001 | 0.4175 |
| N2Concreteness | 0.4968 | 0.1211 | 4.10 | <0.0001 | 0.4722 |
| N1isTimePlace=yes | 2.2119 | 0.2901 | 7.62 | <0.0001 | 0.4183 |
| N1isMaterial=yes | 1.1785 | 0.2119 | 5.56 | <0.0001 | 0.4183 |
| SpellingRatio | 0.6297 | 0.3144 | 2.00 | 0.0452 | 1.0678 |
| N1Concreteness | 0.2710 | 0.1916 | 1.41 | 0.1573 | 0.4690 |
| SynsetCountN2 | -0.1023 | 0.0357 | -2.87 | 0.0041 | 2.5329 |
| SpellingRatio*N1Concreteness | -0.1525 | 0.0691 | -2.20 | 0.0275 | 4.6248 |
| SpellingRatio*SynsetCountN2 | -0.0301 | 0.0132 | -2.29 | 0.0221 | 8.8942 |

The main effects in the model are shown in Figure 3. It can be seen that an increase in the family size ratio of the first constituent is associated with a rise in the probability of the compound having a basic relation. This is in keeping with Prediction 3 in Section 2.4, namely that the probability of a compound having a basic relation will be positively correlated with the productivity of N1 as a modifier. Nouns with high family size ratios in N1 position typically behave as modifiers, and are more productive as modifiers than they are as heads. Similarly in line with our hypotheses, and specifically with Prediction 1 in Section 2.4, it can be seen that as the concreteness of N2 increases, so too does the probability that a compound will have a basic relation. Since N2 is the semantic head in these compounds, this indicates that basic relations are more likely in compounds that represent entities that can be perceived through one or more of the senses, as opposed to those that need to be explained in terms of other concepts. Also as expected, and in accordance with Prediction 4, basic relations are more likely when N1 represents a material or a time or place.

**Figure 3.** Main effects in the final model for basic relations

The remaining predictors enter into significant interactions, in keeping with Prediction 5, and these are shown graphically in Figure 4. In this figure, the lines on the graphs are like contour lines on a topographic graph, with the figures on the lines representing the probability of a compound having a basic relation. The left-hand graph shows the interaction between spelling ratio and the concreteness of N1 in predicting the type of relation. It can be seen that basic relations are most likely when spelling ratio is low and N1 concreteness is high. When the spelling ratio is low, that is to say when the compound is not lexicalised, there is a strong effect of N1

concreteness, such that the likelihood of a basic relation increases with increased concreteness of N1, in line with our first prediction. However, when the spelling ratio is high, that is to say when the compound is highly lexicalised, the effect of N1 concreteness becomes negligible. The right-hand graph shows the interaction between spelling ratio and synset count of N2. It can be seen that, as expected, the probability of a basic relation generally falls with increasing spelling ratio, i.e. with increasing lexicalisation. However, this effect is most marked when the synset count of N2 is high, i.e. when N2 is most highly polysemous. Overall, basic relations are least likely when a compound is highly lexicalised and has a highly polysemous head.



**Figure 4.** Interaction effects in the final model for basic relations

## 5. Discussion and Conclusion

Overall, the results show that the likelihood of any compound involving one of Fanselow's so-called basic relations can be successfully predicted on the basis of properties of the compound's constituents as well as the degree to which the compound as a whole is lexicalised. What does this tell us? Firstly, the fact that the type of semantic relation can be predicted on the basis of properties of the constituents supports models of compound interpretation in which the relation is attached to or arises from the concepts associated with the constituent nouns. In the RICE theory of conceptual

combination (Spalding, Gagné, Mullaly & Ji 2010), for example, the modifier noun is thought to suggest possible relations which the head noun either accommodates or rejects. Gagné & Spalding (2014) report evidence that recent use of a relation in connection with a particular modifier concept increases the availability of that relation in subsequent processing, but they do not find evidence of relation priming in cases where the target modifier is neither identical nor closely semantically related to the prime. This suggests that semantic relations are conceptually connected only to particular lexical entries. On the other hand, the finding reported here, that relation type can be predicted not only on the basis of individual constituents but also on the basis of more general distributional properties and, in at least the cases of material and location nouns, in terms of semantic classes, suggests that compound relations are likely to have representations at higher levels of generalisation than individual lexical items.

The fact that groups of semantic relations can be predicted on the basis of general categories supports the finding by Maguire, Wisniewski & Storms (2010) that relational preferences of particular constituents might be based on semantic categories rather than applying only to individual words. In fact, the results of the present study suggest that the statistically significant categories in predicting compound relations may be even more general than those proposed by Maguire, Wisniewski & Storms (ibid.). In other words, not only the semantic category of the constituents might be relevant but also their more general distributional patterns, i.e. typically modifier or typically head, and more general semantic type, i.e. concrete as opposed to abstract.

In view of the correlation between compound semantic relations and constituent semantic categories, it is perhaps surprising that relational priming effects have not been more convincingly demonstrated in the absence of repeated constituents. Estes (2003) and Estes & Jones (2006) claimed to have demonstrated purely relational priming in the absence of a repeated constituent, but Gagné & Spalding (2014) argue that these effects were confounded by semantic similarity of the constituents used in the experimental items in the first case and by skewing of the relation frequencies in the second case. Nevertheless, Gagné & Spalding (ibid.) do report relational priming in cases where the modifier concepts of the prime and target are not identical, but semantically similar. This suggests

that there is some activation of the relational concept during processing and that this can spread to other modifiers, but the fact that the effect is not found in more distantly-related modifiers, suggests that it has what might be regarded as a short half-life in terms of spreading activation and perhaps diminishes more rapidly than other semantic activation. If particular concepts have strong preferences for particular relations, as shown by Maguire, Wisniewski & Storms (2010), it could well be that any priming effect due to recent exposure is very small compared with the background preference for a particular relation and that such effects are therefore very hard to detect experimentally, at least with the techniques currently available. An alternative would be to account for the priming effects found between closely related items by assuming some sort of binary distinction between closely-related concepts and not-closely-related concepts, but this seems intuitively implausible.

What of the question as to whether certain compound relations are more basic than others? We have seen that those relations classified by Fanselow (1981) as basic can to a significant extent be predicted on the basis of properties of the compound constituents. So-called basic relations are more likely when the constituents are more concrete, and when the compound itself is not semantically lexicalised. So if we regard concreteness as more basic than abstractness and compositional meaning as more basic than lexicalised meaning, then these relations may indeed be more basic than others.

But why should basic relations be associated with compounds in which N1 typically occurs in that position? One possibility is that the connection is related to productivity, since these nouns will tend to be productive modifiers, and certainly more productive as modifiers than they are as heads. A number of studies have established a correlation between productivity in word formation processes and the semantic and phonological transparency of the forms produced (see Plag 2006 for an overview). In terms of phonology, right-stressed English compounds might be regarded as more transparent than left-stressed forms, since so-called 'right stress' actually consists of two pitch accents, one on each constituent (Kunter 2011), and right-stressed compounds therefore clearly consist of two phonological words. If this phonological transparency is associated with the productivity

of the modifiers, then we might expect from the literature that right-stressed compounds will also be more semantically transparent than left-stressed types. And since right-stressed types are associated with basic relations, we might hypothesise that basic relations are more semantically transparent than stereotype relations; in fact, 'basic' in this context might mean 'more transparent', but this is an idea which requires further investigation.

If the basic relations do indeed involve the most productive modifiers, this might account for Fanselow's characterisation of them as offering fall-back interpretations for compounds where a more specific relation is not available, and even for his choice of the term 'basic'. The most basic characteristics of any class might be defined as those which apply most generally to the members of that class. For example, growth is a more basic property of living things than is sexual reproduction, since all living things grow, but not all living things reproduce sexually. By definition, productive modifiers can be used to modify a wide range of concepts and might therefore be regarded as representing more general characteristics than less productive modifiers. If some semantic relations are indeed associated with more productive modifiers, then these relations may also be perceived as representing more generalisable and hence more basic properties of the head nouns.

If the basic relations are basic in the sense of being more transparent or compositional, this might help explain why lexicalisation interacts with semantic specificity of the head as illustrated in Figure 4. If a noun has only one sense then it will retain that sense even in lexicalised compounds. With highly polysemous nouns, on the other hand, it is likely that lexicalised compounds will use less common and possibly unique senses of the noun in question. If the degree of compound compositionality reflects the extent to which the meaning of a compound can be deduced from the meanings of its constituents in isolation, and if polysemous constituents are likely to have different senses in isolation from their senses in compounds, then it follows that compounds involving highly polysemous nouns are likely to be less compositional than those including nouns with fewer senses. Hence, if basic relations are associated with compositionality, we would predict that they will also be associated with less polysemous heads.

However, it is not completely clear why this particular set of relations

constitutes the 'basic' group, or whether it is the most appropriate set. It could be that designating a different set of relations as 'basic' would lead to an even more successful model. I have hypothesised that the five relations designated as basic by Fanselow (1981) might be united by reference to states as opposed to actions, or by reference to natural as opposed to man-made entities. But just as there is significant variation in how individual compounds are understood, which makes their semantic classification so difficult, so too there is variation in how the relations themselves can be perceived. For example, I have argued that the MADE OF relation results from the laws of physics, since every physical entity has to be composed of some substance, but MADE OF could also be taken to refer to human activity, since artefacts are created out of substances. Further research could investigate whether changing the set of basic relations improves the predictive power of the models, and whether similar patterns can be found in languages other than English. A logical place to start would be with German, since that was the language on which Fanselow (1981) based his original ideas. It might also be fruitful to investigate the possibility that basic and stereotype relations do not represent a binary distinction, but that different relations fall along a cline, e.g. in terms of constituent concreteness.

Another area where further investigation might be fruitful is the relationship between particular semantic classes of modifiers and heads and particular semantic relations. Although it seems intuitively obvious that e.g. material modifiers will be associated with the MADE OF relation, the salient question is not so much which patterns seem plausible but rather which patterns actually occur and with what relative frequency. The models in this study only included two classes of modifier; in contrast, Maguire, Wisniewski & Storms (2010) and Maguire, Maguire & Cater (2010) provided evidence for a much wider set of associations between classes of head, modifier and semantic relation, and it would be interesting to include more of these classes in statistical models.

In summary, this paper has provided some initial empirical support for the notion that compound semantic relations can be regarded as falling into two main classes, but the details of this division, including for example what factors underpin it, whether it is gradient or categorical, and how it relates

to other possible taxonomies, remain to be elucidated.

## References

Baayen, R. H. 2010. The directed compound graph of English. In S. Olsen (Ed.), *New impulses in word-formation (Linguistische Berichte Sonderheft 17)*, 383-402. Hamburg: Buske.

Bauer, L. 1979. On the need for pragmatics in the study of nominal compounding. *Journal of Pragmatics* 3, 45-50.

Bell, M. J. 2013. *The English noun noun construct: its prosody and structure*. Cambridge: University of Cambridge PhD thesis.

Bell, M. J., & Plag, I. 2012. Informativeness is a determinant of compound stress in English. *Journal of Linguistics* 48(3), 485-520.

Bell, M. J., & Plag, I. 2013. Informativity and analogy in English compound stress. *Word Structure* 6(2), 129-155.

Brysbaert, M., Warriner, A., & Kuperman, V. 2014. Concreteness ratings for 40 thousand generally known English word lemmas. *Behavior Research Methods* 46(3), 904-911.

Davies, M. 2004-. BYU-BNC. (Based on the British National Corpus from Oxford University Press). Available online at http://corpus.byu.edu/bnc/.

Downing, P. 1977. On the creation and use of English compound nouns. *Language* 53, 810-42.

Estes, Z. 2003. Attributive and relational processes in nominal combination. *Journal of Memory and Language* 48(2), 304-319.

Estes, Z., & Jones, L. 2006. Priming via relational similarity: A copper horse is faster when seen through a glass eye. *Journal of Memory and Language* 55(1) 89-101.

Fanselow, G. 1981. Zur Syntax und Semantik der Nominalkomposition. Ein Versuch praktischer Anwendung der Montague-Grammatik auf die Wortbildung im Deutschen, *Linguistische Arbeiten* 107. Tübingen: Niemeyer.

Fudge, E. 1984. *English word-stress*. London: George Allen & Unwin.

Gagné, C. L., & Spalding, T. L. 2014. Conceptual composition: The role of relational competition in the comprehension of modifier-noun phrases and noun-noun compounds. *Psychology of Learning and Motivation* 59, 97-130.

Girju, R., Moldovan, D., Tatu, M., & Antohe, D. 2005. On the semantics of noun compounds. *Computer Speech and Language* 19, 479-96.

Hatcher, A. G. 1960. An introduction to the analysis of English noun compounds. *Word* 16, 356-73.

Jespersen, O. 1942. A modern English grammar on historical principles, part VI: *Morphology*. London: George Allen & Unwin Ltd.

Kunter, G. 2011. Compound stress in English: The phonetics and phonology of prosodic prominence, *Linguistische Arbeiten* 539. Berlin/New York: Walter de Gruyter.

Ladd, D. R. 1996. Intonational phonology, *Cambridge Studies in Linguistics* 79. Cambridge: Cambridge University Press.

Langacker, R. W. 1987. *Foundations of cognitive grammar*. Stanford: Stanford University Press.

Lauer, M. 1995. *Designing statistical language learners: Experiments on compound nouns*. Sydney: Macquarie University PhD thesis.

Lees, R. B. 1960. *The grammar of English nominalizations*. The Hague: Mouton.

Levi, J. N. 1978. *The syntax and semantics of complex nominals*. New York: Academic Press.

Liberman, M., & Sproat, R. 1992. The stress and structure of modified noun phrases in English. In I. A. Sag & A. Szabolcsi (Eds.), *Lexical matters* (CSLI Lecture Note 24), 131-181. Stanford: CLSI Publications.

Lieber, R. 2004. *Morphology and lexical semantics*. Cambridge: Cambridge University Press.

Lieber, R. 2009. A lexical semantic approach to compounding. In R. Lieber & P. Stekauer (Eds.), *The Oxford handbook of compounding*, 78-104. Oxford: Oxford University Press.

Marchand, H. 1969. T*he categories and types of present-day English word-formation: A synchronic-diachronic approach*, 2nd edn. Munich: C.H. Beck'sche Verlagsbuchhandlung.

Maguire, P., Maguire, R., & Cater, A. W. 2010. The influence of interactional semantic patterns on the interpretation of noun–noun compounds. *Journal of Experimental Psychology: Learning, Memory, and Cognition* 36(2), 288.

Maguire, P., Wisniewski, E. J., & Storms, G. 2010. A corpus study of semantic patterns in compounding. *Corpus Linguistics and Linguistic Theory* 6(1), 49-73.

Miller, G. A. 1995. WordNet: A Lexical Database for English. *Communications of the ACM* 38(11), 39-41.

Nastase, V., & Szpakowicz, S. 2003. Exploring noun-modifier semantic relations. *Proceedings of the 5th International Workshop on Computational Semantics* (IWCS-5), Tilburg, The Netherlands.

Ó Séaghdha, D. 2007. Designing and evaluating a semantic annotation scheme for compound nouns. *Proceedings of the 4th Corpus Linguistics Conference,* Birmingham, UK.

Ó Séaghdha, D. 2008. Learning compound noun semantics, *Technical Report* 735.

Cambridge: Computer Laboratory, University of Cambridge.

Ostendorf, M., Price, P. & Shattuck-Hufnagel, S. 1996. *Boston University Radio Speech Corpus*. Philadelphia: Linguistic data consortium, University of Pennsylvania.

*Oxford Advanced Learner's Dictionary* 1995. Oxford: Oxford University Press.

Plag, I. 1999. Morphological productivity: structural constraints in English derivation, *Topics in English Linguistics* 28. Berlin: Mouton de Gruyter.

Plag, I. 2006. Productivity. In B. Aarts & A. McMahon (Eds.), *The handbook of English linguistics,* 537-556. Oxford: Blackwell.

Plag, I. 2010. Compound stress assignment by analogy: The constituent family bias. *Zeitschrift für Sprachwissenschaft* 29(2), 243-282.

Plag, I., Kunter, G., Lappe, S., & Braun, M. 2007. Testing hypotheses about compound stress assignment in English: a corpus based investigation. *Corpus Linguistics and Linguistic Theory* 3(2), 199-233.

Plag, I., Kunter, G., Lappe, S., & Braun, M. 2008. The role of semantics, argument structure, and lexicalisation in compound stress assignment in English. *Language* 84(4), 760-794.

Pustejovsky, J. 1991. The generative lexicon. *Computational Linguistics* 17, 409-41.

Spalding, T. L., Gagné, C. L., Mullaly, A. C., & Ji, H. 2010. Relation-based interpretations of noun-noun phrases: a new theoretical approach. In S. Olson (Ed.), *New impulses in word-formation (Linguistische Berichte Sonderheft 17)*, 283–315. Hamburg: Buske.

Sweet, H. 1891. *A new English grammar logical and historical part 1: Introduction, phonology and accidence*. Oxford: The Clarendon Press.

Tarasova, E. 2013. *Some new insights into the semantics of English N+N compounds*. Wellington: Victoria University of Wellington PhD thesis.

Warren, B. 1978. *Semantic patterns of noun-noun compounds*. Göteborg: Acta Universitatis Gothoburgensis.