

# Between syntax and morphology: German noun-verb units as reluctant compounds

Roland Schäfer  
*Deutsche Sprache und Linguistik,  
Humboldt-Universität zu Berlin*  
Dorotheenstraße 24, 10117 Berlin  
roland.schaefer@hu-berlin.de

Ulrike Sayatz  
*Deutsche und niederl. Philologie,  
Freie Universität Berlin*  
Habelschwerdter Allee 45, 14195 Berlin  
ulrike.sayatz@fu-berlin.de

**Abstract** ...

**Keywords:** universion, prototypes, corpus data, experiments, German

## 1 Introduction

Over the past one to two decades, investigations into the probabilistic nature of grammar have gained attention from a growing number of linguists. The core theoretical question underlying such investigations is whether grammar itself is a probabilistic phenomenon (probabilism attributed to competence), or whether indeterminacies in speakers' output as well as their acceptability judgements can be attributed to problems accessing an otherwise binary non-probabilistic grammar (probabilism attributed to performance). Both inter-speaker and intra-speaker variability (such as alternations between case forms, lexical near-synonyms, or syntactic constructions, etc.) could be caused by either of these mechanisms, and it is our goal to contribute to this discussion by arguing that a certain alternation in written Modern German is better compatible with a probabilistic competence model. While many studies compare the competition between two or more lexical options, morphological constructions, or syntactic constructions, the alternation discussed here is one between a syntactic and a morphological construction.

The alternation we are going to explore affects a construction containing a verb and a noun, and it alternates between a syntactic manifestation

References re probabilism to be added later. Take from habil.

(where the noun combines with the verb via a syntactic relation) and a morphological one (where the noun is incorporated into the verb), and we argue that alternations in spelling provide evidence for the grammatical status of the instances of the construction. In this construction, a noun that is either a direct object of the verb (normally in the accusative case) or an adjunct to the verb which would normally take the form of a prepositional phrase occur in strictly bare form (the adjuncts occurring even without the preposition) and acquire an unspecific generic reading. While singular indefinite mass nouns typically occur without an article in German, this is the only frequent construction in German wherein bare count nouns occur. In the examples in (1), *Rad fahren* ('bike riding') refers to the concept of riding any bike, and the unspecific reading of *Rad* is obligatory, which is not the case for the English translations with the indefinite article.

- (1) a. Remy **fährt** gerade **Rad**.  
 Remy rides<sub>PRES</sub> right now bike  
 Remy is riding a bike right now.
- b. Yael weiß, dass Remy **Rad fährt**.  
 Yael knows that Remy bike rides<sub>PRES</sub>  
 Yael knows that Remy is riding a bike.
- c. Remy ist gestern **Rad gefahren**.  
 Remy is yesterday bike ridden<sub>PART</sub>  
 Remy rode a bike yesterday.
- d. Remy hat keine Lust, **Rad zu fahren**.  
 Remy has no motivation bike to ride<sub>INF</sub>  
 Remy doesn't feel like riding a bike.
- e. Remy ist am **Rad fahren**.  
 Remy is at the bike ride<sub>INF/NOUN</sub>  
 Remy is riding a bike.
- f. Remy singt beim **Rad fahren**.  
 Remy sings upon the bike ride<sub>INF/NOUN</sub>  
 Remy is singing while riding a bike.
- g. \* Remy lobt das **Rad fahren**.  
 Remy praises the bike riding<sub>NOUN</sub>  
 Remy praises the riding of bikes.

Such N + V units occur flexibly in all types of syntactic contexts: with finite verbs in verb-second order (1a), with finite verbs in verb-last order (1b), in the analytical perfect where the lexical verb takes the form of a participle (1c), in infinitives with the particle *zu* (1d), in a progressive-like

Best reference?

And here, too: Best reference?

construction with the preposition *an* fused with the dative article to *am* where the infinitive is potentially nominalised (1e), and in regular prepositional phrases (1f). In (1g), the spelling of the N + V unit as two words is impossible, hence the asterisk. In this case, we can assume that the noun and a fully nominalised infinitive form a regular nominal compound.<sup>1</sup> The spelling as two words for (1e) and (1f) is not accepted by all native speakers, a fact to which will return throughout the paper.

Good reference for the footnote. Maybe English?

In the examples (1c) through (1g), the noun and the verb occur in sequence with no intervening material. In these cases, the noun and the verb alternate between the spelling as multiple words seen in (1) and spellings as one word shown in (2).

- (2) c. Remy ist gestern **radgefahren**.
- d. Remy hat keine Lust, **radzufahren**.
- e. Remy ist am **Radfahren/radfahren**.
- f. Remy singt beim **Radfahren/radfahren**.
- g. Remy lobt das **Radfahren**.

In (2e) and (2f), additional variation is introduced in the form of upper-case and lower-case initials.<sup>2</sup> The compound with the nominalised infinitive in (2g) is fine if spelled as one word.

English reference for the footnote?

We call cases where a multi-stem unit is spelled as two words such as in (1) ‘disjunct spellings’ and cases where a unit is spelled as one word as in (2) ‘compound spelling’. We see that N + V units potentially undergo graphemic *univerbation* in the form of compound spelling. Lehmann (2021: 2) calls univerbation “the union of two syntagmatically adjacent word forms in one”. We follow this terminology and assume univerbation to be the directly observable phenomenon, i. e., compound spelling of adjacent words that could potentially (or historically) also be used in disjunct spelling. Historically and—as we’re going to show especially in Section 4—individually, univerbation is a gradual process, and it is thus a probabilistic phenomenon. However, univerbation per se is not necessarily the result of a regular grammatical pattern.<sup>3</sup> Thus, a major aim of this paper is to show whether and how the univerbation of N + V units in German is based on an established morphological construction wherein a noun is incorporated into a verb, forming a new verb expressing a new event concept.

It was unclear which aspects the three references pertain to. Hence, I put them in the footnote. Please fix/make more explicit.

Reintroduce some good examples of random words that happened to undergo univerbation at some point.

<sup>1</sup> Infinitives in German are routinely nominalised as an action noun.

<sup>2</sup> In German, all nouns are capitalised anywhere in a sentence.

<sup>3</sup> See also Gallmann (1999: 294), Jacobs (2005: 107), Lehmann (2021: 4).

We will argue that such a morphological construction exists, but that the alternative syntactic construction remains available to speakers because N + V units have properties of both morphological as well as syntactic prototypes. This double nature and the resulting alternation is interpreted as favouring probabilistic competence models over deterministic competence models. In Section 2, we lay the theoretical and descriptive foundations. Section 2.1 introduces probabilistic grammar and how graphemic evidence can be interpreted in probabilistic grammar. In Section 2.2, the nature and status of spaces and their loss (univerbation) are discussed briefly, and Section 2.3 provides a detailed account of N + V units in German. We finish Section 2.3 by summing up our hypotheses before presenting a large-scale corpus study and an elicitation experiment in Sections 3 and 4, both representing tests our particular hypotheses about N + V units and the overarching hypothesis regarding the probabilistic nature of grammar. We conclude with a summary, further interpretation and discussion in Section 5.

## 2 Theoretical background

### 2.1 Probabilistic grammar and graphemic evidence

### 2.2 Spaces, words, and univerbation

In this paper, we use graphemic evidence—both from corpora and from controlled experiments—and argue that it allows us to draw conclusions about writers’ cognitive grammars. More specifically, we assume that compound spellings of N + V units indicate that writers conceive of those units as single syntactic words, whereas disjunct spelling indicates that they conceive of the unit as two syntactic words. Therefore, we briefly introduce the status of the space in German writing and how it pertains to N + V units.

German writing uses an alphabetic script with a strong correlation between underlying phonological forms (the phonemic level) and characters (graphemes). A common fundamental principle of such scripts is the separation of syntactic words by spaces (Jacobs 2005: 22). Also, stems and affixes are never separated from one another in German, which reinforces the status of the space as a demarcation of syntactic words.<sup>4</sup> These factors

<sup>4</sup> There is a class of verbal particles which does not follow this principle. Verbs like *aufessen* (‘eat up’) formed from a verb stem (*essen*) and a prefixed particle (*auf*) are spelled as one word when they are adjacent in verb-last order, but they are separated in verb-second

A brief 1–2 page overview of prototype as applied to grammatical categories, including the role of contexts in alternation modelling.

Mention item-specific effects.

Introduce usage-based graphemics.

Fix MuellerXYZ in footnote.

facilitate the reader's ability to decode the sequence of syntactic words, and they constitute a crucial principle in the encoding and conventionalisation of meanings associated with word forms (Jacobs 2005: 22).

Unlike in English, compound spelling of syntactic words comprising more than one stem was also established in the history of German writing, especially for the case of the highly productive noun + noun (N + N) compound pattern (Fuhrhop 2007: 182, Jacobs 2005: 34, Section 2.3 below), for which compound spelling is the dominant graphemic realisation. However, there is a heterogeneous group of multi-word constructions for which only a tendency towards compound spelling can be observed (Szczepaniak 2009: 95, Wurzel 1998: 335). As opposed to N + N compounds, these constructions typically consist of words with different parts of speech, such as *mithilfe (von)* ('with the help (of)') from *mit der Hilfe (von)* or *zu Hause* ('at home') from *zu Hause*.<sup>5</sup> For such cases, Lehmann (2021: 2) posits a "downgrading of a syntactic to a morphological boundary" between the two words. When writers use compound spelling in these cases, they choose to encode the construction as a single word with a morphological boundary instead of a sequence of words with a syntactic boundary. If many speakers consistently make this choice over a significant period of time, the unit might become conventionalised as a single lexical word or—in other words—lexicalised (Lehmann 2021: 7). Until such a diachronic process is complete and one of the spellings has become clearly dominant, conventionalisation does not provide a very strong input to writers, and they alternate between a syntagmatic and a morphological realisation. For many of these constructions, this is the case both in non-standard as well as standard written German, albeit assumedly to different degrees.<sup>6</sup>

N + V units with different affinities towards compound spelling like *Rad fahren* ('bike riding', often also spelled *radfahren*) and *eislaufen* ('ice skating', infrequently also spelled *Eis laufen*) represent different levels of diachronic re-conventionalisation as single words.<sup>7</sup> This indeterminacy means that

Reference needed.

order where the verb is moved to sentence-second position and the particle remains in sentence-last position through obligatory long-distance movement **MuellerXYZ**.

<sup>5</sup> Normative approaches as well as individuals display a lot of variation with respect to at least some of those constructions (cf. below).

<sup>6</sup> We are not aware of any published research comparing the alternation tendencies in standard and non-standard written German.

<sup>7</sup> The orthographic norm is notoriously unstable with respect to N + V units, which contributes to their unclear status. Before the significant reform of the orthographic norm in 1996, both *radfahren* and *eislaufen* were supposed to be spelled as one word. After the reform, both units were supposed to be written as two words (*Eis laufen* and *Rad fahren*). After a revision of the reform in 2006, *eislaufen* was again supposed to be spelled as one

speakers have both the syntagmatic realisation (disjunct spelling) and the morphological realisation (compound spelling) in their graphemic input, which subsequently leaves them with quite a free choice to be made based on how a concrete token is classified according to their individual grammar. It is the task of usage-based probabilistic graphemics (Schäfer & Sayatz 2016) to uncover factors influencing such decisions and decode the principles at work in speakers' internal grammar by analysing their writing habits.

## 2.3 The status of noun-verb units in German

In Section 1, we showed that N + V units alternate between compound spelling and disjunct spelling when they occur in sequence. In this section, we explain why the existence of this alternation is not surprising considering the morphosyntactic system of German. Furthermore, we argue that in each concrete case where an N + V is written, the strength of the tendency towards either compound or disjunct spelling can be derived from the overall syntactic and morphological patterns available in present-day German. These patterns are shown to have prototypical properties which are matched by individual N + V units and their syntactic contexts more or less well, which leads to either compound or disjunct spelling being the preferred realisation. The hypotheses put forward here are then tested in Sections 3 and 4.

In order to achieve this end, we need to shed some light on the productive N + N compound construction in Section 2.3.1 before turning to N + V units as reluctant compounds in Section 2.3.2. We sum up our arguments and derive our hypotheses for the empirical studies in Section 2.3.3.

### 2.3.1 N+N compounds

For an N + V unit to undergo graphemic univerbation (i. e., , a downgrading of a syntactic to a morphological construction in the sense of Lehmann 2021: 2) systematically, it must resemble established morphological constructions enough to be classified as an instance of that construction itself.<sup>8</sup> The prototypical and by far most productive morphological construction combining more than one stem is noun + noun compounding (N + N), to

---

word, whereas *Auto fahren* was supposed to be spelled as two words exclusively (Eisenberg 2013: 327).

<sup>8</sup> Random isolated univerbations like *zu Hause* 'at home' from *zu Hause* are not systematic in this sense. They are merely the result of idiosyncratic diachronic developments.

which we turn now in some detail.<sup>9</sup> German N + N compounds instantiate a morphological construction proper and are therefore inseparable. Syntactically, nothing can intervene in between the two stems of the compound, and they cannot be rearranged. With minor exceptions (often exaggerated in normative discussions), they are also inseparable graphemically, i. e., they are always written as one word. Furthermore, they are always head-final, mostly determinative, and they allow recursive formation wherein an N + N compound enters into another N + N compound, resulting in [[N + N] + N] or [N + [N + N]] structures (Fleischer & Barz 2012: 13). Some examples are given in (3) and (4) for *Haustür* and *Haustürschlüssel*, the latter being recursively formed from the former.<sup>10</sup>

Reference needed for the footnote which is not Eisenberg.

Reference needed, including maybe a reference talking about the normative discussions.

- (3) Haus.tür  
house.door  
front door
- (4) Haus.tür.schlüssel  
[[house.door].key]  
key to the front door

The semantic relation between the first noun (N<sub>1</sub>) and the second noun (N<sub>2</sub>) is highly unspecific, rendering many compounds semantically ambiguous unless they are strongly lexicalised (Klos 2011: 252).<sup>11</sup> The historic development of the stable N + N compound construction was furthered during the Early High German period (approximately from the 14th to the 17th century AD) by a syntactic change. The dominant pattern of noun-noun attribution had been a prenominal genitive as in the now obsolete (5), which swiftly changed to a postnominal genitive as in (6).

Check years/centuries.

- (5) † des Hauses Tür  
the<sub>Gen</sub> house<sub>Gen</sub> door  
the door of the house

<sup>9</sup> Adjectives also enter compounds as the head, such as in *feuerrot* 'red like fire', literally 'fire red'. However, they are much less frequent than N + N compounds, and we consequently don't discuss them here.

<sup>10</sup> If necessary, we present compound spelling with a minimal analysis of the morphological structure. Affixes are separated from stems by hyphens, an lexical stems are separated from each other by a period. Within compounds containing more than two stems, structure is using using square brackets as in examples (7) and (8) below.

<sup>11</sup> Obviously, once they are strongly lexicalised, they cannot help to establish a more canonical type of semantic relation between N<sub>1</sub> and N<sub>2</sub>, either, simply because lexicalised compounds are often intransparent to the language user (Klos 2011: 59), such as *Kammerjäger* ('pest controller', literally *chamber hunter*).



- (6) die Tür des Hauses  
 the door the<sub>Gen</sub> house<sub>Gen</sub>  
 the door of the house

To the extent that prenominal attribution in syntax became more and more obsolete, the prenominal position was used to establish the highly productive morphological construction of N + N compounds as in (3) (see Nübling et al. 2017: 132, Schlücker 2012), which showed a tendency to be written in compound spelling very early on. The N + N compound construction is semantically at least as unspecific as the syntactic genitive construction to which it is diachronically related (see also Eisenberg 2013: 220). Its recursive application is virtually unrestricted (Wurzel 1994: 504). N<sub>1</sub> and N<sub>2</sub> are just concatenated as bare stems in most cases, but there are also so-called linking elements, which are sometimes positioned in between the stems.<sup>12</sup> Diachronically, linking elements stem from diverse sources, but the overall pattern of inserting them is related to the former morphological marking in prenominal genitives.

N + N compounds as described in this section are clearly the prototype for morphological constructions combining more than one stem in German. In the next section, we show how N + V units deviate from this prototype, and how this leads to them alternating between a syntactic and a morphological construction.

### 2.3.2 N+V units as reluctant compounds

In this section, we argue that N + V units are *reluctant compounds*. While in principle the morphological N + V construction as a kind of compound written as one word has existed for centuries, we show why and how it remains in competition with a syntactic construction. At the same time, we argue why—at least under the right circumstances—morphological compounding (and consequently the spelling as one word) become the preferred realisation.

As opposed to compounding with nominal heads (as discussed in the previous section), compounding with verbal heads is not a highly productive pattern in German.<sup>13</sup> A major difference compared to N + N compounds is the fact that N + V units are usually not inseparable as was already shown in Section 1. There can be intervening material in between the noun and the

<sup>12</sup> A recent large-scale study (Schäfer & Pankratz 2018: 339) showed that 60% of all N + N compound types have no linking element, whereas 40% do.

<sup>13</sup> Eisenberg (2013: 224) finds that there are roughly 400 lexicalised N + V compounds.

Reference needed.

Reference needed.

Examples needed for the footnote. Are we talking exclusively about our type of N + V, or does this include other types of N + V?



verb in some contexts (the infinitival particle *zu* and the inflectional affix *ge-* of past participles). Additionally, the noun and the verb can be separated and rearranged as in verb-second order, where N + V units resemble particle verbs (Fortmann 2015). This fact alone means that N + V units do not fit the compounding prototype well. This likely introduces great resistance in speakers to classify them as compounds and consequently use compound spelling.

Pages missing for Fortmann. Also, please provide full reference for Fortmann (2015).

Maybe add details from Ulrike's page 3.

Another major difference between N + N compounds and N + V units is that the morphological N + V construction is not recursive. Nominalised N + V units marginally occur as N<sub>1</sub> in N + N compounds (contrary to claims by Fuhrhop 2007: 54) as in (7).<sup>14</sup> However, an N + V unit cannot function as the verbal head in another N + V unit (i. e., a [N + [N + V]] structure) as illustrated in (8). Native speakers will readily acknowledge that such constructions are outright absurd.

- (7) a. *Energie.spar.messe*  
       [[energy.save].fair]  
       trade fair for products useful in saving energy
- b. *Endlager.such.gesetz*  
       [[final storage.search].law]  
       law about the search for a permanent repository for nuclear waste
- c. *Feuer.lösch.boot*  
       [[fire.extinguish].boat]  
       fire-fighting boat
- (8) \* *Rad.fahr.mach-en*  
       [[bike.ride].make-INF]  
       make bike riding

We posit that the lack of core properties of prototypical (productive) German compounding constructions (separability, potential reordering, lack of recursive application) is a major factor in keeping the formation of N + V units from establishing a fully productive morphological compounding construction, thus keeping it from reliably requiring graphemic univerbation.

Another noticeable difference between the N + N and the N + V construction is the specificity of the internal relation. While the relation in N + N

Any references that specifically go with this paragraph?

<sup>14</sup> The examples in (7) are attested and taken from the DECOW16B web corpus (see Section ??). Their document frequencies are 218 for *Energiesparmesse*, 416 for *Endlagersuchgesetz*, and 414 for *Feuerlöschboot* in a corpus of 17.1 million documents. The document frequency is the number of documents the lemma occurs in, not counting multiple occurrences within each document.

compounds is vastly underspecified (see Section 2.3.1), there are only two possible relations within N + V units, and these relations are determined by—and above all decodable through—the distributional properties of the verb (including its argument structure). It's either an object relation or an adjunct relation where all distributional restrictions apply that would apply in a syntactic realisation of the same verb. As a consequence, there is always a syntactic paraphrase for N + V units with an adjunct relation where the noun occurs in a prepositional phrase which is an adjunct to the verb.<sup>15</sup> See (9) for an example.

- (9) a. Kim will die Corvette probefahren.  
 Kim wants the Corvette test.drive  
 Kim wants to test-drive the Corvette.
- b. Kim will die Corvette zur Probe fahren.  
 Kim wants the Corvette to the test drive.  
 Kim wants to test-drive the Corvette.

The relation is decodable except in rare cases which underwent full lexicalisation a long time ago such that the meanings of the lexemes or their distributions have changed significantly. However, the decodable relations (direct object or prepositional adjunct) are prototypically realised syntactically as German is a language with a very weak (if any) tendency towards noun incorporation (see below in this section). The arguments of a verb as well as its adjuncts are usually realised as syntactic dependants of the verb. Even if the verb is nominalised, direct objects are realised as genitives in the noun phrase, and adjuncts remain the same with nominalised verbs.

The fact that the relation can be decoded for almost all N + V units means that the morphological construction marked by graphemic univerbation almost always remains in competition with a syntactic construction with distinct syntactic words separated by spaces in writing. This competition between a morphological construction and a syntactic construction was pointed out with varying terminology by—among others—Fleischer & Barz (2012: 12), Schlücker (2012: 13), and Morcinek (2012: 88). Whereas the parallel syntactic construction for N + N compounds (prenominal genitives) disappeared within a relatively short period of time, the ambiguity between syntax and morphology of N + V units remains intact. This is true although univerbation of N + V units with an object relation dates back to

<sup>15</sup> Pragmatically, these paraphrases might often be subject to blocking because of the availability of the N + V construction. However, this does not make them syntactically or semantically unacceptable.

Maybe a reference here?

On your p. 5, you mention homonymous morphological and syntactic constructions, not giving examples. I fail to see where homonymous constructions exist. It seems to be the point that in N + V units, N loses its genuine syntactic properties (article, modifiability, ...). The constructions are therefore never homonymous.

Middle High German (*lobpreisen* ‘praise’, literally ‘compliment praise’) and even Old High German (*hals-werfōn* ‘turning around’, literally ‘neck turning’), see Wurzel (1994: 517), Wurzel (1998: 334). For N + V units with an adjunct relation, Morcinek (2012: 89) notices that dictionaries from between 1750 and 1993 AD list novel N + V units with an adjunct relation with increasing frequency. For centuries or even more than a millennium, N + V units have been co-existing in syntax and morphology. We assume that the stable availability of an alternative syntactic realisation is yet another major factor in preventing N + V formation from becoming a more clearly morphological construction in language users’ cognitive grammars, making N + V units *reluctant* compounds.

However, we still have to ask why and under which conditions true compounding and graphemic univerbation might be preferred. As mentioned in Section 1, the morphological construction for N + V units is a type of noun incorporation. N + V units are usually seen as the only cases of potential incorporation in Modern German (Eisenberg 2013: 224), which is why we postulated above that object and prepositional adjunct relations are prototypically realised syntactically. According to Mithun (1984: 848), incorporation is “a particular type of a compounding in which a V and N combine to form a new V”.<sup>16</sup> As Mithun (1984: 848–849) points out, incorporation happens when the verb denotes a new and independent event concept in combination with the incorporated noun the semantics of which are determined by the previous syntactic relation between the noun and the verb. Typically, the noun loses its referential autonomy as well as its specificity, and it acquires a generic reading, which is indeed the case for N + V units. In sentences like (10), no specific bike is referenced, and *radfahren* refers to the whole concept of riding any bike. This is true for both compound and disjunct spelling.

Don't forget to mention it in the introduction.

(10) Friedel kann radfahren/Rad fahren.

Friedel can bike.ride

Friedel knows how to ride a bike.

As a result of the semantic degradation of the noun, it loses its modifiability (also regardless of spelling), as illustrated in (11).

(11) \* Friedel kann schnelles Rad fahren.

Friedel can quick bike ride

Friedel knows how to ride a quick bike.

<sup>16</sup> From Mithun’s types of noun incorporation, German N + V units clearly represent type 1 *lexical compounding*. Since nothing could be gained from it, we do not discuss the literature on the typological classification of incorporation further.

Such losses of referential autonomy and syntactic combinatorics are referred to as ‘noun stripping’ by Gallmann (1999: 287). The loss of specificity and referential autonomy as well as the acquisition of a generic reading are part of the semantics of the N + V construction (see also Gallmann 1999: 287, Bredel & Günther 2000: 108, Eisenberg 2013: 325). Functionally speaking, the construction exists in order to express the new event concept which requires the generic/unspecific reading of the noun. Thus, the noun has the properties typical of nouns that are subject to incorporation of the lexical compounding type. Hence, N + V units have a tendency to form proper compounds and subsequently undergo univerbation despite the factors mentioned above that pull them towards a syntactic construction. This is why we call them reluctant compounds.

In the next section, we will summarise the factors that influence the tendency of N + V units to incorporate and undergo univerbation. We also formulate testable hypotheses for the empirical work to be reported in Sections 3 and 4.

### 2.3.3 Hypotheses for the empirical studies

What favors univerbation? Greater similarity to N + N prototype and reduced competition from syntactic realisation. Hence:

- (1) Nominal contexts. Say what they are.
- (2) Reduced competition from syntax through adjunct relation between N and V. Also diachronic results from Morcinek (2013:??). The adverbial relation would need more explicit marking via preposition in the syntactic realisation. Also, there is the productive and functionally similar type of government compounds in the N + N world. However, there is no such pattern for adjunct relations.
- (3) Semantically weak/generic verbs like haben, machen, fahren, ...because they only denote a specific concept together with the noun.
- (4) Idiosyncratic diachronic status. Each N + V unit may have undergone conventionaisation/lexicalisation to different degrees. Tendencies are therefore likely to vary across individual N + V units.

The individual diachronic tendencies as reflected in Modern German are measured in the corpus (= degree of conventionalisation/lexicalisation) and confirmed in the experiment (= individual grammar).

Mention nominal/verbal contexts.

Mention linking elements!

Reuse from corpus part: The goals of the corpus study was (i) to assess which V + N units are used in written German, (ii) to corroborate that morphosyntactic contexts, internal relations, and linking elements influence the probability of univerbation, and finally (iii) to show how strongly the individual N + V units are attracted by univerbation. The operationalisations relied on the fact that the major graphemic principles in German are clear and dominant, and that they are both deeply rooted in diachrony and well entrenched in writers' usage (Reference?). The relevant major principle for the present study was compound spelling of syntactic words, and we interpreted compound spelling as a direct indication of univerbation in the writers' grammars.

We only look at cases where the alternation occurs, i.e., , contact position.

### 3 Analysing the usage of noun-verb units

In this section, we apply both exploratory and confirmatory methods of analysing the univerbation of N + V units using corpus data. We motivate our choice of corpus and describe the sampling and annotation procedure in Section 3.1.

We perform exploratory analysis using association measures in Section 3.2 in order to gauge the individual tendencies of N + V units to incorporate and undergo univerbation in written language usage. The tendencies calculated here will also be used as a control variable in the experiment reported in Section 4.

Finally, the results of estimating the parameters of a multilevel model explaining the variation in the univerbation of N + V units are reported in Section 3.3.

#### 3.1 Choice of corpus, sampling, and annotation

As a first step, we adopted a data-driven approach in order to find close to all N + V units in contemporary written usage. In a second step, we had to count their occurrences in compound and separate spelling in the relevant morphosyntactic contexts enumerated in Section 2.3.3: fully nominalised as the heads of noun phrases, in *am* progressives, as participles in analytical verb forms, and as infinitives in a range of verbal constructions (for example with modal verbs).

Clearly, a large corpus with rich morphological and morphosyntactic annotations containing texts written in a broad variety of registers and styles (including ones written under low normative pressure) was required. We chose the DECOW16B corpus (Schäfer & Bildhauer 2012) because it fulfils all the aforementioned criteria.<sup>17</sup> Much like the SketchEngine corpora (Kilgarriff et al. 2014), the COW corpora contain web documents from recent years. However, the German DECOW (containing 20.5 billion tokens in 808 million sentences and 17.1 million documents) offers a much wider range of annotations compared to SketchEngine corpora, including morphological annotations and several levels of syntactic annotation (dependencies and topological parses). For our purpose, the fully internal analysis of nominal compounds described in Schäfer & Pankratz (2018) was particularly of interest. It allows for searches of roots within nominal compounds. For example, we could query compounds with a deverbal head such as *Zeitnehmen* ('time

<sup>17</sup> <https://www.webcorpora.org>

taking’). Furthermore, the interface offered by the creators of the COW corpora allows for automated queries controlled by Python scripts using the open-source SeaCOW interface.<sup>18</sup> The scripts we used to make the queries are released on a curated open-data server along with all data as well as the  $\LaTeX$ , knitr, and R scripts created in the writing of this paper.<sup>19</sup>

The list of actually occurring N + V units was obtained by querying for compounds with a nominal non-head and a deverbal head.<sup>20</sup> The rationale behind this approach is that any N + V unit of interest should occur at least once in compound spelling as a fully nominalised compound. Since this step relied on automatic annotation already available in the corpus, the results contained erroneous hits which we removed manually. The resulting list contained 819 N + V units.<sup>21</sup>

In the second step, we created lists of all relevant inflectional forms of the verb in each V + N unit and used these to query all possible compound and separate spellings (including variance in capitalisation) of each of the 819 N + V unit types. In total, 28,665 queries were executed to create the final data set used here, a number which clearly demonstrates the necessity of script-based corpus access in data-driven methods. The queries were matched by 958,118 compound spellings and 1,288,768 separate spellings, which results in a total sample size of 2,246,886 tokens.

For each N + V unit in the sample, the following variables were annotated automatically: (i) the verb lemma, (ii) the noun lemma, (iii) whether a linking element is used in the use as a full noun, (iv) the overall frequency in the corpus. Additionally, we manually coded all 819 N + V units for the relation holding between the verb and the noun. The codes used in clear-cut cases were *Object* (441 units) and *Adjunct* (286 units). For 92 units, both relations were conceivable, and those cases were coded as *Undetermined*. This class is illustrated by *Daumenlutschen* (“thumb sucking”), which could correspond to the paraphrase either in (12a) or in (12b).

- (12) a. [den Daumen]<sub>NP<sub>Acc</sub></sub> lutschen  
the thumb suck

<sup>18</sup> <https://github.com/rsling/seacow>

<sup>19</sup> The DOI of the data set will be revealed in the accepted version of this paper.

<sup>20</sup> See the scripts available under the abovementioned DOI for concrete queries and further details.

<sup>21</sup> Notice that three highly frequent N + V units were excluded because they could be considered outliers, having an overly strong tendency to be used in compound spelling. They are *Teilnehmen* (“take part”), *Maßnehmen* (“take measure”), and *Teilhaben* (“have part” = “participate”).

- b. [am Daumen]<sub>pp</sub> lutschen  
on the thumb suck

Glue paragraph missing.

### 3.2 Results 1: Association strengths

In this section, we report an analysis of the item-specific affinities of N + V units towards univertation. The method we use seems superficially similar to collocational analysis (Evert 2008 for an overview) or collostruc-tional analysis (Stefanowitsch & Gries 2003). However, there are major differences to be explained momentarily (see also Schäfer & Pankratz 2018; Schäfer n.d.).

The order of this section and the next section have been switched without adapting the text, yet. Do it!

Our goal was to quantify how strongly each N + V unit tends towards univertation vis-a-vis all other N + V units. Thus, we need to compare the counts of cases with and without univertation of the unit in question with the total counts for all other N + V units. Such comparisons must be made relative to the overall number of the specific N + V unit as well as the num-ber of all other units. The counts needed for each N + V unit are nicely summarised in a 2×2 contingency table as shown in Table 1.

|                     | Univertation | No univertation |
|---------------------|--------------|-----------------|
| Specific N+V unit   | $c_{11}$     | $c_{21}$        |
| All other N+V units | $c_{21}$     | $c_{22}$        |

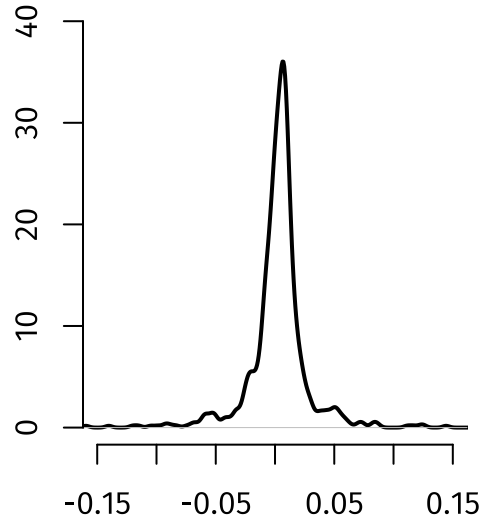
**Table 1:** 2×2 contingency table as used in the calculation of the strengths of the associations of N+V units with univertation.

With these counts, we are able to quantify how strongly the proportion in the first row differ from those in the second row, and there is a range of statistical measures for that. For example, one could use odds ratios or effects strengths from frequentist statistical tests.<sup>22</sup> We chose Cramér’s  $\nu$  derived from standard  $\chi^2$  scores ( $\nu = \sqrt{\chi^2/n}$ ). The  $\nu$  measure quantifies for each individual N + V unit how strongly its counts (cells  $c_{11}$  and  $c_{21}$ ) deviate from its counts that we would expect if there were no difference between this unit and all other N + V units (cells  $c_{21}$  and  $c_{22}$ ) with respect to their tendency to univertate. Since Cramér’s  $\nu$  is always in the range between 0

<sup>22</sup> p-values from frequentist statistical tests are measures of evidence, and therefore not appropriate in such situations (Schmid & Küchenhoff 2013; Küchenhoff & Schmid 2015) although they were used in early collostruc-tional analysis. However, even collostruc-tional analysis is now often used with measures of effect strength (Gries 2015).



and 1, it allows us to compare analyses where the samples differ. In itself,  $\nu$  does not tell us whether the deviation is negative (for a N + V unit with less than average compound spellings) or positive (for a N + V unit with more than average compound spellings). The information about the direction of the deviation is added by multiplying  $\nu$  with the sign of the upper left cell of the residual table of the  $\chi^2$  test.<sup>23</sup>



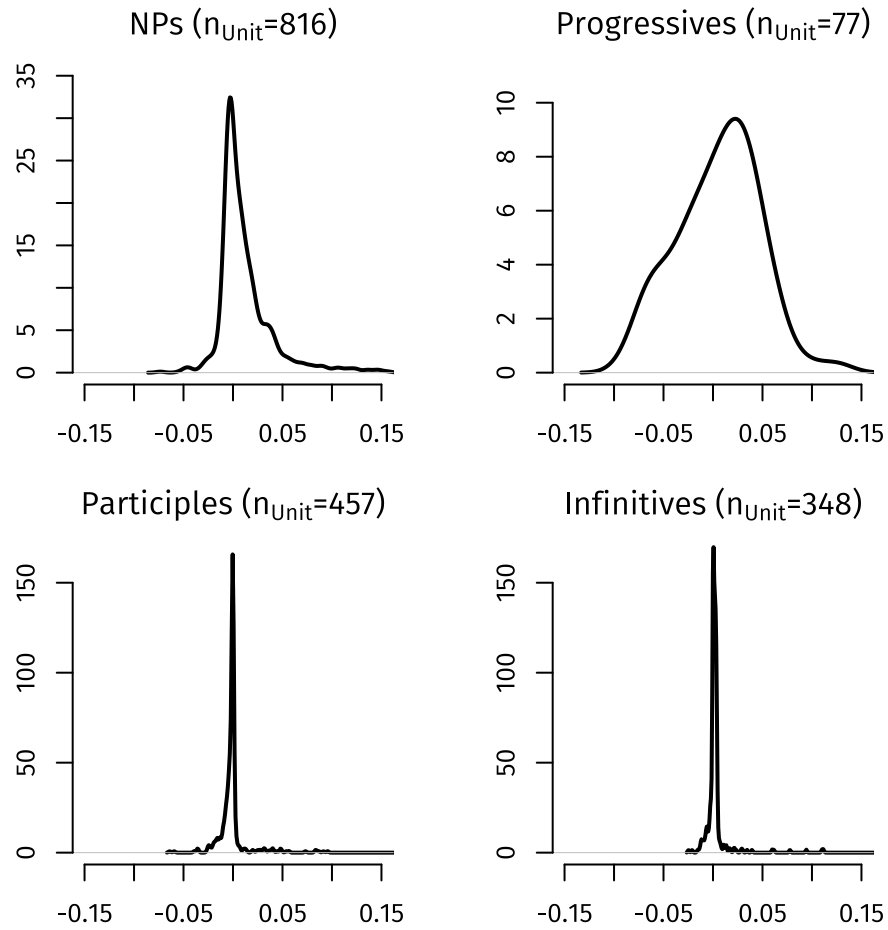
**Figure 1:** Density estimate of the distribution of the 819 association scores (across all morphosyntactic conditions).

We calculated the signed  $\nu$  for each of the 819 N + V units. Their distribution is plotted in the form of a density estimate in Figure 1.<sup>24</sup> Based on the annotations in the corpus data set, we can also compare the association strengths for specific morphosyntactic contexts. The counts as shown in Table 1 are simply reduced to the counts in each of the four contexts in turn. With the resulting lower sample sizes, the  $\chi^2$  measure can no longer be calculated for a number of infrequent N + V units, leading to lower  $n_{Unit}$  (= the number of N + V units analysed). The resulting distributions are shown in Figure 2.

The context-wise distributions of the association scores corroborate the results from the GLMM reported in Section 3.3. In the NP context (top

<sup>23</sup> The association scores encode almost the same information as the second-level model in the GLMM reported in Section 3.3, but they have a much more accessible interpretation.

<sup>24</sup> It approximates a scaled symmetric  $\chi^2$  distribution with  $df = 1$  squashed between -1 and 1.



**Figure 2:** Density estimates of the distribution of the association scores in the specific morphosyntactic conditions..

left panel of Figure 2), the right tail of the curve is much heavier than the left tail, which means there are mostly higher than usual tendencies towards univertation. In the morphosyntactically similar progressive (top right panel), the distribution is approximately symmetric, but given the low number of 77 N + V units for which  $\nu$  could be calculated, the result cannot be seen as stable.<sup>25</sup> Both prototypically verbal contexts (lower two panels) show heavier left tails, meaning that N + V units tend to resist univertation in these contexts. Once again, this is just another (and maybe more intuitive) look at the data in addition the GLMM analysis.

For the selection of stimuli in the experiment, the overall association strength (Figure 1) is relevant, because it truly represents the effect of the unit, independently of the context. The context effect will be controlled independently in the experiment. To illustrate how the data analysis allows for a selection of N + V units based on their affinity towards univertation, we show the top ten units with the highest negative and highest positive association in Table 2.

| V+N Unit        | Assoc. | Rel.    | V+N Unit       | Assoc. | Rel.   |
|-----------------|--------|---------|----------------|--------|--------|
| Radfahren       | 0.190  | N/D     | Gedankenmachen | -0.160 | Object |
| Computerspielen | 0.144  | Adjunct | Geldverdienen  | -0.140 | Object |
| Zeitreisen      | 0.125  | Adjunct | Rechtgeben     | -0.120 | Object |
| Skifahren       | 0.123  | Adjunct | Spaßhaben      | -0.115 | Object |
| Autofahren      | 0.117  | N/D     | Rechthaben     | -0.105 | Object |
| Probefahren     | 0.111  | Adjunct | Kinderhaben    | -0.099 | Object |
| Bogenschießen   | 0.087  | N/D     | Zeitnehmen     | -0.093 | Object |
| Schifffahren    | 0.085  | N/D     | Auftraggeben   | -0.092 | Object |
| Windsurfen      | 0.084  | Adjunct | Fehlermachen   | -0.088 | Object |
| Bergsteigen     | 0.082  | Adjunct | Urlaubmachen   | -0.083 | Object |

**Table 2:** Top ten V+N units with a strong tendency for univertation (left panel) and top ten V+N units with a strong tendency against univertation (right panel).

The tables illustrate that units with the strongest tendencies against univertation are predominantly ones with an object relation. The ones which most strongly favour univertation are mostly ones with an adjunct relation or an ambiguous relation. The ten items with the least clear tendency in

<sup>25</sup> The low number is one the one hand due to the fact that progressives are rare compared to NPs, participles, and infinitives. On the other hand, it is likely that many N + V units cannot be used in the progressive for semantic or pragmatic reasons.

either direction are shown in Table 3. They mostly come with an object relation.

| V+N Unit          | Assoc. | Rel.    |
|-------------------|--------|---------|
| Autowaschen       | 0.009  | Object  |
| Zigarettenrauchen | 0.007  | Object  |
| Haarewaschen      | 0.005  | Object  |
| Notenlesen        | 0.003  | Object  |
| Golfspielen       | 0.001  | Object  |
| Haarschneiden     | -0.007 | Object  |
| Wasserholen       | -0.008 | Object  |
| Feuermachen       | -0.009 | Object  |
| Blutabnehmen      | -0.009 | Object  |
| Schlangestehen    | -0.010 | Adjunct |

**Table 3:** Top ten V+N units without any tendency for or against univertation.

Among the units with an object relation, it is difficult to tell based on native-speaker intuition why the ones in Table 3 should have no preference and the ones in Table 2 (right panel) should resist univertation. This goes to show that, while we can model the tendencies to some extent using linguistic features, there are obvious item-specific effects which should be taken seriously from a theoretical perspective, and which must be accounted for in behavioural experiments. We now turn to such an experiment in Section 4.

Glue paragraph missing.

### 3.3 Results 2: Multilevel model

In this section, we present the parameter estimates (and predictions of conditional modes) for a binomial multilevel model (or generalised linear mixed model, GLMM) which models the relevant factors influencing writers' choice of the compound and the separate spelling.<sup>26</sup> Given the grand total of 2,246,886 observations in the sample (see Section ??), we will completely refrain from an interpretation in terms of inferential statistics. For samples of such magnitude in data-driven approaches, frequentist significance tests are the wrong tool. Therefore, we provide standard likelihood ratio confidence intervals for parameter estimates and prediction intervals for conditional modes as an approximate measure quantifying the precision

Mention that results between assoc. and GLM converge, but mention Arppe & Järvikivi 2007 etc.

References!

<sup>26</sup> See (Schäfer 2020, to appear) for an overview of the method and our philosophy in modelling.

of the parameter estimates and predictions. The models we specify reflect theoretically motivated decisions, and we therefore reject all types of model selection by means of step-up or step-down procedures.

As argued in Section `sec:thestatusofnounverbunitsingerman`, we expect the probability of the univerbation of N + V units to depend on the morphosyntactic context, the relation holding between the verb and the noun, the presence of absence of a linking element in the nominal compound (as a marker of a stronger lexicalisation) and on the specific N + V unit (a lexical tendency).<sup>27</sup> Accordingly, the response variable was chosen to be the proportion of compound spellings among all the spellings of the N + V unit. In the input data provided to the estimator, the response variable was thus a vector of 819 proportions, one for each N + V unit.<sup>28</sup> We specified four regressors. The only first-level (or observation-level) fixed effect regressor is the morphosyntactic context (a four-way categorical variable). As there is a huge number of 819 N + V units, the lexical indicator variable for the individual N + V unit should not be used as a fixed effect (Gelman & Hill 2006: 244–247, Schäfer 2020, to appear). We specified a generalised linear mixed model with the N + V unit variable as a random effect. The variables encoding the internal relation and the presence/absence of a linking element are nested inside the levels of the random effect, and they are therefore treated as second-level fixed effects in a multilevel model. In R notation, the specification is shown in (13).<sup>29</sup>

$$(13) \quad \text{Univerbation} \sim (1|\text{NVUnit}) + \text{Context} + \text{Relation} + \text{Link}$$

The estimated parameters of the model are given in Table 4. Additionally, effect plots for *Context* and *Relation* are given in Figure 3.<sup>30</sup> As

<sup>27</sup> Make sure linking elements were mentioned sufficiently prominently above!

<sup>28</sup> Binomial models can be specified in this manner (Zuur et al. 2009: 245–260). In the estimation of such models, the influence of each proportion is weighted according to the number of cases observed to calculate it. Without the weighting, highly frequent observed proportions would have too small an influence on the estimation, and infrequent ones would have an inappropriately high influence. In the case at hand, such a model on proportion data is also a convenient way of getting around the practical difficulties of estimating a model on the raw 2,246,886 observations.

<sup>29</sup> See Appendix A for a precise specification in mathematical notation.

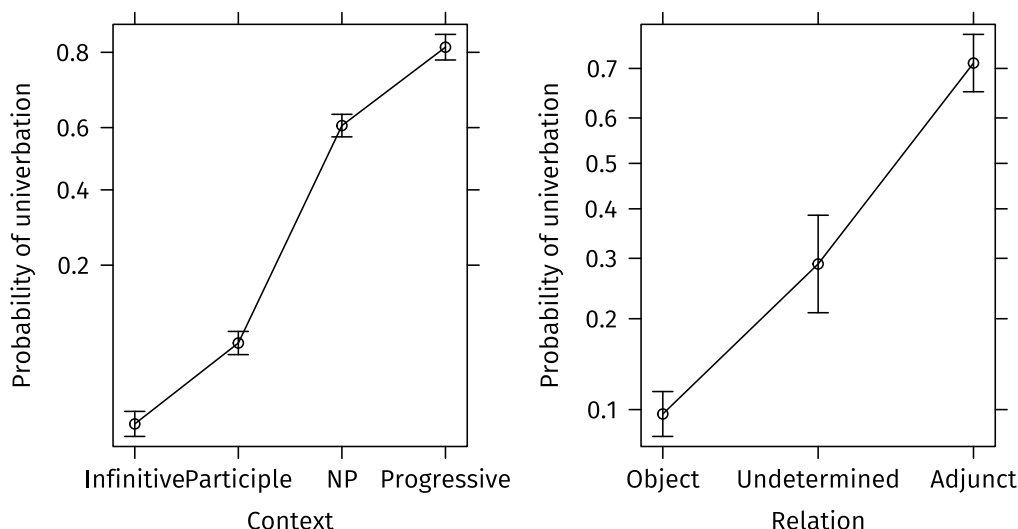
<sup>30</sup> Effect plots for binomial GLM(M)s (Fox & Weisberg 2018) plot the probability of the outcome across values of a regressor assuming default values for all other regressors. While model coefficients in binomial (and other) models have no direct interpretation in terms of probability, effect plots allow a more intuitive interpretation in terms of changes in probability.

|                         | Estimate | CI low | CI high |
|-------------------------|----------|--------|---------|
| (Intercept)             | -4.787   | 2.000  | 2.218   |
| Context = Participle    | 1.054    | -5.020 | -4.555  |
| Context = NP            | 3.886    | 0.976  | 1.133   |
| Context = Progressive   | 4.907    | 3.815  | 3.959   |
| Relation = Undetermined | 1.339    | 4.801  | 5.015   |
| Relation = Adjunct      | 3.132    | 0.862  | 1.816   |
| Link = Yes              | 0.361    | 2.808  | 3.456   |

**Table 4:** Coefficient table for the binomial GLMM modelling the corpus data with 95% profile likelihood ratio confidence intervals. The horizontal line separates first-level and second-level effects. Weighting was used to account for the bias in models on proportion data. Random effect for V+N lemma: Intercept = 4.430, sd = 2.105. The intercepts model the fixed effects Relation = Object and Link = No. Nakagawa & Schielzeth's  $R_m^2 = 0.577$  and  $R_c^2 = 0.999$ .

expected, the prototypically verbal contexts (infinitives and participles in analytic verb forms) are associated with a low probability of compound spelling (the infinitive is on the intercept, which is estimated at  $-4.787$ , and participles have a coefficient of  $1.054$ ). NPs and progressives as prototypically nominal contexts clearly favour compound spelling (coefficients of  $3.886$  and  $4.907$ , respectively). Both the coefficients and the effect plot (right panel in Figure 3) show a low probability of compound spelling when the relation between the verb and the noun (on the intercept) is that of an object, and a high probability when the relation is that of an adjunct (coefficient  $3.132$ ). The undetermined cases are in between the two clear-cut cases (coefficient  $1.339$ ). The presence of a linking element in fully nominalised compounds favours compound spelling only slightly (coefficient  $0.361$ ).

Given the narrow confidence intervals and the high marginal measure of determination  $R_m^2 = 0.577$ , we consider the hypotheses regarding fixed effects as well corroborated by the data, especially the effects of the context and the internal relation. Based on our commitment to a usage-based probabilistic view of language, we also predicted differences between N + V units not explainable by the fixed effects. These effects would show up as the residual variance in the random effects (in the form of the conditional modes) not modelled by the second-level effects. The conditional modes are centred around a second-level intercept of  $4.430$  with a standard devia-



**Figure 3:** Effect plots for the regressor encoding the morphosyntactic context of the N+V unit and the regressor encoding the syntactic relation within the N+V unit in the GLMM modelling the corpus data.

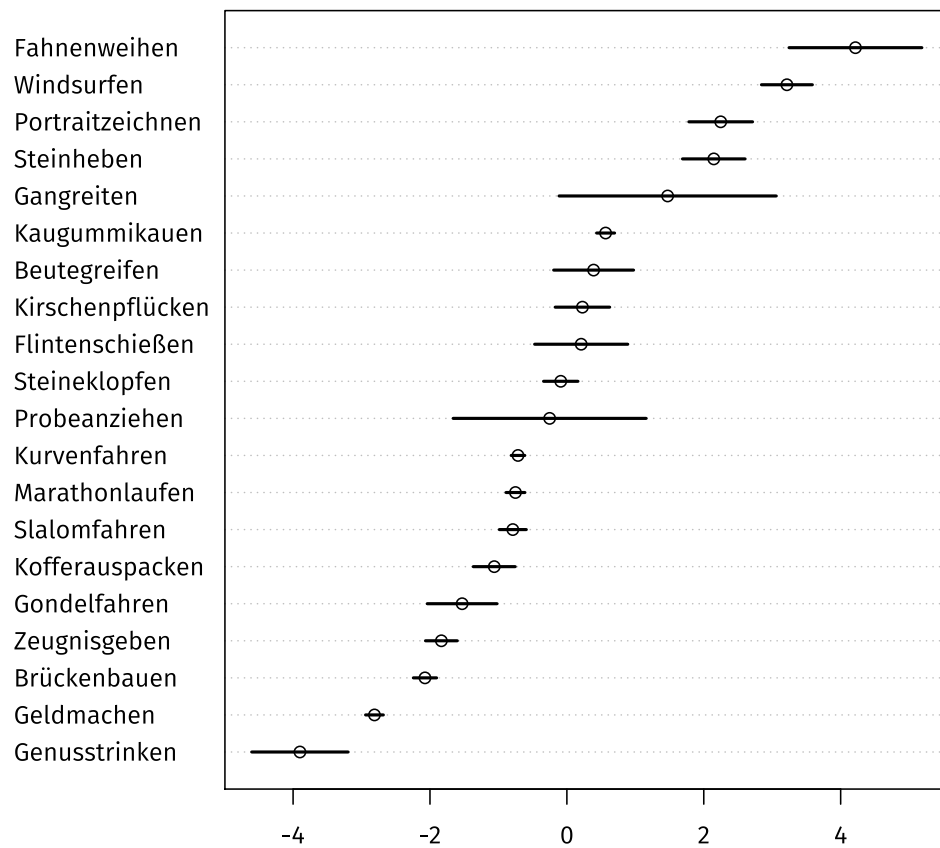
tion of 2.105. The relatively high standard deviation is a sign that there is considerable variation across the individual N + V units. Furthermore, the conditional  $R_c^2$  is as high as 0.999. This is commonly interpreted as saying that the fixed effects and the idiosyncratic effect of concrete N + V units almost fully explain the variance in the data. A random selection of 20 conditional modes, which illustrates the relevance of lexical idiosyncrasies through obvious differences with mostly very narrow prediction intervals, is shown in Figure 4.

The individual V + N unit thus plays a major role in writers' tendency to univerbate V + N units, which conforms the results from Section 3.2. The results obtained from that method will be used to predict participants' behaviour in the controlled experiment reported in Section 4.

## 4 Elicited production of noun-verb units

In this section, we corroborate the findings from Section 3 in a controlled experiment. We describe the rationale behind the experiment, the methods used, the design, and the group of participants in Section 4.1. Section 4.2





**Figure 4:** A random selection of conditional modes with 95% prediction intervals for the levels of the random effect in the GLMM modelling the corpus data.

reports the results descriptively and in the form of a generalised linear mixed model.

## 4.1 *Design and participants*

The goal of the experiment was to corroborate the findings from the corpus study and to test whether writers' behaviour under controlled experimental conditions is similar to the behaviour of writers under uncontrolled circumstances as found in corpora. We used pre-recorded auditory stimuli in order to elicit spellings of given N + V units. The stimuli were chosen based on theoretically motivated criteria and the information about item-specific tendencies obtained from the exploratory part of the corpus study in Section 3.2. We constructed eight sentences instantiating the four contexts. For each of the four contexts, we chose one N + V unit with a high and one with a low attraction strength according to the analysis of the usage data.<sup>31</sup> The sentences were constructed in a ways such that all N + V units were the predicate of a subordinate clause. This consistently ensured verb-last constituent order and avoided interfering verb-second effects, which are typical of independent sentences in German. The full sentences are given in Appendix C.

We added 32 fillers, resulting in a total of forty sentences being read to the participants.<sup>32</sup> Of the forty sentences, twenty (including the target items) had to be written down by the participants. The order of the target items was randomised, but it was made sure that there were at least three sentences in between pairs of items. There were nine distractors in the form of yes-no questions related to random sentences previously heard by the participants. An overview of the item design is shown in Table 5, where each line represents the features of one of the eight items.

In total, 61 participants took part in the experiment. All of them were first-semester students of German Language and Literature at Freie Universität Berlin. They were between 18 and 44 years old with a median age of 22 years. There were two separate groups (32 and 29 participants, respectively), and the randomisation of the stimuli was different between the two groups.

<sup>31</sup> Given the overall constraints on the choice of the items, “low” and “high” had to be interpreted as quite relative terms. However, we made sure that all low attractions scores are lower than zero and all high attraction scores are greater than zero. Also, for each context, the low and the high attraction score differ by at least 0.05.

<sup>32</sup> Of the forty fillers, six were actually items from an unrelated experiment.

| Context     | N+V unit      | Attraction | Binary |
|-------------|---------------|------------|--------|
| Infinitive  | Platzmachen   | -0.052     | Low    |
|             | Seilspringen  | 0.011      | High   |
| Participle  | Mutmachen     | -0.069     | Low    |
|             | Probehören    | 0.055      | High   |
| Progressive | Teetrinken    | -0.037     | Low    |
|             | Bogenschießen | 0.087      | High   |
| Clitic      | Spaßhaben     | -0.115     | Low    |
|             | Bergsteigen   | 0.082      | High   |

**Table 5:** Items from the experiment, chosen by context and attraction score.

## 4.2 Results

The distribution of responses of the experiment is shown in the form of a mosaic plot in Figure 5. It shows the number of compound spellings (univerbation) and separate spellings in each of the four contexts and for N + V units with high and low attraction score.

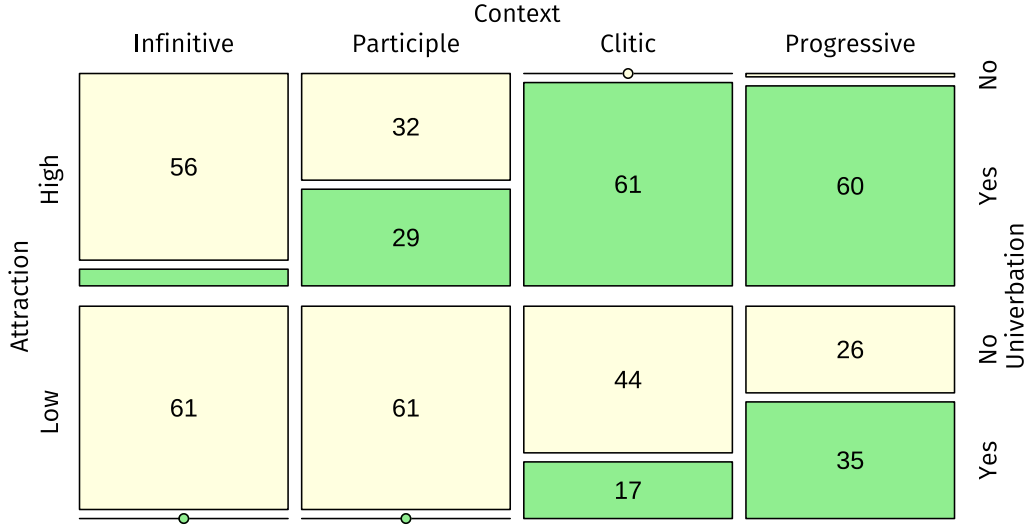
The overall number of positive responses (i. e., compound spellings) rises across the four contexts. It's 5 for the infinitive, 29 for the participle, 78 for NPs with a cliticised article, and finally 95 for the progressive (in each case out of 122). For the N + V units with a high attraction score, participants (almost) always use compound spelling in NPs with a cliticised article (61 out of 61) and in the progressive (60). Between the infinitive (5) and the participle (29), there is a clear differentiation in positive responses, however.

For the N + V units with a low attraction score, the items with an infinitive (0) or a participle (0) don't seem to allow univerbation at all. However, with NPs (17) and progressives (35), we see a considerable number of positive responses.

Clearly, both independent variables are highly useful in predicting the behaviour of participants. However, among the items with low association scores, we would expect the NPs as highly prototypical nominal contexts to trigger univerbation most strongly, while in the experiment they lose to the progressive (17 out of 61 for the NPs, 35 for the progressives).

Fix the nomenclature with "NPs" and "Clitic".

To examine the results further, we proceeded to estimate the parameters of a GLMM with additional control for individual participants in the form of a random effect. Instead of using a grouping variable for the N + V units,



**Figure 5:** Mosaic plot of the responses in the production experiment (vertical right) grouped by the morphosyntactic context (horizontal) and the binned N+V unit's attraction strength calculated from the corpus (vertical left).

we included their association strengths directly in the model. The model specification in R notation is given in (14). Appendix B provides the specification mathematical notation. The coefficient estimates for the GLMM are reported in Table 6.

$$(14) \quad \text{Univerbation} \sim (1|\text{Participant}) + \text{Attraction} + \text{Context}$$

There is some variation between writers as captured in the standard deviation of the conditional modes (1.722), but the small difference between the marginal  $R^2$  (0.803) and the conditional  $R^2$  (0.896) suggests that speaker variation does not explain much of the variance in the data. The coefficients indicate that the attraction strength derived from the corpus is positively correlated with the participants' tendency to univerbate (48.740). There seems to be no evidence that the participle has a different effect than the infinitive (which is on the intercept) given the large confidence interval  $([-0.375, 2.574])$ . On the other hand, progressives (6.224) and NPs with cliticised articles (8.166) clearly have a much more positive effect on the probability of univerbation. Thus, in the GLMM analysis NPs appear to have a stronger tendency to favour univerbation than progressives. This is

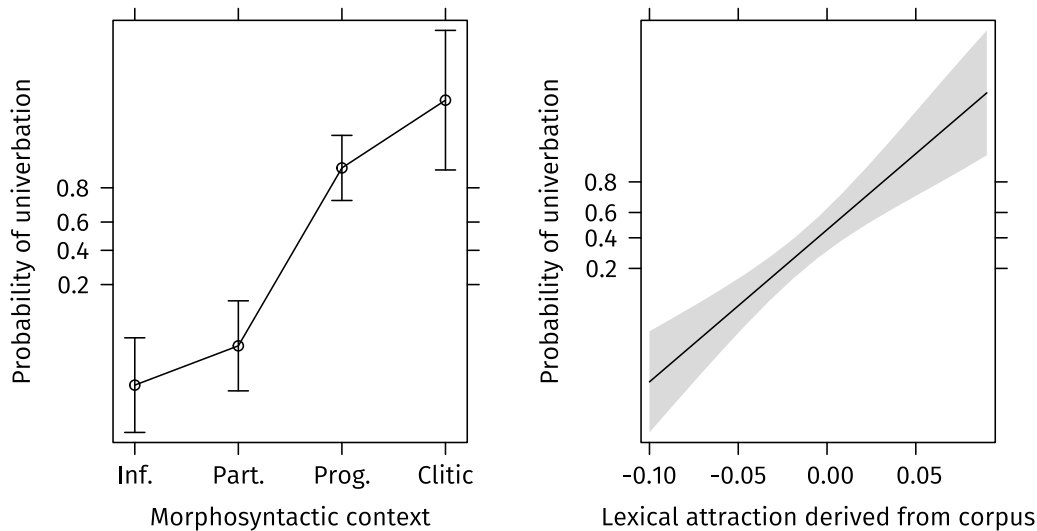
|                       | Estimate | CI low | CI high |
|-----------------------|----------|--------|---------|
| (Intercept)           | -4.026   | -5.534 | -2.851  |
| Attraction            | 48.740   | 34.657 | 73.476  |
| Context = Participle  | 1.126    | -0.375 | 2.574   |
| Context = Progressive | 6.224    | 4.691  | 8.184   |
| Context = Clitic      | 8.166    | 5.978  | 11.513  |

**Table 6:** Coefficient table for the GLMM modelling the experiment data with 95% confidence intervals. Nakagawa & Schielzeth's  $R_m^2 = 0.803$  and  $R_c^2 = 0.896$ . Random effect for participant: Intercept = 2.966, sd = 1.722 The intercept models the fixed effect Context = Infinitive as well as Attraction = 0.

in line with our theoretical predictions but seems to contradict the descriptive analysis (see Figure 5).

The effect plots in Figure 6 show the same picture as the coefficient table. The prototypically verbal contexts are associated with low probabilities of univertation, the two prototypically nominal ones with high probabilities of univertation. While progressives and NPs with clitics show the order predicted by theory, there is only very weak to no support for assuming a substantial difference judging by the large and mostly overlapping confidence intervals. The attraction scores are neatly correlated with the probability of univertation.

The apparent paradox with respect to the order of the effects of NP and progressive contexts that we see between Figure 5 on the one hand and Table 6 and Figure 6 on the other hand can be explained by looking at the concrete attraction strengths in Table 5. The unit with “low” attraction used in the progressive context (*Teetrinken*) has a numeric attraction score of  $-0.037$ , which is much closer to 0 than the one used in the NP context (*Spaßhaben*) with  $-0.115$ . At the same time, the high attraction counterparts are rather close to each other numerically ( $0.087$  for *Bogenschießen* and  $0.082$  for *Bergsteigen*). Figure 5 therefore shows a positive bias for the progressive context, which is likely due to the concrete choice of items for this experiment and the binary binning of the attraction scores into “low” and “high”. The GLMM compensates for this because we used the numerical attraction scores rather than just a binned “low” and “high” classification. As the selection of the few stimuli for a given experiment is virtually never possible with perfect control over all variables, the more advanced statistical analysis protects us against misinterpretation.



**Figure 6:** Effect plots for the regressor encoding the morphosyntactic context and the attraction strength as calculated from the corpus in the GLMM modelling the experimental data.

In sum, the experiment supports our theoretically motivated hypotheses, and it corroborates the results from the corpus study. We proceed to a final analysis of the phenomenon at in the light of our findings hand in Section 5.

## 5 Explaining noun-verb univerbation

### A Corpus study: full specification of the model

In Section 3.3, the specification of the model was given in R notation as (13), repeated here as (15).

$$(15) \quad \text{Univerbation} \sim (1|\text{NVUnit}) + \text{Context} + \text{Relation} + \text{Link}$$

This notation blurs the difference between first-level and second-level fixed effects. The model specification is the crucial step in statistical modelling since it encodes the researchers' commitment to a causal mechanism

Does not work with prepositional objects. \*vorurteilsgelitten

Not very productive with adjuncts. \*brückengestanden

Mention difference: object is related to argument structure of the verb, adjunct is not.

How productive is the N + V incorporation? Further studies required.

Es geht ja um die „Umdeutung/Modifizierung“ eines verbalen Ereigniskonzepts, bei dem entweder eine Objektrelation oder eine Adjunktrelation zu einer Veränderung der lexikalischen Markierungsstruktur des – wie auch immer – entstehenden komplexen Verbstammes führt. Nicht unwichtig scheint dabei auch zu sein, wie spezifisch das Verb in seiner lexikalischen Bedeutung ist. Ist das verbale Ereigniskonzept semantisch spezifiziert, wie bei fahren oder schwimmen, scheint eine weitere Modifizierung leichter möglich (was dann am Ende zu komplexen N-V-Verbindungen führt oder wenigstens zu Reihenbildungen oder Noun-

controlling the phenomenon to be modelled (in this case, writers' mental grammars with respect to the univertation of N + V units). Model specification thus deserves more attention than (15) has to offer. Mathematically and thus more transparently, the model is given in (16). The notation with angled brackets in  $\alpha_{NV_j[i]}$  should be read as "the value of the random effect  $\alpha_{NV}$  for the factor level  $j$ , chosen appropriately for observation  $i$ ."

$$(16) \quad Pr(Univ_i = 1) = \text{logit}^{-1}[\alpha_0 + \alpha_{NV_j[i]} + \vec{\beta}_{Context} \cdot \vec{x}_{Context_i}]$$

The proportion of compound spellings  $Prop_{Comp_i}$  is the logit-transformed sum of the overall intercept  $\alpha_0$ , the random intercept for the  $j$ -th N + V unit  $\alpha_{NV_j[i]}$  (whichever is found in observation  $i$ ) and the dot product of the vector of dummy-coded binary value for the morphosyntactic context  $\vec{x}_{Context_i}$  and the vector of their corresponding regressors  $\vec{\beta}_{Context}$ . Since it is a multilevel model,  $\alpha_{NV}$  has its own linear model, which is given in (17).

$$(17) \quad \alpha_{NV_j} = \gamma_j + \vec{\delta}_{Relation} \cdot \vec{x}_{Relation_j} + \delta_{Link} \cdot x_{Link_j}$$

It is also assumed that 18 holds, i. e., that the random intercepts for individual N + V units are normally distributed.

$$(18) \quad \alpha_{NV} \sim Norm$$

The random effects are assumed to be a normally distributed variable  $\alpha_{NV}$  which is for each N + V unit  $j$  given as the sum of the conditional mode of unit  $i$  (often wrongly called the *random effect* per se), the dot product  $\vec{\delta}_{Relation} \cdot \vec{x}_{Relation_j}$  of the vector of binary variables encoding the relation and the vector of their corresponding coefficients, and finally the product  $\delta_{Link} \cdot x_{Link}$  of the binary variable encoding the presence of a linking element and its coefficient.

## B Experiment: full specification of the model

In this appendix, we provide the mathematical notation of the model specified in Section 4.2 as (14) and repeated here as (19).

$$(19) \quad Univertation \sim (1|Participant) + Attraction + Context$$



The model is specified in the same notation as in Appendix A in (20). The regressor  $x_{Attract_i}$  is numeric (the attraction score), whereas  $\vec{x}_{Context_i}$  is a dummy-coded vector of binary variables.

$$(20) \Pr(Univ_i = 1) = \text{logit}^{-1}[\alpha_0 + \alpha_{part_j[i]} + \beta_{Attract} \cdot x_{Attract_i} + \vec{\beta}_{Context} \cdot \vec{x}_{Context_i}]$$

It is also assumed that 18 holds, i. e., that the random intercepts for individual participants are normally distributed.

$$(21) \alpha_{participant} \sim Norm$$

## C Sentences used in the experiment

The N + V units are typeset in smallcaps and spelled as separate words. The order of the sentences corresponds to Table 5.

- (22) Lara trat zur Seite, um PLATZ zu MACHEN.  
Lara stepped to the side in order room to make  
Lara stepped aside to make way.
- (23) Sarah ging auf den Spielplatz, um SEIL zu SPRINGEN.  
Sarah went onto the playground in order rope to jump  
Sarah went to the playground to do some skipping.
- (24) Leon konnte nur deshalb gewinnen, weil Johanna ihm  
Leon could only therefore win because Johanna him  
MUT GEMACHT hat.  
courage made has  
Leon could win only because Johanna encouraged him.
- (25) Maria hat einen Kopfhörer gekauft, nachdem sie ihn PROBE  
Maria has a headphone bought after she it test  
GEHÖRT hatte.  
listened had  
Maria bought a headphone after doing a listening test.
- (26) Melanie mag Fußball, weil es ein Sport zum SPASS HABEN ist.  
Melanie likes soccer because it a sport to the fun have is  
Melanie likes soccer because it's a fun sport.

- (27) Benjamin ruft seinen Freund an, weil er eine Frage zum  
 Benjamin calls his friend on because he a question to the  
 BERG STEIGEN hat.  
 mountain climbing has  
 Benjamin calls his friend because he has a question about mountain  
 climbing.
- (28) Kim sah sich das Tennisspiel an, solange sie am TEE  
 Kim watched herself the tennis match on while she at the tea  
 TRINKEN war.  
 drink was  
 Kim watched the tennis match while drinking some tea.
- (29) Simone hört ein Hörbuch, während sie am BOGEN  
 Simone listens an audiobook while she at the bow  
 SCHIESSEN ist.  
 shoot is  
 Simone listened to an audiobook while practicing archery.

## Funding information

Roland Schäfer's work was funded by the Deutsche Forschungsgemeinschaft (DFG, German Research Foundation) – SFB 1412, 416591334.

## Acknowledgments

We thank Luise Rissmann for her help annotating and cleaning the corpus data as well as conducting most of the experiments.

## References

- Arppe, Antti & Juhani Järvikivi. (2007). Every method counts: combining corpus-based and experimental evidence in the study of synonymy. *Corpus Linguistics and Linguistic Theory* 3(2). 131–159. <http://dx.doi.org/10.1515/cllt.2007.009>.
- Bredel, Ursula & Hartmut Günther. (2000). Quer über das Feld das Kopfad-junkt. Bemerkungen zu Peter Gallmanns Aufsatz Wortbegriff und Nomen-Verb-Verbindungen. *Zeitschrift für Sprachwissenschaft* 19(1). 103–110.

- Eisenberg, Peter. (2013). *Grundriss der deutschen Grammatik: Das Wort*. 4th edn. Stuttgart: Metzler. <http://dx.doi.org/10.1007/978-3-476-03762-6>.
- Evert, Stefan. (2008). Corpora and collocations. In Anke Lüdeling & Maria Kytö (eds.), *Corpus linguistics. an international handbook*, 1212–1248. Berlin: Mouton. <http://dx.doi.org/10.1515/9783110213881.2.1212>.
- Fleischer, Wolfgang & Irmhild Barz. (2012). *Wortbildung der deutschen Gegenwartssprache*. Marianne Schröder (ed.). 4th edn. Berlin, Boston: De Gruyter.
- Fortmann, MISSING. (2015). Missing. *MISSING MISSING(MISSING)*. MISSING.
- Fox, John & Sanford Weisberg. (2018). Visualizing fit and lack of fit in complex regression models: effect plots with partial residuals. *Journal of Statistical Software* 87(9). 1–27.
- Fuhrhop, Nanna. (2007). *Zwischen Wort und Syntagma. Zur grammatischen Fundierung der Getrennt- und Zusammenschreibung*. Tübingen: Niemeyer.
- Gallmann, Peter. (1999). Wortbegriff und Nomen-Verb-Verbindungen. *Zeitschrift für Sprachwissenschaft* 18(2). 269–304.
- Gelman, Andrew & Jennifer Hill. (2006). *Data analysis using regression and multilevel/hierarchical models*. Cambridge: Cambridge University Press. <http://dx.doi.org/10.1017/cbo9780511790942>.
- Gries, Stefan Th. (2015). More (old and new) misunderstandings of collocation analysis: on Schmid and Küchenhoff (2013). *Cognitive Linguistics* 26(3). 505–536.
- Jacobs, Joachim. (2005). *Spatien. zum system der getrennt- und zusammenschreibung im heutigen deutsch*. Berlin: de Gruyter. <http://dx.doi.org/10.1515/9783110919295>.
- Kilgariff, Adam, Vít Baisa, Jan Bušta, Miloš Jakubíček, Vojtěch Kovář, Jan Michelfeit, Pavel Rychlý & Vít Suchomel. (2014). The Sketch Engine: ten years on. *Lexicography*. 1–30. <http://dx.doi.org/10.1007/s40607-014-0009-9>.
- Klos, Verena. (2011). *Komposition und Kompositionalität. Möglichkeiten und Grenzen der semantischen Dekodierung von Substantivkomposita*. Berlin, New York: De Gruyter.
- Küchenhoff, Helmut & Hans-Jörg Schmid. (2015). Reply to “More (old and new) misunderstandings of collocation analysis: on Schmid & Küchenhoff” by Stefan Th. Gries. *Cognitive Linguistics* 26(3). 537–547.
- Lehmann, Christian. (2021). Univerbation. *Folia Linguistica Historica* 42. MISSING.
- Mithun, Marianne. (1984). The evolution of noun incorporation. *Language* 60(4). 847–894.

- Morcinek, Bettina. (2012). Getrennt- und Zusammenschreibung: Wie aus syntaktischen Strukturen komplexe Verben wurden. In Livio Gaeta & Barbara Schlücker (eds.), *Das Deutsche als kompositionsfreudige Sprache. Strukturelle Eigenschaften und systembezogene Aspekte*, 83–100. Berlin: De Gruyter.
- Nübling, Damaris, Antje Dammel, Janet Duke & Renata Szczepaniak. (2017). *Historische Sprachwissenschaft des Deutschen. Eine Einführung in die Prinzipien des Sprachwandels*. Tübingen: Narr.
- Schäfer, Roland & Ulrike Sayatz. (2016). Punctuation and syntactic structure in “obwohl” and “weil” clauses in nonstandard written German. *Written Language and Literacy* 19(2). 212–245. <http://dx.doi.org/10.1075/wll.19.2.04sch>.
- Schäfer, Roland. (2020, to appear). Generalized linear mixed models. In Stefan Gries & Magali Paquot (eds.), *The practical handbook of corpus linguistics*. Berlin, Heidelberg: Springer.
- Schäfer, Roland. (N.d.). Statische inferenz in der linguistik. in preparation.
- Schäfer, Roland & Felix Bildhauer. (2012). Building large corpora from the web using a new efficient tool chain. In Nicoletta Calzolari (Conference Chair), Khalid Choukri, Thierry Declerck, Mehmet Uğur Doğan, Bente Maegaard, Joseph Mariani, Asuncion Moreno, Jan Odijk & Stelios Piperidis (eds.), *Proceedings of the eighth international conference on language resources and evaluation (LREC '12) international conference on language resources and evaluation (LREC 12)*, 486–493. Istanbul: European Language Resources Association (ELRA).
- Schäfer, Roland & Elizabeth Pankratz. (2018). The plural interpretability of German linking elements. *Morphology* 28(4). 325–358. <http://dx.doi.org/10.1007/s11525-018-9331-5>.
- Schlücker, Barbara. (2012). Die deutsche Kompositionsfreudigkeit: Übersicht und Einführung. In Livio Gaeta & Barbara Schlücker (eds.), 1–25. Berlin: De Gruyter. <http://dx.doi.org/10.1515/9783110278439.1>.
- Schmid, Hans-Jörg & Helmut Küchenhoff. (2013). Collostructional analysis and other ways of measuring lexicogrammatical attraction: theoretical premises, practical problems and cognitive underpinnings. *Cognitive Linguistics* 24(3). 531–577. <http://dx.doi.org/10.1515/cog-2013-0018>.
- Stefanowitsch, Anatol & Stefan Th. Gries. (2003). Collostructions: investigating the interaction between words and constructions. *International Journal of Corpus Linguistics* 8(2). 209–243. <http://dx.doi.org/10.1075/ijcl.8.2.03ste>.
- Szczepaniak, Renata. (2009). *Grammatikalisierung im Deutschen. Eine Einführung*. Tübingen: Narr.

- Wurzel, Wolfgang Ullrich. (1994). Inkorporierung und “Wortigkeit” im Deutschen. In Wolfgang U. Dressler (ed.), *Natural morphology: perspectives for the nineties*, 109–125. Wien: Unipress.
- Wurzel, Wolfgang Ullrich. (1998). On the development of incorporating structures in German. In Richard M. Hogg & Linda van Bergen (eds.), *Historical linguistics 1995*, 331–344. Amsterdam, Philadelphia: Benjamins.
- Zuur, Alain F., Elena N. Ieno, Neil Walker, Anatoly A. Saveliev & Graham M. Smith. (2009). *Mixed effects models and extensions in ecology with R*. Berlin etc.: Springer. <http://dx.doi.org/10.1007/978-0-387-87458-6>.