

Zero marking in inflection: A token-based approach

Anonymous

ABSTRACT

This study examines zero marking, i.e. the absence of an overt exponent, in adjectival, nominal and verbal inflectional morphology across languages. The first part of the study provides an overview of the distribution of zero markers in inflection paradigms using the UniMorph dataset. The results show that there is a general preference against zero marking. The distribution of zero markers varies to a great extent across languages and lemmas, the only robust trend being that they are avoided in cells that express a high number of grammatical values. The second part of this study examines the association between marker frequencies and phonological length using the Universal Dependencies treebanks. While token frequency is a good predictor for the length of overt markers, it does not account for the occurrence of zero markers. This is taken as evidence to support a differential non-development scenario of zero marking rather than a phonetic reduction scenario.

Keywords:
token-based
typology,
corpus typology,
zero marking,
zero exponence

INTRODUCTION

1

The present study examines the distribution of zero forms in adjectival, nominal and verbal inflectional morphology. In typology, zero marking plays an important role for coding efficiency or form-frequency effects in morphosyntax. The analysis of form-frequency effects go

back to the early findings by Zipf (1935) that more frequent lexical elements tend to be shorter than less frequent ones. There is crosslinguistic evidence that also in inflectional morphology, more frequent or predictable markers tend to be shorter or at least not longer than comparable less frequent markers (Greenberg 1966; Guzmán Naranjo and Becker 2021; Stave *et al.* 2021; Haspelmath 2008b; Haspelmath *et al.* 2014; Haspelmath 2021; Haspelmath and Karjus 2017). Such effects can be subsumed under the term of coding efficiency. The coding of grammatical expressions is efficient, because it saves effort in the production and processing of speech but maintains the successful transfer of information (cf. Levshina 2022, for an overview of efficiency in language and communication).

Usually, zero markers (in the sense of zero exponence) are grouped with shorter markers as opposed to longer ones. It is often explicitly or implicitly assumed that zero markers are used to express highly frequent morphosyntactic functions similarly to shorter markers (e.g. Bybee 2011; Croft 2003, Ch. 4; Diessel 2019, Ch. 11; Greenberg 1966, 32-37; Haspelmath 2008a, 2008b, 2021; Song 2018, Ch. 7). However, a quantitative crosslinguistic overview of the distribution of zero marking in inflection is still not available. The objective of this paper is to start filling this gap.

To do so, I analyze the distribution of zero markers in the UniMorph dataset (McCarthy *et al.* 2020), which is a crosslinguistic database of inflectional paradigms for individual lemmas. I first provide some theoretical background on zero marking and coding efficiency and introduce a working definition of zero markers in Section 2. Section 3 describes the dataset as well as the marker extraction, and it discusses examples of zero markers. I then analyze the probability of zero marking using the UniMorph dataset in Section 4. As will be seen, zero marking is generally dispreferred across languages and parts of speech. Section 5 then zooms in on those cells and values of adjectival, nominal and verbal inflectional paradigms that are most likely to be zero marked across languages. In Section 6, I turn to the distribution of zero markers in language use. Using corpus data from the Universal Dependencies treebanks (Zeman *et al.* 2023), I analyze the association between token frequencies of inflection markers and their phonological length, including the distribution of zero markers. As we will see, frequency does not affect zero markers in the same

way as overt markers. Section 7 discusses the findings of this study with a special focus on the role of coding efficiency to account for the distribution of zero marking. Section 8 concludes.

ZERO MARKING

2

This section introduces the relevant theoretical notions related to zero marking. Section 2.1 introduces zero marking and its relation to coding efficiency in typology. In Section 2.2, I then propose a working definition of zero markers for the purposes of the present study. Throughout the paper, I use zero marking to refer to the absence of phonetic exponence (“zero exponence”) of a morphosyntactic function.

Zero marking and coding efficiency

2.1

The modern understanding of coding efficiency began with Zipf (1935), who showed that more frequent words tend to be shorter than less frequent words. Greenberg (1966, 1963) was one of the first typologists to relate the token frequencies of grammatical values to their formal markedness. An “unmarked” value in this sense is characterized by the absence of an exponent, which is contrasted with a “marked” value that is expressed by an overt exponent. For instance, Greenberg (1966, 32-37) showed how the markedness of singular, plural and dual forms of nouns, verbs, and adjectives is reflected in their distribution in corpora from various languages. He noted that the “unmarked” number value, singular, is substantially more frequent than the usually “marked” number values of plural and dual in corpus data from different languages.

Taking up Greenberg's findings and doing away with the concept of markedness, Haspelmath (2008a,b) argued that the length, complexity or availability of grammatical markers can be accounted for by their frequency in language use. In a more recent study, Haspelmath (2021, 2) proposed the following form-frequency correspondence hypothesis:

- (1) *The grammatical form-frequency correspondence hypothesis*
When two grammatical construction types that differ minimally (i.e. that form a semantic opposition) occur with significantly different frequencies, the less frequent construction tends to be overtly coded (or coded with more segments), while the more frequent construction tends to be zero-coded (or coded with fewer segments), if the coding is asymmetric. (Haspelmath 2021, 2)

This hypothesis includes the assumption that zero forms pattern with shorter forms in being used for coding comparatively frequent expressions. Applied to inflectional morphology, we should thus expect zero marking for highly frequent values of morphosyntactic features. By now there is indeed much evidence for effects of coding efficiency between comparable grammatical expressions. However, examples usually only involve a difference in lengths, i.e. shorter vs. longer forms.¹ The participation of zero forms has not yet been the focus of any systematic crosslinguistic study. There are some indications from the literature, though, which suggests that coding efficiency and frequency may not be a suitable explanation for the distribution of zero markers.

Stolz and Levkovich (2019) provide a qualitative overview of the distribution of zero marking in inflection (“absence of material exponence, AOME”) from the perspective of canonical morphology. They note that “[f]rom the small number of cases discussed above it transpires that frequency might not always be the most powerful factor to make a given word-form or category a candidate for AOME” (Stolz and Levkovich 2019, 396-397).

Guzmán Naranjo and Becker (2021) come to a similar conclusion based on a quantitative analysis of the association between the length of nominal inflection markers and their distribution across paradigms. They also use the UniMorph database, but focus on nominal inflec-

¹ A few examples of quantitative approaches to form-frequency effects in grammar are: Guzmán Naranjo and Becker (2021) for the length and paradigmatic distribution of nominal inflection markers, Stave *et al.* (2021) for the length and frequency of morphemes in general, Haspelmath *et al.* (2014) for the expression of causal and non-causal alternations, Haspelmath (2008c) for reflexive marking, Haspelmath and Karjus (2017) for number marking and Ye (2020) for (in)dependent possessor marking.

tion and test different distributional factors for their association with marker length. Although they find that marker length is associated with their type frequency, their results suggest that other measures such as the entropy of the marker are better predictors for their length. With their main focus being on predicting marker length from distributional measures, one detail of their analysis concerns zero marking and is highly relevant for the present study. Guzmán Naranjo and Becker (2021) note that a simple Poisson model to predict marker length strongly overestimates the occurrence of zero markers. This suggests that the distribution of zero markers does not simply follow the pattern of shorter ones.

Another area in which zero marking has been mentioned to behave differently is the occurrence of zero for person and number marking on verbs. Several quantitative typological studies (Bickel *et al.* 2015; Cysouw 2003; Siewierska 2010) find that zero for person marking is rather uncommon across languages. In contrast to the traditional view in typology, these studies do not find evidence for a paradigmatic preference of third person (singular) being zero marked on the verb. However, all three studies show that if a person marker is zero, it more likely expresses third person (singular) than first or second person.

Seržant and Moroz (2022) also mention zeros in verbal person-number marking. Analyzing the length of person-number markers in a typological sample, they argue for an attractor state in which the lengths of different indexes are associated with their frequencies in language use. Seržant and Moroz (2022, 6) note that “[...] articulatory efficiency plays an important role here: the more expected the sign is the shorter it is. Nevertheless, zero is not preferred.” Yet, they motivate the crosslinguistic avoidance of zero forms by invoking two types of efficiency: processing and planning efficiency. Seržant and Moroz (2022, 7) hypothesize that an overt exponent facilitates processing on the addressee's side. They also propose that avoiding zero marking makes planning more efficient on the speaker's side, “[...] because it provides a straightforward link from meaning to coding, while zero is inherently ambiguous by being linked to various meanings and domains” (Seržant and Moroz 2022, 7). Whether or not the avoidance of zero marking can indeed be accounted for by processing or planning efficiency requires proper psycholinguistic testing. The relevant point is that coding efficiency does not seem to be applicable to the

frequency distribution of zero markers in person indexing in the same way as it is for overt markers.

2.2

A working definition of zero markers

The discussion and use of zero has a long tradition in morphology and in linguistics in general. It goes back to Pāṇini, who introduced the idea of zero morphs for morphemes that lack a phonetic representation as the outcome of morphological rules (Robins 1997, 181-182). The concept of zero morphs for linguistic analysis was also widely applied in later work of structuralists, e.g. Bloch (1947); Bloomfield (1933); Jakobson (1983) and Saussure (1916).² Starting with Haas (1957), linguists began to criticize the assumption of zero morphs in the structuralist tradition and argued for stricter criteria to define zero morphs in order to avoid the assumption of excessive linguistic structure (e.g. Sanders 1988; Mel'čuk 2002; McGregor 2003). The potential danger being that the linguist may postulate a zero morph for any single morphosyntactic function that does not correspond to an overt exponent. As Anderson (1992, 30) notes, it “leads to the formal problem of assigning a place in the structure (and linear order) to all of those zeros”.³ Others, such as Arkadiev (2016); Contini-Morava (2006) and Mithun (1986), used data from typologically diverse languages to show that the absence of phonetic material can also correspond to the absence of a morphosyntactic feature rather than to zero marking.

In line with those more cautious approaches to zero morphs, this study uses the notion of “zero marker” as a descriptive shorthand for the absence of material exponence of a given morphosyntactic function (cf. Stolz and Levkovych 2019). In other words, I do not assume the presence of a zero morph. Instead, I understand zero markers as

²For more details, see Meier (1961). See also Al-George (1967); Diehl (2008) and McGregor (2003) for more details on the history of linguistic zero.

³For examples and discussions on issues related to the use of zero morphs in morpheme-based, segmental approaches to morphology, see Anderson (1992); Pullum and Zwicky (1991); Blevins (2016); Bank and Trommer (2015). For overviews of zero exponence in morphological theories, see Trommer (2012) and Dahl and Fábregas (2018).

the absence of exponence which expresses a certain morphosyntactic function in addition to the lexical content of a word form. This also means that zero markers can only occur in contrast to at least one other, overtly coded morphosyntactic function of the same inflectional paradigm.

To analyze the distribution of zero markers in inflectional morphology, we need to identify the invariable, lexical parts (stems) as well as the potential exponents of a morphosyntactic function in an inflected word form. This conforms with the basic intuition that we want to separate the segments that convey the word's lexical meaning from the segments that convey morphosyntactic information (cf. Matthews 1972).⁴ For the purposes of the present study, I define stems, markers, and zero markers as shown in (2), (3) and (4), respectively. These definitions are motivated by both theoretical as well as practical considerations regarding the dataset and annotations available.

(2) *Stem*

The stem expresses the lexical content of a word form; it corresponds to the longest common subsequence shared by all inflected forms of a word.

(3) *Marker*

A marker encodes the morphosyntactic function of a word form, i.e. a value of some morphosyntactic feature defined for that word or a bundle of values of several such features. The marker corresponds to the phonetic material outside of the stem of a word form.

(4) *Zero Marker*

A zero marker occurs when the word form does not feature any overt marker (as defined in (3)) to encode its morphosyntactic function. If the morphosyntactic function of the word consists of several morphosyntactic features, zero marking applies

⁴In reality, the identification of stems is not always this straightforward. There are many different ways in which the lexical parts of inflected words can vary in their phonological shape. Baerman and Corbett (2012) provide a number of examples and introduce a canonical approach to stems to capture that variation.

to the combination of feature values and not to single feature values in isolation.

Consider a simple example of stem and marker identification. The paradigm of English nouns consists of two cells: the singular form and the plural form. Given the paradigmatic relation between the singular form /*dei*/ (*day*.SG) and the plural form /*dez*/ (*day*.PL), we can identify the string /*dei*/ as the stem, i.e. the phonetic material that both forms of the paradigm share. Since the form filling the plural cell includes the additional material /*z*/, we can establish /*z*/ as a plural marker. In the singular cell, the form does not include any material other than what was identified as the stem. We can therefore treat the form of the singular cell of *day* in English as zero marked.

However, as will be described in detail in Section 3.3, I automatically adjusted the stems extracted according to the definition in (2) in order to account for stem allomorphy to a certain extent. This is motivated by the fact that many stem alternations are phonologically driven, which means that they do not necessarily provide meaningful insights about the inflectional properties of a system in general and about the distribution of zero marking in particular. Ignoring such alternations allocates additional material to the marker segments and runs the risk of systematically underestimating the number of zero markers. The adjusted marker_A and zero marker_A, which take into account stem alternations, are operationalized as described in (5) and (6), respectively.

- (5) *Marker_A*
A marker_A is extracted from a marker as defined in (3) by removing all material from those affix position that the system does not use for inflection.
- (6) *Zero Marker_A*
A zero marker occurs when the word form does not feature any overt marker (as defined in (5)) to encode its morphosyntactic function.

This operationalization of stems, markers_(A) and zero markers_(A) has the practical advantage that it does not require any morphological analysis particular to a single language or paradigm. It is a solution

to identify the segments that contribute inflectional information that can be applied automatically and consistently to the crosslinguistic UniMorph dataset used in this study.

Besides practical considerations, this method is also based on theoretical grounds and follows the definition of stems by Beniamine and Guzmán Naranjo (2021); Bonami and Beniamine (2021) and Guzmán Naranjo and Becker (2021). Despite much theoretical work on the role and identification of stems in morphology, Bonami and Beniamine (2021) note that “there is no agreed upon method for identifying which part of an inflected word is a stem, and that the heuristics used by morphologists in that area are neither systematic nor principled enough”.⁵ The authors compare two types of stem identification based on prioritizing two different principles, namely to avoid stem allomorphy and to avoid discontinuous stems. Since those two principles are in conflict with each other many times, every approach to stem identification needs to rank them in some way to resolve such conflicts. Bonami and Beniamine (2021) compare the two methods of either adhering to the first or the second principle, resulting in what they call “unique discontinuous stems” (no stem allomorphy allowed) and “continuous stem sets” (no discontinuous stems allowed). While the first method of unique discontinuous stems allocates all the variation of word forms to the exponents, leading to more exponent allomorphy, the second method of continuous stem sets keeps exponent allomorphy minimal, but leads to a high degree of stem allomorphy, since all variation that is enclosed by stem segments has to be included in the stems. What this shows is that neither approach creates more allomorphy; they simply allocate it differently. Of course, which of the two approaches is more useful depends on the research question at hand.

One of the questions discussed by the authors is what types of stems are more helpful in addressing the ‘Inflected Word Recognition Problem’ (IWRP), i.e. understanding what allows speakers to draw inferences from a word's form about its content. This results in the task

⁵ Cf. Blevins (2003); Bonami (2012); Brown (1998); Maiden (1992); Montermini and Bonami (2013); Pirrelli and Battista (2000); Spencer (2012); Stump (2001); Stump and Finkel (2013) for work on stem identification and stem allomorphy.

of separating the lexical and the inflectional parts of a word form, and Bonami and Beniamine (2021) note that “[i]n terms of the IWRP, the answer is quite simple. Sets of continuous stems are by definition less useful than a unique discontinuous stem: the unique discontinuous stem identifies exactly that part of the word that has no exponential value, while stem allomorphs blur the distinction between exponential and nonexponential material.” As the identification of zero forms relies on separating lexical segments from exponents of morphosyntactic information in word forms, the IWRP is of high relevance to this study and provides the theoretical grounds for the definition of stems given in (2).

Furthermore, this study will largely follow a word and paradigm approach to inflection (cf. Anderson 1992; Blevins 2016; Matthews 1972; Stump 2001; Zwicky 1985). This approach bases morphological analyses on the paradigmatic relation between different word forms, representing the different morphosyntactic functions a given word can have. The exponent of a cell in an inflectional paradigm is determined through the relation of the word form to the forms used for the other cells of the paradigm. The word and paradigm approach has a very important practical advantage. It allows us to refrain from further segmentation of exponents into morphemes which may require language-specific insights and which may not always be desirable or useful (cf. Blevins 2005, 2006).

Although morphological segmentation analyses may sometimes be uncontroversial, there are many cases where a morpheme analysis is less than clear. Various examples are given in Spencer (2012), one of them being the Spanish subjunctive verb form *cantaríamos* ‘we would sing’. A number of theoretical motivations exist to segment this word form into morphemes in five different ways: (i) *cant-a-r-í-a-mos*, (ii) *canta-ríamos*, (iii) *cant-a-ría-mos*, (iv) *canta-r-í-a-mos* and (v) *cantar-íamos* (Spencer 2012, 93). The fact that these profoundly varying morphological analyses are motivated in the literature suggests that such morpheme segmentations are always, whether explicitly or implicitly, theoretically guided. Moreover, it is likely that the segmentation into morphemes in lesser-studied languages involves even more theoretical uncertainty, given that we may know much less about the morphological structure and its diachrony than for languages like Spanish.

As will be shown in more detail in Sections 3.3 and 3.4, cells of

paradigms are defined by (a combination of) values of morphosyntactic features. For instance, the inflectional paradigms of German nouns combine the morphosyntactic features of case and number. While nouns are inherently specified for gender, each word form in context is also specified for number and case so that each cell of the paradigm corresponds to a number-case combination, e.g. dative plural.

For the purposes of this study, I do not distinguish between an exponent for plural number and one for dative case. Instead, I treat the material in addition to the stem in the dative plural cell as the marker of the dative-plural function. When no additional phonetic material is used, this cell is then analyzed as being zero marked (cf. Table 9). Put differently, I do not assign zero markers to single abstract morphosyntactic values but to the relevant value combinations of the inflectional paradigms. The theoretical reason to do so lies in exponents of morphosyntactic functions being defined based on the relations between the forms of the different cells of the inflection paradigm, which combine these function. This also reflects the morphological reality of many if not most languages in that morphosyntactic functions are usually not marked in isolation but often occur in combinations. As mentioned above, it is not always trivial to justify a segmental analysis. The practical reason is that there is still no language-independent and theory-independent way of segmenting distinct morphosyntactic exponents, and those segmentations are not (yet) automatizable. Since automatic processing is indispensable for the purposes of the present study, no further segmentation of different morphosyntactic exponents will be carried out.

The segmentation into stems and markers is often additionally complicated by inflection classes, which use different types of markers. Sections 3.3 and 3.4 show in more detail how the present approach deals with variation in the exponents due to inflection classes, with stem alternations and with suppletive forms.

The data used in this study comes from the UniMorph database (McCarthy *et al.* 2020), a large-scale crosslinguistic database of complete inflectional paradigms of adjectives, nouns and verbs for individual lexemes from different languages. The present study includes adjectival, nominal and verbal paradigms of 39, 62 and 96 languages, respectively. Some languages are featured with paradigms for more than one part of speech; the total number of languages analyzed in this study is 114. Figure 1 shows the geographical distribution of the languages in the dataset.⁶

While the dataset is not a balanced typological sample in the strict sense, it does include languages from all six macro areas (Africa, Eurasia, Papunesia, Australia, North America and South America), which ensures that typological and areal diversity is captured at least to some degree. Table 1 provides an overview of the final dataset with the number of languages, lemmas, paradigm cells, marker types and observations by part of speech.

Table 1:
Dataset overview

	N lang	N lemma	N cell	N markers	N obs
adjectives	39	157355	961	5552	6348198
nouns	62	610242	727	19537	6261881
verbs	96	129377	2753	47457	4407743

The morphosyntactic annotation in the UniMorph dataset follows the guidelines described in Sylak-Glassman (2016). Sylak-Glassman (2016, 3) notes: “This paper presents the Universal Morphological Feature Schema (UniMorph Schema), which is a set of morphological features that functions as an interlingua for inflectional morphology by defining the meaning it conveys in language-independent

⁶More details about the languages, the part of speech and the number of lexemes is provided in the files `affixation.csv` and `lemmas.csv` in the supplementary materials. All supplementary materials referred to in this paper can be found here: https://osf.io/e48qc/?view_only=eb53b1e02e034a459c335a0736941f9b

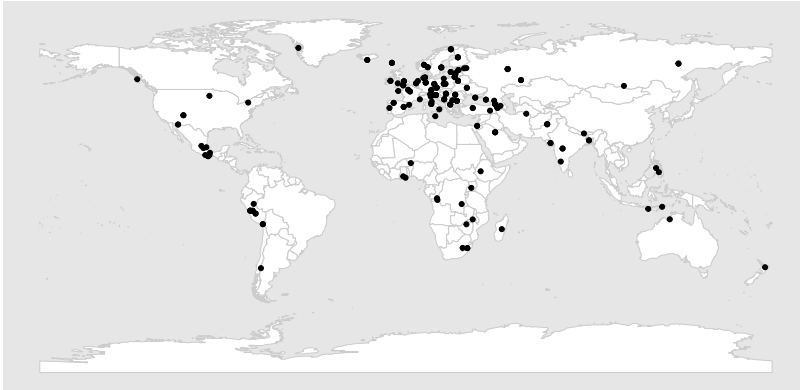


Figure 1:
Location of the
languages in the
dataset

terms. The features of the Universal Morphological Feature Schema have precise definitions based on attested cross-linguistic patterns and descriptively-oriented linguistic theory, and can capture the maximal level of semantic differentiation within each inflectional morphological category.” Annotations thus do not necessarily follow the linguistic traditions of particular languages but are defined and used in the sense of comparative concepts in typology (cf. Haspelmath 2018).

Data pre-processing

3.2

I excluded a number of languages available in UniMorph from the final analysis on the basis of unclear or insufficient annotations in the original datasets, some of which were annotated only automatically with no manual checks. Since the database is somewhat biased towards languages spoken in Eurasia (mostly Indo-European languages), I only included languages with paradigms for more than 30 lemmas from this area. For languages from other macro areas, especially from Africa or the Americas, I did not apply this threshold of 30 lemmas in order to include more non-Indo-European languages and to keep the dataset as diverse as possible.⁷

⁷For adjectives, only Zulu has fewer than 30 lemmas (17); for nouns, this is the case in Kalaallisut (23). For verbal paradigms, languages with fewer than 30 lemmas are Sotho (26), Mapudungun (26), Murrinpatha (29) and Zarma (27).

The next step was to pre-process the data to remove errors and to make annotations more consistent across languages.⁸ The pre-processing consisted of different global as well dataset-specific corrections. Global corrections included resolving inconsistencies in the annotations across languages. For instance, the value “indefinite” was coded as “INDF” in some languages and as “NDEF” in others. Similarly, the annotation of person-number combinations in verbs varied, e.g. between “SG;1”, “1;SG”, “1SG” for first person singular. In such cases, I adjusted the annotation to a single label across all languages. I also removed complex lemmas containing a space or “-”. On the one hand, this removed some erroneous lemmas that were complex expressions rather than nouns, adjectives or verbs. On the other hand, in some languages both parts of a complex noun or adjective are inflected. Leaving such lemmas in would have caused the marker extraction to detect infixation for complex lemmas with suffixes on two or more parts. Removing them avoided the artificial creation of more complex inflection patterns. Similarly, periphrastic forms were removed in cases of inflected auxiliaries, which would equally have led to the erroneous analysis of infixation. Complex forms were also removed if they contained a separate marker that occurred before or after the inflected verb form, depending on the cell of the paradigm. This was especially common with verbal paradigms, e.g. verbal particles in German or reflexive markers in Italian and Macedonian.

Dataset-specific cleaning steps included, for instance, deleting “?” following the interrogative verb forms in the Turkish data or deleting the indefinite article from Romanian nominal forms. Other cleaning steps had to do with the alphabetic scripts used. For instance, the Serbian-Bosnian-Croatian dataset contained forms in the Latin script with a handful of forms in Cyrillic. The latter were removed to allow for consistent processing. Some datasets, e.g. Old French or Yoloxochitl Mixtec, contained alternative forms for certain cells. In those cases, I systematically left the first form in and removed the other ones.⁹ Other dataset-specific operations included deleting single forms

⁸A detailed documentation of all pre-processing steps can be found in `preprocessing.txt` in the supplementary materials. For the implementation, see `code-preprocessing.R`.

⁹It would have been insightful to include overabundance in a systematic way.

containing obvious errors (e.g. misalignment, cells with missing data).

Following the data cleaning, I added phonological transcriptions to the inflected forms whenever possible. For some languages, e.g. Palantla Chinantec, the UniMorph database already provided the inflected forms in a phonological transcription. For most other languages, however, forms were given in the standard orthographic representation. This can of course be problematic, especially in languages such as French, where the orthographic representation continues to make many distinctions that are no longer realized in the spoken language. For this reason, whenever possible, I replaced the orthographic forms by a phonological transcription using Epitran (Mortensen *et al.* 2018). Epitran currently has modules to transcribe 31 of the languages used here.¹⁰

While not perfect, Epitran offers a more realistic representation of the forms occupying the different cells of inflectional paradigms. Table 2 illustrates this by showing the transcriptions generated with Epitran for the French verb *allumer* ‘light something, turn on (light)’. The rows show seven TAM combinations; for each of these, the first row contains the form in their orthographic representation, and the second row shows the phonological transcriptions generated with Epitran.

For the remaining 81 languages, the forms in UniMorph are given in their orthographic representation, which reflect phonological shapes to a varying degree. To consider the potential influence that the type of phonological representation may have on the detection of zero forms, I manually coded whether or not the representation was phonological.¹¹ Orthographic representations that systematically reflected phonology were treated as phonological representations. This led to 31 languages with a transcription generated using Epi-

Overabundance refers to the phenomenon of two distinct forms being available to express a single cell in a paradigm (cf. Thornton 2012). However, alternative forms are not systematically annotated in the UniMorph datasets. If provided, their relation differs greatly across datasets and it is usually not documented in the dataset descriptions. Alternatives can represent diachronic, dialectal or stylistic variants; in other cases their alternation behavior remains unclear. It is also unclear how many overabundant forms are not provided in UniMorph. Including overabundance is thus not possible with the approach used in this study.

¹⁰ For details, see `epitran.py` in the supplementary materials.

¹¹ For details by language, see `affixation.csv` in the supplementary files.

Table 2:
Phonological
transcription of
the French verb
allumer ‘turn on
(light)’

	1SG	2SG	3SG	1PL	...
PRS.IND	<i>allume</i> alym	<i>allumes</i> alym	<i>allume</i> alym	<i>allumons</i> alymõ	
PST.IPFV.IND	<i>allumais</i> alyme	<i>allumais</i> alyme	<i>allumait</i> alyme	<i>allumions</i> alymiõ	
PST.PFV.IND	<i>allumai</i> alyme	<i>allumas</i> alyma	<i>allumat</i> alyma	<i>allumâmes</i> alymam	
FUT	<i>allumerai</i> alymre	<i>allumeras</i> alymra	<i>allumera</i> alymra	<i>allumerons</i> alymreõ	
PRS.COND	<i>allumerais</i> alymre	<i>allumerais</i> alymre	<i>allumerait</i> alymre	<i>allumerions</i> alymriõ	
PRS.SUBJ	<i>allume</i> alym	<i>allumes</i> alym	<i>allume</i> alym	<i>allumions</i> alymiõ	
PST.SUBJ	<i>allumasse</i> alymas	<i>allumasses</i> alymas	<i>allumât</i> alyma	<i>allumassions</i> alymasiõ	
...					

tran, 63 languages with original representations that systematically reflect phonological shapes, and 20 languages with orthographies that do not always reflect phonological shapes. The type of phonological representation was then added as a control variable in the analysis.

3.3

Extracting stems and zero markers

In order to analyze the distribution of zero markers, I automatically segmented the inflected word forms following the method developed in Beniamine and Guzmán Naranjo (2021) and Guzmán Naranjo and Becker (2021). As mentioned in Section 2.2, the segmentation follows a word and paradigm approach to morphology in that whole forms are paired with morphosyntactic functions according to their distribution across the inflectional paradigms. This means that the subsequence shared by all cells of the paradigm is automatically extracted and taken as the stem according the working definition given in (2). All material not included in this subsequence is analyzed as the marker of a given cell, as defined in (3). If the form corresponds to the longest common subsequence (i.e. the stem), the marker is analyzed as zero according to the definition in (4). This automated detection of stems and markers is necessary for two reasons. First, it is not feasible to

apply manual, language-specific segmentations to this dataset. Second, this method allows for a single, consistent way of detecting zero marking across languages, which is necessary for the crosslinguistic comparisons made in this study.¹²

cell	form	stem	marker
PRS.IND.1SG	alym	alym	-
PRS.IND.2SG	alym	alym	-
PRS.IND.3SG	alym	alym	-
PRS.IND.1PL	alymɔn	alym	-ɔn
PRS.COND.1SG	alymere	alym	-ere
PRS.COND.2SG	alymere	alym	-ere
PRS.COND.3SG	alymere	alym	-ere
PRS.COND.1PL	alymɛrjɔn	alym	-ɛrjɔn
PRS.SUBJ.1SG	alym	alym	-
PRS.SUBJ.2SG	alym	alym	-
PRS.SUBJ.3SG	alym	alym	-
PRS.SUBJ.1PL	alymjɔn	alym	-jɔn
...

Table 3:
Marker
extraction for the
French verb
allumer ‘turn on
(light)’

To give an example of the segmentation into stems and markers and of the detection of zero markers, Table 3 shows parts of the present tense paradigm of the French verb *allumer* from Table 2. Comparing the forms of the different cells of the paradigm, the string *alym* is detected as the longest common subsequence between all forms of the paradigm. For the purposes of the present paper, this subsequence is analyzed as the stem. All remaining material is analyzed as the marker of a particular cell. In cells where the form corresponds to the stem, markers are analyzed as zero. Here, this is the case for some of the present tense forms; the respective markers are shaded in gray in Table 3.

Aymara (Aymaran) is a language with nominal inflection known for its subtractive morphology. The accusative singular cell is usually analyzed as being expressed by the subtraction of the final vowel of the nominative singular form (cf. Coler 2015). Table 4 illustrates this with parts of the paradigms of two Aymara nouns. However, for the purposes of this study, the accusative singular form corresponds to the

¹²Stem alternations are not accounted for by this extraction method; Section 3.4 shows how they are included in the present study.

stem, because it equals the longest common subsequence of all forms of the lexeme. Compared to the accusative form, the nominative form has an additional final vowel, which is also found in all other forms of the paradigm except for the inessive (INESS) and equative (EQTV) forms.

Table 4:
Marker
extraction for the
Aymara nouns
anu ‘dog’ and
chaski
‘messenger’

cell	form	stem	marker	form	stem	marker
NOM.SG	anu	an	-u	chask	chask	-i
ACC.SG	an	an	-	chask	chask	-
GEN.SG	anuna	an	-una	chaskina	chask	-ina
COM.SG	anumpi	an	-umpi	chaskimpi	chask	-impi
ABL.SG	anuta	an	-uta	chaskita	chask	-ita
ALL.SG	anuru	an	-uru	chaskiru	chask	-iru
INESS.SG	anpacha	an	-pacha	chaskpacha	chask	-pacha
EQTV.SG	anjama	an	-jama	chaskjama	chask	-jama
...

Traditionally, the nominative form with the final vowel is analyzed as the stem of the noun, while the accusative is argued to be a subtractive form, i.e. consisting of less material than the stem of the lexeme (Coler 2015, 2018; Baerman *et al.* 2017). Diachronically speaking, there are valid arguments to support such an analysis. Coler (2018) provides examples of historical Aymara with accusative forms that still have the final vowel. In addition, vowel deletion is a common phonological process in Aymara. Nevertheless, aiming at a synchronic and comparable analysis across languages, I treat the accusative form as the stem of the lexeme here. In the Aymara data, the accusative is zero marked in all 1522 nouns of the dataset with no exception.

Another rather unusual case of zero marking can be found in Georgian (Kartvelian) verbs. Besides a number of other theoretically interesting patterns, Georgian verbs have been cited in the typological and morphological literature for their crosslinguistically unusual 2nd person singular zero marker (e.g. Stolz and Levkovych 2019; Anderson 1992; Blevins 2016). However, not all lexemes express the second person singular form with a zero marker in the sense of the present study. Only 1 out of 118 verbal lexemes in the dataset feature a zero marker in the second person singular present tense cell. Table 5 shows

this for the verb *ts'ers* 'write' in opposition to *ak'eteb* 'make'.¹³

cell	form	stem	marker	form	stem	marker
PRS.1SG	vts'er	ts'er	v-	vak'eteb	k'et	va-eb
PRS.2SG	ts'er	ts'er	-	ak'eteb	k'et	a-eb
PRS.1PL	vts'ert	ts'er	v-t	vak'etebt	k'et	va-ebt
IMPF.1SG	vts'erde	ts'er	v-de	vak'etebdi	k'et	va-ebdi
IMPF.2SG	ts'erde	ts'er	-de	ak'etebdi	k'et	a-ebdi
IMPF.1PL	vts'erdet	ts'er	v-det	vak'etebdit	k'et	va-ebdit
FUT.1SG	davts'er	ts'er	dav-	gavak'eteb	k'et	gava-eb
FUT.2SG	dats'er	ts'er	da-	gaak'eteb	k'et	gaa-eb
FUT.1PL	davts'ert	ts'er	dav-t	gavak'etebt	k'et	gava-ebt
AOR.1SG	davts'ere	ts'er	dav-e	gavak'ete	k'et	gava-e
AOR.2SG	dats'ere	ts'er	da-e	gaak'ete	k'et	gaa-e
AOR.1PL	davts'eret	ts'er	dav-et	gavak'etet	k'et	gava-et
...

Table 5:
Marker
extraction for the
Georgian verbs
ts'ers 'write' and
ak'eteb 'make'

In general, Georgian verbs take a so-called preverb in some but not all of the tenses (Hewitt 1995, 148-169). When it occurs, it precedes the prefixal part of agreement marking on the verb. As we can see in Table 5, present and imperfect forms occur without the verbal prefix, while the future, aorist and perfect forms all make use of the prefix (*da-* and *ga-* in the examples in Table 5). In most TAM series, many Georgian verbs also have so-called thematic suffixes (Hewitt 1995, 143-147), as e.g. *-eb* in *ak'eteb* 'make'. The presence of those thematic suffixes in the present tense results in the absence of zero marking in most of the verbs. The thematic suffix *-eb/-ob* is part of the second person singular present form, but as it is not used in the aorist forms, the former does not correspond to the longest common subsequence of the verb forms. The second person singular present tense cell can thus only be expressed by a zero form with verbs that generally do not use any of the thematic suffixes, like the verb *ts'ers* 'write' in Table 5.

Arabic (Semitic) is well known for having roots that consist of discontinuous consonants, with prefixed, infixes and suffixed vowels and other consonants to mark the grammatical values of a given form in the paradigm (e.g. Ratcliffe 1998; Schramm 1962; Yip 1988; Boude-laia and Marslen-Wilson 2001). The automatic extraction of the longest

¹³ The segment *-a-* is not part of the verb stem of *ak'eteb* 'make', as it does not occur in all forms of the paradigm, e.g. the imperfective masdar form *k'etebi*.

common subsequence detects these consonants and assigns all additional material to the markers. This is shown for two verbs, *ʔarsala* ‘send’ and *iktašafa* ‘discover’ in Table 6.

Table 6:
Marker
extraction for
Arabic verbs
ʔarsala ‘send’
and *iktašafa*
‘discover’

cell	form	stem	marker	form	stem	marker
IPFV.1SG	ʔursilu	rsl	ʔu-i-u	ʔaktašifu	ktšf	ʔa-a-i-u
IPFV.2SG.F	tursilina	rsl	tu-i-īna	taktašifina	ktšf	ta-a-i-īna
IPFV.3PL.M	yursilūna	rsl	yu-i-ūna	yaktašifūna	ktšf	ya-a-i-ūna
PFV.1SG	ʔarsaltu	rsl	ʔa-a-tu	iktašaftu	ktšf	i-a-a-tu
PFV.2SG.F	ʔarsalti	rsl	ʔa-a-ti	iktašafti	ktšf	i-a-a-ti
PFV.3PL.M	ʔarsalū	rsl	ʔa-a-ū	iktašafū	ktšf	i-a-a-ū
...

Another pattern that is interesting from the point of view of marker extraction is Tohono O'odham (Uto-Aztecan, Mexico, USA). Some nouns in Tohono O'odham mark plural using partial reduplication of the stem (Hill and Zepeda 1998). Table 7 shows this for the two nouns *ban* ‘coyote’ and *ceoj* ‘boy’ using the phonological transcription generated by Epitran.

Table 7:
Tohono O'odham
nouns *ban*
‘coyote’ and *ceoj*
‘boy’

cell	form	stem	marker	form	stem	marker
SG	ban	ban	-	$\overline{tʃiɪndʒ}$	$\overline{tʃiɪndʒ}$	-
PL	ba:ban	ban	-:ba-	$\overline{tʃiɪndʒ}$	$\overline{tʃiɪndʒ}$	-tʃ-

Applying the automatic stem extraction for the purposes of this study, the reduplicated stem is analyzed as infixation, i.e. the marker of the plural cell occurs within the sequence shared by both cells.

3.4

Stem alternations and suppletion

All examples shown in the previous section had stems corresponding to continuous strings with no internal alternation across cells. This is not necessarily the case; alternations within stems are common across languages. Stem alternations can be defined as phonological changes within the material expressing the lexical meaning of a word across cells of a paradigm (cf. Paster 2016; Baerman and Corbett 2012).

As was mentioned in Section 2.2, such alternations do not necessarily provide meaningful insights about the inflectional properties

of a system. For inflected forms with stem alternations, the stem and marker extraction method shown in Section 3.3 would result in material being analyzed as part of the marker, which could otherwise be considered as belonging to the stem. Therefore, this method runs the risk of detecting fewer zero markers than there potentially are.

To gauge the effect of marker material resulting from stem alternations, I extracted another set of zero markers_A according to the definition given in (5) by removing material that could be analyzed as a stem alternation. To do so, I determined the position(s) of inflectional affixation for all language and part of speech combinations in the dataset. This was done based on language descriptions and based on the extracted stems and markers used in this study. Given the observed patterns, I distinguished between the following five categories of affix position: prefix, suffix, prefix + suffix, infix + suffix, prefix + infix + suffix.¹⁴ According to this classification, all material that had originally been assigned to the marker but did not occur in a regular affix position for a given language and part of speech was removed. A schematic overview of this step is shown in Table 8.

affix position	removal	marker	marker _A
pfx	remove infixes and suffixes	pfx-ixf-sfx	pfx-
sfx	remove prefixes and infixes	pfx-ixf-sfx	-sfx
pfx + sfx	remove infixes	pfx-ixf-sfx	pfx-sfx
ixf + sfx	remove prefixes	pfx-ixf-sfx	-ixf-sfx
pfx + ixf + sfx	/	pfx-ixf-sfx	pfx-ixf-sfx

Table 8:
Marker_A
extraction

Similarly to the first step of stem and zero marker extraction, these marker adjustments were automated so that they could be applied systematically for all the languages in the dataset without any additional manual annotations. Only for the type prefix + infix + suffix, no additional material could be removed from markers, because all available affix positions were in use by inflectional morphology. The three languages in this category are Arabic, Hebrew and Maltese; I applied no further changes to the markers in this case.

The following paragraphs provide a few examples of how markers_A were extracted in the presence of stem alternations. One example is

¹⁴The list of languages and affix position values can be found in `affixation.csv` in the supplementary materials.

a vowel change in German nouns, where a back stem vowel in the singular cells is opposed to a front stem vowel in the plural cells. This is shown for the German noun *Kloß* ‘dumpling’ in Table 9. All forms are given in the phonological transcription generated with Epitran.

Table 9:
Marker
extraction of the
German noun
Kloß ‘dumpling’

cell	form	stem	marker	marker _A
NOM.SG	klos	kls	-o-	-
ACC.SG	klos	kls	-o-	-
DAT.SG	klos	kls	-o-	-
GEN.SG	kloses	kls	-o-es	-es
NOM.PL	kløɐ̯	kls	-ø-ə	-ə
ACC.PL	kløɐ̯	kls	-ø-ə	-ə
DAT.PL	kløɐ̯ən	kls	-ø-ən	-ən
GEN.PL	kløɐ̯	kls	-ø-ə	-ə

In the case of *Kloß*, the longest common subsequence is not continuous. Due to the umlaut process in the plural forms, the automatically extracted stem of *Kloß* consists of the three consonants *kls*. The vowel changes from /o/ in the singular to /ø/ in the plural is analyzed as a part of the cells’ markers, respectively. Therefore, lemmas such as *Kloß* in German do not have zero marking according to the first method of marker extraction. Adjusting the markers by removing all material that is not a suffix takes into account that the alternation between /o/ and /ø/ is a stem alternation. The markers_A now no longer contain infixal material and are analyzed as zero for the nominative, accusative and dative singular cells.

Another process of stem alternation is metathesis. Table 10 shows how this is dealt with in the case of the Hungarian noun *gyomor* ‘stomach’.

Table 10:
Marker
extraction for the
Hungarian noun
gyomor ‘stomach’

cell	form	stem	marker	marker _A
NOM.SG	jomor	jomr	-o-	-
ACC.SG	jomrot	jomr	-ot	-ot
DAT.SG	jomornøk	jomr	-o-nøk	-nøk
INSTR.SG	jomor:ɔl	jomr	-o:ɔl	-:ɔl
TERM.SG	jomorig	jomr	-o-ig	-ig
ON.ESS.SG	jomron	jomr	-on	-on
ON.ALL.SG	jomor:ɔ	jomr	-o:ɔ	-:ɔ
ON.ABL.SG	jomor:ɔ:l	jomr	-o:ɔ:l	-:ɔ:l
...

In this example, the final segment *-or* is metathesized when certain affixes are added to the stem. Again, this leads to a situation where the stem does not include the segment undergoing metathesis, and the discontinuous string *jomr* is analyzed as the stem. This in turn leads to the infixal marker *-o-* in the nominative singular cell. The nominative singular is usually (81% in this dataset) not overtly marked in Hungarian. The adjusted markers_A no longer feature material that is infixal; this means the nominative singular is zero marked for the noun *gyomor* as well.

Another example of stem-internal alternations is epenthesis, the addition of phonological material in the stem in some but not all cells of the paradigm. One example of epenthesis is found with certain types of adjectives in Slovenian, which feature stem-final consonant clusters. This can be seen with the adjective *absúrden* ‘absurd’ in Table 11.

cell	form	stem	marker	marker _A
NOM.SG.M.INDEF	absúrden	absúrdn	-e-	-
NOM.SG.N	absúrdno	absúrdn	-o	-o
NOM.SG.F	absúrdna	absúrdn	-a	-a
DAT.SG.M	absúrdnemu	absúrdn	-emu	-emu
DAT.SG.N	absúrdnemu	absúrdn	-emu	-emu
DAT.SG.F	absúrdni	absúrdn	-i	-i
...

Table 11:
Marker
extraction for the
Slovenian
adjective
absúrden ‘absurd’

In all but one inflected form the stem ends in the cluster */rdn/*, and an overt suffix is added to the stem. The indefinite nominative singular masculine cell, however, is not marked by an additional suffix. Instead, the epenthetic vowel */-e-/* is inserted between the stem-final consonants to break up the consonant cluster. The adjusted markers_A remove all infixal material for Slovenian adjective markers, and the indefinite nominative singular masculine cell is analyzed as zero marked.

Another language where stem alternations are relevant in yet a different way is Tlatepuzco Chinantec (Otomanguan). Tlatepuzco Chinantec has a complex inflectional paradigm because it combines various patterns of stem and tone changes. Table 12 shows the inflectional paradigm of the verb *køgʔ²* ‘eat’. The forms of *køgʔ²* have different tones for first vs. second and third person forms in all three tenses. Given that the tones are represented by superscript numbers following the tone-bearing unit, they are taken into account by the

extraction and the detection of zero markers. While present and future tense forms do not make use of an additional segmental marker, the tone annotations are extracted as marker material. Given that otherwise, Tlatepuzco Chinantec verbs only use prefixation, I removed all infixal and suffixal material for the adjusted markers_A. As can be seen in Table 12, the adjusted markers_A now capture tonal changes as changes to the stem, and the present and future tense cells are now taken to be zero marked.

Table 12:
Marker
extraction for the
Tlatepuzco
Chinantec verb
*køgʔ*² ‘eat’

cell	form	stem	marker	marker _A
PRS.1SG	køgʔ ¹²	køgʔ	₋₁₂	-
PRS.1PL	køgʔ ¹²	køgʔ	₋₁₂	-
PRS.2	køgʔ ²	køgʔ	₋₂	-
PRS.3	køgʔ ²	køgʔ	₋₂	-
PST.1SG	mi ³ -køgʔ ¹²	køgʔ	mi ³⁻¹²	mi ³⁻
PST.1PL	mi ³ -køgʔ ¹²	køgʔ	mi ³⁻¹²	mi ³⁻
PST.2	mi ³ -køgʔ ²	køgʔ	mi ³⁻²	mi ³⁻
PST.3	mi ³ -køgʔ ²	køgʔ	mi ³⁻²	mi ³⁻
FUT.1SG	køgʔ ¹³	køgʔ	₋₁₃	-
FUT.1PL	køgʔ ¹³	køgʔ	₋₁₃	-
FUT.2	køgʔ ³	køgʔ	₋₃	-
FUT.3	køgʔ ¹	køgʔ	₋₁	-

Although this automated way of accounting for stem alternations is able to deal with almost all of the relevant cases, there is one type of alternation that this method cannot capture. If a stem alternation occurs at the edge between stem and affix, then the extraction methods used for this study are not able to detect that the boundary between marker and stem should occur in a different position.

Table 13:
Marker
extraction for the
Northern Saami
adjectives
aiddolaš ‘exact’
and *bahá* ‘angry’

cell	form	stem	marker _(A)	form	stem	marker _(A)
NOM.SG	aiddolaš	aiddola	-š	bahá	bahá	-
ACC.SG	aiddolačča	aiddola	-čča	bahá	bahá	-
GEN.SG	aiddolačča	aiddola	-čča	bahá	bahá	-
ILL.SG	aiddolažžii	aiddola	-žžii	bahái	bahá	-i
COM.SG	aiddolaččain	aiddola	-ččain	baháin	bahá	-in
FRML.SG	aiddolažžan	aiddola	-žžan	bahán	bahá	-n
PRP.SG	aiddolaččas	aiddola	-ččas	bahás	bahá	-s

One example is the so-called consonant gradation in Northern Saami. It can be described as an alternation of the final stem consonants

across cells of the paradigm, leading to their weakening or strengthening (cf. Bakró-Nagy 2022). An example from Northern Saami adjectives is shown in Table 13. We see that the final stem consonant of the adjective *aiddolaš* ‘exact’ alternates between /-š/, /-čča/ and /žž/. The extraction process used here analyzes this alternation as part of the marker. The adjective *bahá* ‘angry’, on the other hand, shows the marker extraction for adjectives with no stem alternations. For such adjectives, the nominative, accusative and genitive singular cells are zero marked. Thus, in cases of alternations at the edge between the stem and the inflectional affix, this method of marker extraction is unable to detect zero marking.

In its most extreme form, a stem alternation that includes the edge segments of stems is suppletion. Suppletion refers to stem alternations where maximally different phonological forms are used to express the same lexical component of an inflected word form across different cells of the paradigm (cf. Mel’čuk 1994; Corbett 2007). Suppletive forms go beyond alternations that can be described in terms of phonological or prosodic relations between forms (at least synchronically). Consider the English examples given in Table 14, where we see the verbs *think* and *go*, both with suppletive stems. In the case of *think*, the suppletion does not affect the entire stem, as the initial segment *θ-* is found in all cells of the paradigm. As a consequence, the extracted marker ends up with all the remaining material (which would usually be analyzed as being part of a suppletive stem). In the case of *go*, suppletion is complete in that no segment is shared between all cells of the paradigm. The complete phonological strings of each form are thus extracted as markers of their respective cells.

cell	form	stem	marker _(A)	form	stem	marker _(A)
NFIN	θɪŋk	θ	-ɪŋk	gow	-	gow
PST	θɔt	θ	-ɔt	went	-	went
PTCP.PST	θɔt	θ	-ɔt	went	-	went
PTCP.PRS	θɪŋkɪŋ	θ	-ɪŋkɪŋ	gowɪŋ	-	gowɪŋ
PRS.3SG	θɪŋks	θ	-ɪŋks	gowz	-	gowz

Table 14:
Marker
extraction for the
English verbs
think and *go*

As the examples from Northern Saami and English showed, neither marker extraction method used for this study has a principled way of removing alternating stem segments that are adjacent to affixal material from the marker. Therefore, neither method detects potential

zero marking with suppletive forms, as it will always assign phonological material to the marker. While it is possible to exclude markers that occur only once per cell (cf. Section 3.5), many suppletive forms do not correspond to such hapax legomena markers. Especially larger datasets often include complex lemmas such as *overthink* or *undergo* in the case of English. For instance, the extracted markers *-gow* and *-ɪŋk* from Table 14 occur 11 times in the verbal paradigms of English. Also the stem alternation pattern shown for Northern Saami in Table 13 occurs systematically (26 times) in the dataset. In such cases, I do not have any principled way of excluding markers from the analysis.

To remain agnostic about the effect of stem alternations and to apply a systematic approach to all languages, I performed the analyses in Sections 4 and 5 for both sets of markers and markers_A. Since the results are very similar with no substantial differences, I only report the results of using markers_A for reasons of brevity. Details about the results based on the originally extracted markers can be found in the supplementary materials as indicated in the respective sections. Given that no substantial differences were found for the distribution of zero markers in inflection paradigms, I only use the markers_A set for the analysis of their distribution in corpus data in Section 6. Whenever markers are mentioned in the following sections, I refer to markers_A if not stated otherwise.

3.5

Hapax legomena markers

The dataset includes a number of markers that occur only once per cell for a given language and part of speech combination. Some of these hapax legomena markers are the result of stem alternations, but most of them go back to remaining errors in the dataset. In total, I identified the following number of hapax legomena markers: 9223 for adjectives, 23539 for nouns and 54768 for verbs. In terms of marker types, hapax legomena markers make up a large proportion, namely 0.45, 0.46 and 0.42 for adjectives, nouns and verbs, respectively. In terms of the total number of occurrences, however, they only amount to a proportion of 0.003 for adjectives, 0.008 for nouns and 0.03 for verbs.

One example of a hapax legomena markers as the result of stem alternation comes from Northern Saami. The adjective *čáppat* ‘pretty’

features gradation similarly to the example shown in Table 13. In this case, stem-final *-pp* alternates with *-bb* across cells of the paradigm. This type of alternation is only attested once in the dataset, making all markers extracted from the lemma *čáppat* hapax legomena markers.

Most hapax legomena markers, however, result from remaining material that is not part of the inflected word forms or from errors in the automatic phonological transcription performed by EpiTran. To give one example, in the Hungarian dataset, the impersonal verb *fái* ‘hurt’ features the annotation of ‘only3rdpersonforms’ as the verb form with some cells. Such linguistic material that does not belong to the word forms causes the extraction of the longest common substring to find non-sensical strings and hence hapax legomena markers.

Visual inspection of the hapax legomena markers suggests that most result from the automatic phonological transcription using EpiTran. For instance, the German adjective *makaberə* ‘macabre’ shows an alternation between stem-final *-b* and *-p* in the phonological transcription. All forms except the comparative form have *-b*, while the comparative form *makaberə* has *-p*, which leads to hapax legomena markers.

In order to exclude such markers that do not allow for much insight on the distribution of zero, I removed all hapax legomena markers from the dataset. Given that their proportions of the total number of observations is very low, it is safe to assume that their removal will not artificially distort the distribution of zero markers.

Morphomic paradigms

3.6

Another potential factor influencing the distribution of zero marking is the distribution of inflected word forms across the paradigm. Many paradigms have syncretic cells, where a single form expresses more than one cell. Taking this into account and considering only the different forms that are found in a paradigm may thus lead to different probabilities of zero markers. To examine how much the results change if proportions of zero marking are established using distinct forms only, I collapsed the data into morphomic paradigms (cf. Boyé and Schalchi 2016). Morphomic paradigms consist of the different forms that a given word can have without taking into account

their meaning. Syncretic forms are counted in only once in morphomic paradigms. Section 4 therefore analyzes the distributions of markers in morphomic paradigms in addition to paradigms that include information on cells. The analysis of the effect of token frequency in language use on the distribution of zero marking in Section 6 is also based on forms exclusively.

4 ESTIMATING THE PROBABILITY OF ZERO MARKERS

4.1 *Observed distributions*

In order to examine the probability of zero markers in adjectival, nominal and verbal inflection, Table 15 and Figure 2 provide an overview of the observed distribution of zero marking in inflection.

Table 15:
Observed
proportions of
zero markers

pos	N zero	prop zero	N(prop ₀ = 0)	N(prop ₀ = 1)
adjective	45859	0.007	1439	12
noun	648859	0.104	1227	5
verb	141268	0.032	3771	26

From Table 15, we see that the proportions of zero markers are very low for adjectives; verbs show a somewhat higher proportion and nouns have the highest proportions of zero marking with 0.1. Zero marking is clearly not common in inflection of any of the parts of speech. The last two columns of Table 15 show the number of cells with proportions of zero marking that equal the two extremes 0 and 1. The proportions that equal 0 correspond to cells where zero marking does not occur. Proportions of 1 mean that a given cell is exclusively expressed by zero markers in this dataset. Unsurprisingly, we find a high number of proportions that equal 0 and a very small number of cells with proportions of 1.¹⁵ For both types of proportions, we find

¹⁵The number of 26 cells that are expressed by zero markers exclusively is rather high; this can in part be explained by many cells in the verbal paradigm that only occur in single languages.

an increasing number from nouns to adjectives to verbs. This reflects the number of cells that those three parts of speech distinguish in the dataset, with 727, 961 and 2753 cells for nouns, adjectives and verbs, respectively.

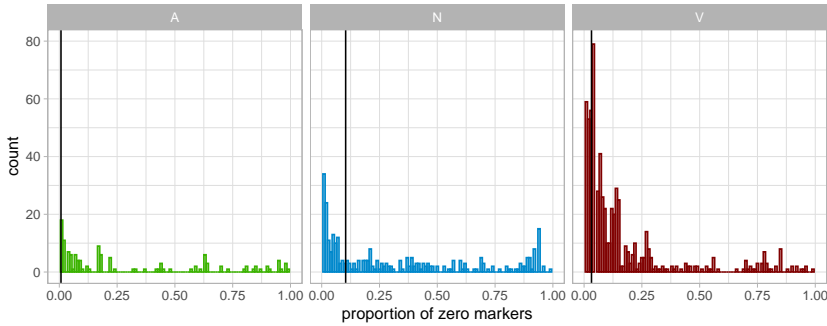


Figure 2:
Observed
proportions of
zero markers

Figure 2 shows a histogram of the proportions of zero marking in adjectival, nominal and verbal inflection. The overall proportions are indicated by a vertical line. We can see that they vary to a great extent across languages and part of speech. All three parts of speech exhibit a preference for proportions of 0 or close to 0. This preference is most pronounced for adjectives and verbs. For nouns, we find a more balanced distribution, with more proportions > 0.5 for zero marking.

There are five additional factors that are relevant for estimating the probability of zero markers in inflection: the number of cells that a paradigm has, the number of values expressed per cell, the number of lemmas for which paradigms are available, the usual affix position and the type of phonological representation.

	min	max	mean	sd
adjective	3	256	44.4	54.6
noun	2	256	25.8	38.7
verb	2	432	49.5	70.1

Table 16:
Number of cells

The number of cells in a paradigm can be taken as a measure of paradigm size. It is an important factor to include, since it is possible that zero markers are less likely to occur in a larger paradigm that makes more morphosyntactic distinctions. Table 16 gives an overview of the number of cells per paradigm in the dataset, showing the minimum, maximum, mean and standard deviation. Because the number

of cells spans several magnitudes, I use log-transformed values for the analysis.

Another important factor for estimating the probability of zero marking is the number of values expressed per cell. For the purposes of this study, we can take the number of values per cell to represent the semantic complexity of the inflectional markers. A summary of the number of values per cell is shown in Table 17. Including this factor in the analysis is important, since one could expect that more complex markers (which express more complex meanings) tend to be encoded by more material.

Table 17:
Number of
values per cell

	min	max	mean	sd
adjective	1	5	2.92	1.01
noun	1	4	2.02	0.587
verb	1	7	2.05	0.887

The average number of lemmas for which inflectional paradigms are available is not inherently related to the probability of zero marking, but it may influence it. As can be seen in Table 18, the average number of lemmas differs greatly across languages. It is therefore an important factor to be controlled for.

Table 18:
Number of
lemmas

	min	max	mean	sd
adjective	17	98464	4035	15666
noun	23	235294	9843	33158
verb	26	30032	1348	3438

Another factor that is included in the analysis for its potential effect on the probability of zero marking is the position of the marker regarding the stem.

Table 19:
Affix position

	pfx	pfx + sfx	pfx + sfx + ifx	sfx	sfx + ifx
adjective	36	259	48	1365	0
noun	8	84	62	1436	2
verb	407	889	164	3093	8

As described in Section 3.3, I distinguish between five affix positions found in the dataset. Table 19 shows the number of cells per part of speech expressed by markers in the five positions. For the analysis, I

merged the two positions that include infixes into one, because the *sfx+ifx* category on its own has too few observations to allow for any meaningful insights. This leaves the following four values of affix position that are considered in the analysis: *pxf*, *pxf+sfx*, *sfx* and *has_ifx*.

Modelling the probability of zero marking

4.2

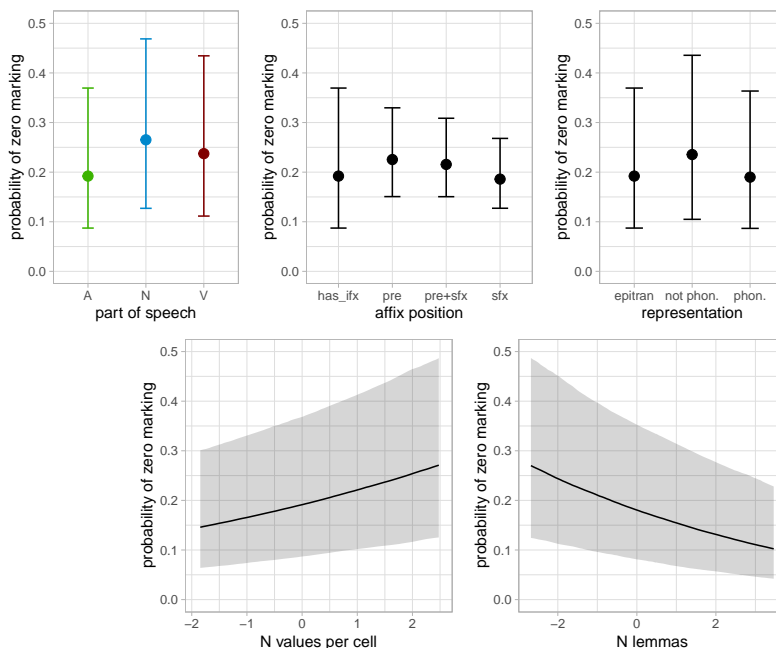
To estimate the probability of zero marking in inflection, I aggregated the data by type of cell, language and part of speech. This means that each datapoint corresponds to a proportion of zero marking (0.81) for a given type of cell (NOM;SG) in a given language (Hungarian) for a given part of speech (noun). As was shown in Table 15, the dataset contains cells with proportions of zero marking that equal 0 or 1. Therefore, I fitted a Bayesian zero-one-inflated regression model. Zero-one-inflated regression models consist of two components. The first component is the regular beta regression model, which deals with proportion values within the interval (0,1). The second component is a logistic regression component that estimates the probability of either of the extremes 0 or 1 as opposed to the proportion data within (0,1).

The models were fitted using Stan (Carpenter *et al.* 2017) with the *brms* package (Bürkner 2017) in R (R Core Team 2021). I additionally controlled for the phylogenetic relations between languages using a phylogenetic regression term following the method described in Guzmán Naranjo and Becker (2022). This term does not model the relations between languages in a categorical way but includes the information of the entire phylogenetic tree and forces the estimates of the single languages to co-vary according to the tree.¹⁶ In other words, if two languages share many nodes of the tree, the model forces their coefficients to be very similar. If two languages are not related at all, the model allows their estimates to vary freely. For instance, if five closely related languages have very high observed proportions of zero markers in a given cell, the model does not take those five observations as independent data points and assigns much less confidence and/or lowers the predicted probability of zero marking in that cell.

¹⁶The phylogenetic tree is taken from Glottolog (Hammarström *et al.* 2021). For details, see `code-phylogeny.R` in the supplementary materials.

The final model predicts the probability of zero marking from the part of speech, affix position, number of values per cell, number of lemmas and the orthographic representation. In addition, I used the type of cell and the phylogenetic relations between languages as group-level effects.¹⁷

Figure 3:
Conditional
effects for the
beta regression
component



Figures 3 and 4 show the conditional effects for the different predictors for the beta and the zero-one-inflation components, respectively.¹⁸

¹⁷To select a reasonable combination of predictors, I fitted several models and compared their performance using approximated leave-one-out cross-validation as described by Vehtari *et al.* (2017). Also, given the low number of proportions of 1, I modelled the conditional one inflation with an intercept-only model. Therefore, I do not discuss conditional one inflation further in this section. See `code-prob.R` in the supplementary files for details.

¹⁸I only report the results of the model based on markers_A that allow for stem alternations. All conditional effects of the model based on markers without stem alternations can be found in `ce-probcheck-mu-<predictor>.pdf` and `ce-probcheck-zoi-<predictor>.pdf` in the supplementary materials.

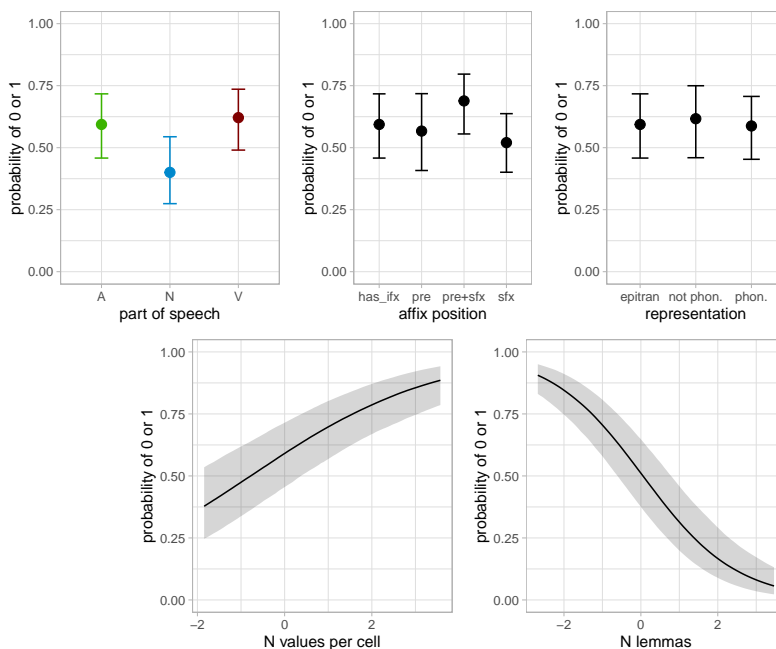
The points and solid lines correspond to the mean values of the posterior distributions; the error bars and error bands show the 95% credible interval. They allow for a straightforward interpretation: given the data and the model, we can be 95% certain that the estimated values lie within that interval. Note that the three numeric predictors are all standardized, so that they have a mean of zero and a standard deviation of 1.

From Figure 3, we see that none of the predictors has a clear impact on the probability of zero marking within the interval (0,1). Across all predictors, the mean predictions lie between 0.15 and 0.3. The results thus show that the probability of zero marking to occur, excluding systematic absence or presence thereof, does not depend much on the predictors explored here. This does not necessarily mean that a better model is needed. It suggests that there is a high degree of idiosyncratic variation across languages, and that no clear association can be drawn to other relevant grammatical properties of the inflectional systems.

Figure 4 shows the model results for the zero-one-inflation component. It predicts the probability of a cell being exclusively zero marked (1) or never zero marked (0) as opposed to probability values in between those two extremes. As was shown in Table 15, no zero marking per cell is common in the data (6437 markers out of 7861), while exclusively zero marked cells are very rare (43 markers out of 7861). This means that zero-one-inflation predictions largely correspond to the probability of no zero marking for a given cell. We can thus interpret the conditional effects shown in Figure 4 as the probability of a very strong trend against zero marking. For the predictors part of speech, affix position and phonological representation, we find no substantial trends regarding a preference against zero marking. For part of speech, adjectives and verbs appear to have a slightly higher probability than nouns to avoid zero marking, but we have little certainty about this difference. The same can be said about the affix order *pfx + sfx*; it has a slightly higher preference to avoid zero marking than the other positions, but no clear picture emerges.

In contrast to the predictions from the beta component, we do find clear effects of the number of values per cell and the number of lemmas. The more lemmas are available, the lower the probability to encounter no single case of zero marking. This is expected and shows

Figure 4:
Conditional
effects for the
zero-one-
inflation
component



that the number of lemmas needs to be controlled for. The number of values per cell has a positive effect on the probability of avoiding zero marking altogether. While cells expressing fewer values show no strong preference for or against zero marking, the model predicts a strong preference against zero marking for cells with many values. This does not restrict where zero marking is likely to occur, but it predicts the total absence of zero marking for complex cells with a high probability of 0.8.

As was mentioned in Section 3.6, it is important to consider the distribution of zero marking in morphomic paradigms as well. I fitted another Bayesian zero-one-inflated regression model using morphomic paradigms with the same predictors as described above. Only the predictors including information on cells (cell, number of values per cell) are no longer included. The model results are generally very similar to one fitted on full paradigms. Therefore, I only show those results that differ and provide new insights.¹⁹

¹⁹See the file `code-morphomic.R` for details. The conditional effects for all predictors of the model using morphomic paradigms are found in

Zero marking in inflection

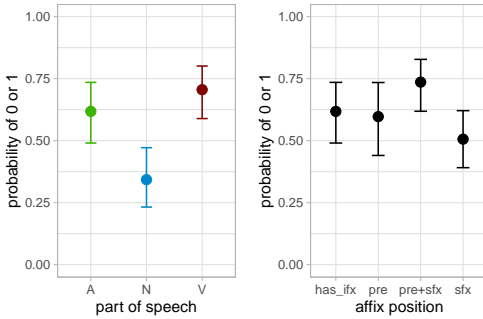


Figure 5: Conditional effects for the zero-one-inflation component of morphomic paradigms

The predictions from the beta regression component are similar to the ones of the full paradigms. The overall predicted probability of zero marking is just below 0.2, which is slightly lower than in full paradigms. This suggests that zero marking is syncretic in a portion of the dataset. As the credible intervals are very wide in both models and overlap, we cannot be very certain about this finding.

For the zero-one-inflation component of the model, the conditional effects of part of speech and affix position allow for additional insights. The model predictions for those two variables are shown in Figure 5. We see that the patterns are similar, only that the differences between parts of speech are much stronger now. With morphomic paradigms, we can be certain that verbs and adjectives have a stronger tendency than nouns to avoid zero altogether. The same holds for affix position. Figure 5 shows that systems with prefixes and suffixes are more likely to avoid zero marking altogether than systems with suffixes only.

FUNCTIONS ASSOCIATED WITH ZERO MARKING

5

Cells with the highest probability of zero marking

5.1

To explore which cells are most likely to be zero marked, I subsetted the dataset to include only those cells with a proportion of zero forms

the supplementary materials as `ce-probmorph-mu-<predictor>.pdf` and `ce-probmorph-zoi-<predictor>.pdf`.

≥ 0.1 in at least 10% of the languages. Subsetting the data in such a way was necessary because of the high number of cell types. The threshold is a heuristic, chosen to restrict the following analysis to the cells with a reasonable crosslinguistic probability of being expressed by zero markers. It leaves in 18 types of cells that show the strongest association with zero marking in the observed distributions.²⁰

In order to estimate the probability of zero marking in these cells, I fitted a Bayesian beta regression model that predicts the probability of zero marking from the type of cell.²¹ In addition, I added the number of values per cell and lemmas as group-level intercepts as well as phylogenetic controls to account for phylogenetic biases in the data.

Figure 6 shows the observed proportions of zero forms (black triangles) together with the model predictions (dots, error bars and error bands).²² Again, the dots represent the mean values of the posterior distribution of the zero probabilities, and the error bars and bands show the 95% credible intervals. The observed proportions of zero forms still differ across cells and parts of speech, ranging from 0.1 (2SG present verb forms and dative singular adjectives) to above 0.7 (indefinite singular nouns). Although adjectives have fewer cells that met the threshold criteria than nouns and verbs, Figure 6 shows that the cells that do meet them have comparatively high proportions of zero marking. In nominal cells, we find a wider range including the highest overall proportions of zero marking. Verbs show the lowest proportions of zero marking compared to the other parts of speech.

When comparing the results of the model with the observed proportions, the predicted probabilities of zero markers reflect the observed proportions for the most part. Figure 6 shows a few differences, though. For some cells, the predicted probability is much lower than

²⁰The exact figures, including the number of languages per cell, are found in `cells-merged.csv` in the supplementary materials.

²¹In this case, I used beta regression instead of zero-one-inflated beta regression for a combined prediction from both processes. To do so, I converted proportions of zero to 0.0000001 and proportions of 1 to 0.9999999. Again, I compared several models using approximated leave-one-out-cross-validation. See `code-cells.R` in the supplementary materials for details.

²²All conditional effects of the model based on markers without stem alternations can be found in `ce-cells-check-<predictor>.pdf` in the supplementary materials.

Zero marking in inflection

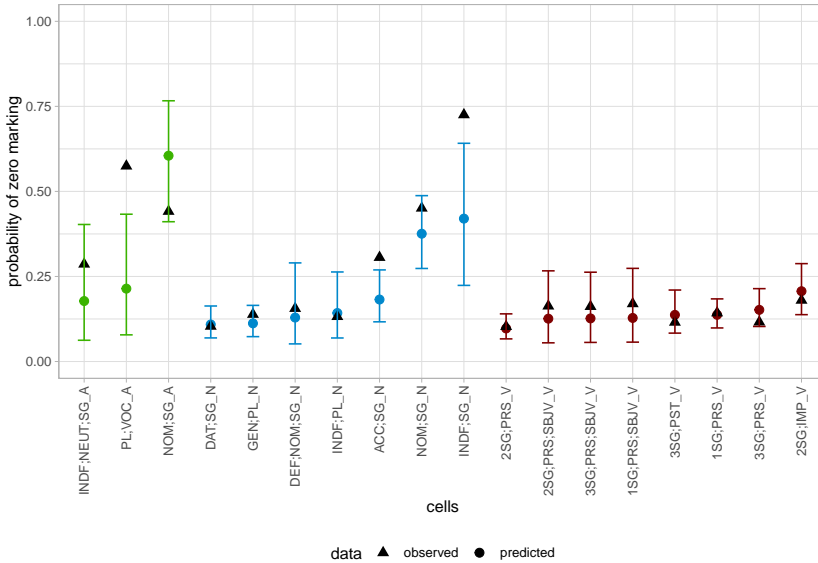
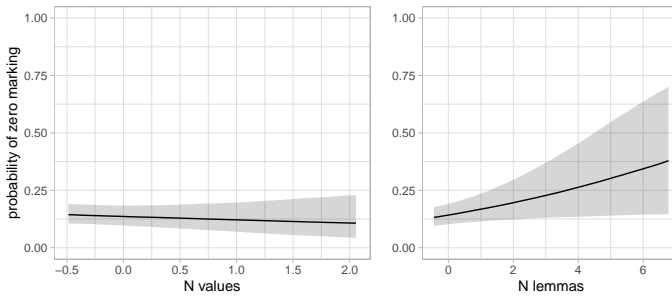


Figure 6: Conditional effects for cells most strongly associated with zero marking



their observed proportions, namely for PL;VOC of adjectives, as well as ACC;SG and INDF;SG for nouns. This points to a bias in the observed distributions, which is also reflected in the large credible intervals of the predictions. The PL;VOC cell is featured in four languages of the dataset, namely in Czech, Georgian, Irish and Sanskrit. In this case, the high proportion of zero marking is mainly an artefact of the data. The PL;VOC cell is exclusively zero marked in the Czech data. Irish has a low proportion of zero marked PL;VOC cells (0.22), and Georgian as well as Sanskrit do not feature zero marking for the VOC;PL cells of adjectives. Thus, in this case, the high overall proportion largely comes from a single language, which is then adjusted to a much lower prediction in the model, together with large credible intervals to in-

dicating the high level of uncertainty. A similar explanation applies to the ACC;SG cell of nouns. It is featured in 26 languages in the dataset, including phylogenetically unrelated languages. However, the higher observed proportion of zero marking is due to high proportions in a few, mostly related, languages with large datasets.²³ For the INDF;SG cell, the lower predicted probability of zero marking is also the consequence of a bias in the observed proportions. Here, the bias comes from Norwegian Bokmål, which makes up more than 50% of all observations for this cell, and which has a very high proportion (0.88) of zero marking.

Comparing the predictions across cells and parts of speech, we see that the adjectival cells have a very high probability of being zero marked. This is noteworthy, as adjectives had only very few cells that made the threshold to begin with. While generally not associated with zero marking, those adjectival cells that are zero marked appear to be the ones with the strongest association with zero marking across parts of speech. Nominal cells are generally predicted to have lower probabilities of zero marking, except for the NOM;SG and the INDF;SG cells, which rank second and third for the predicted probability of zero marking. All verbal cells range between 0.1 and 0.25 for the probability of zero marking. The cell that stands out for having the highest probability of zero marking is the 2SG imperative cell, which will be taken up in the discussion in Section 7.2.

5.2 *Values with the highest probability of zero marking*

The fact that the languages in the dataset differ with respect to the combinations of values in single cells makes it somewhat difficult to assess the association between zero marking and cells that are less common in the dataset. It is therefore important to consider the association of single grammatical values and zero marking as well. Note that due to how zero markers were extracted, pulling apart the values of cells and analysing their association with zero marking does not translate directly into the traditional analysis of an abstract feature

²³This includes German (0.77), Old English (0.50), Finnish (0.37), Russian (0.35), Ukrainian (0.23), Polish (0.22) and Serbian-Croatian-Bosnian (0.30).

value, e.g. singular, as being zero marked. Rather, the singular value being expressed by a zero marker refers to all cells in the dataset that encode singular (potentially besides other feature values) and that are zero marked.

In order to examine the association of single values with zero marking, I applied a similar threshold heuristic as in Section 5.1 to select those values that show the strongest association with zero marking. I only included values with an overall proportion of zero marking ≥ 0.02 that are featured in 10% of the languages per part of speech. This led to the selection of 21 values in total.²⁴ To assess how robust the observed proportions of zero marking are, I fitted a Bayesian beta regression model, adding a phylogenetic control and the number of cells and lemmas as group-level effects.²⁵

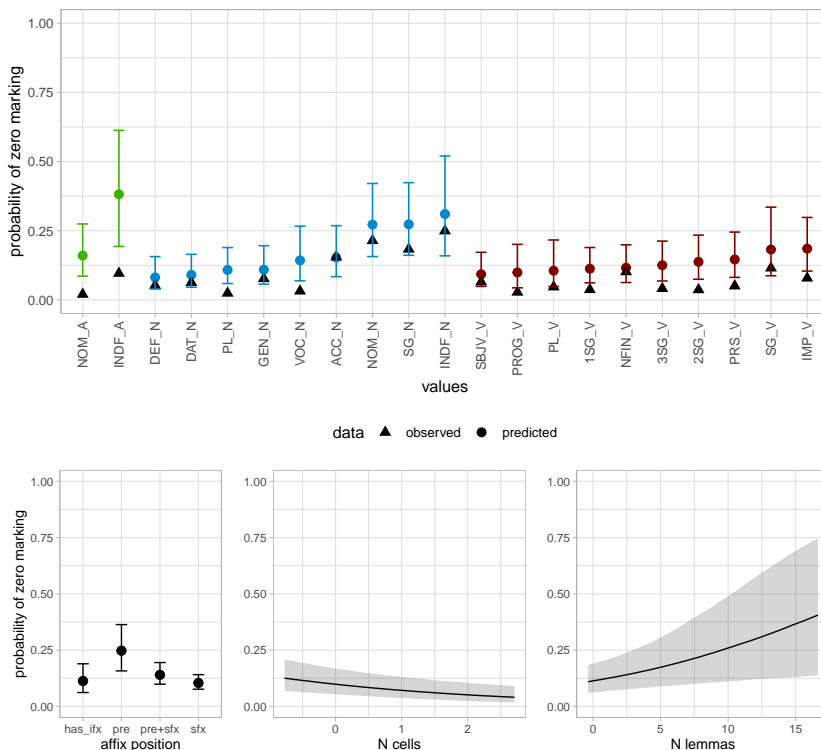
Figure 7 shows the observed proportions (triangles) together with the model predictions (dots, lines).²⁶ The dots represent the mean values of the posterior distribution of the zero probabilities; error bars and bands indicate the 95% uncertainty intervals. The distributions in Figure 7 mostly mirror the tendencies seen in Figure 6 in the previous section. Almost all values that make the threshold and are thus the values with the highest proportions of zero marking have also been part of the cells most likely to be zero marked. Only the nominal value of vocative and the verbal values of progressive, plural and non-finite have not been part of the cells most associated with zero marking. Compared to cells, values show much lower absolute proportions of zero marking. This is expected, since single values potentially occur in many different contexts, not all of which are necessarily zero marked. As for the three parts of speech, we now see the highest proportions for nominal values. Adjectival and verbal values show lower proportions of zero marking.

²⁴The exact figures, including the number of languages per value, are found in `values-merged.csv` in the supplementary materials.

²⁵I used the same method as for the model described in Section 5.1. See `code-values.R` in the supplementary materials for details.

²⁶All conditional effects of the model based on markers without stem alternations can be found in `ce-values-check-<predictor>.pdf` in the supplementary materials.

Figure 7:
Conditional effects for the values most associated with zero marking.



Turning to the model predictions, we see that in the case of values, the probability of zero marking is generally estimated by the model to be higher than the observed proportions. This can be explained by the fact that the model takes into account information on the affix position, the number of cells and the number of lemmas. The effects of single values thus correspond to their effects once all the other predictors are controlled for. Interestingly, the affix position is also relevant in this case. The model predicts a higher probability of zero marking for systems with prefixes as opposed to those with suffixes.

The highest predicted probabilities of zero marking are found for the indefinite value in adjectival and nominal inflection. This mirrors the model results of cells shown in Figure 6. Other values with a comparatively high probability of zero marking are singular and nominative for nouns, as well as imperative and singular for verbs. These results also reflect the tendencies seen with cells in Section 5.1.

THE FREQUENCY OF ZERO MARKERS IN LANGUAGE USE

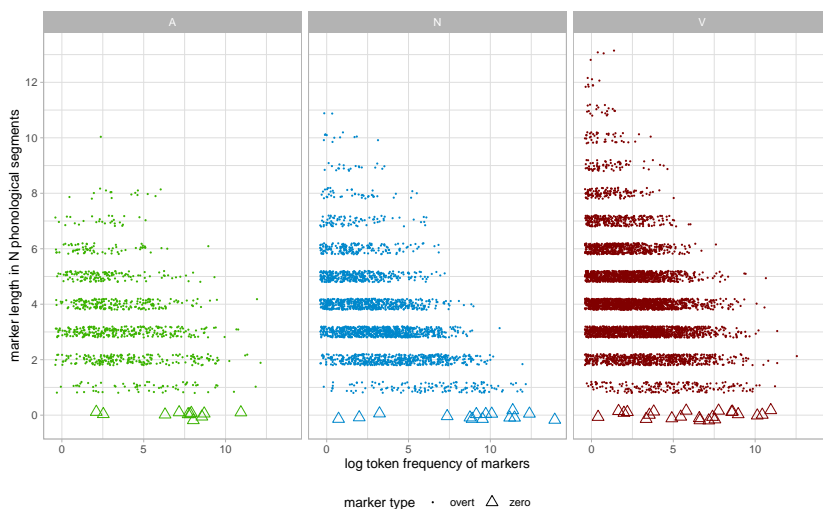
6

To assess the usage frequencies of inflection markers and their phonological length including zero, I analyzed the distribution of zero markers in the Universal Dependencies treebanks (UD) (Zeman *et al.* 2023). To do so, I merged the adjective, noun and verb forms in UniMorph identified as zero forms with the Universal Dependencies data. I only included the languages for which a phonological transcription was available, so that marker length could be approximated in a more realistic way. From the original dataset, 20 languages have phonological transcriptions and are represented in UD. For merging UniMorph forms with forms in UD, I did not include cell information and merged the forms purely based on their orthographic representation. The identification of zero markers, however, was based on the phonological transcriptions and the marker_A extraction as described in Section 3. The resulting dataset contains 9975 types of markers, which are made up of 51 types of zero markers (across different language and part of speech combinations) and 9924 distinct types of overt markers. In terms of token frequencies, zero markers make up 23% of all the marker occurrences (7382497 tokens in total). For the purposes of this study, the distribution of zero and overt markers in UD is measured by their log-transformed token frequencies. The length of the markers corresponds to the number of phonological segments identified with the UniMorph dataset. Figure 8 shows the relation between log token frequencies and marker length for adjectives, nouns and verbs. Overt markers are shown as dots, and zero markers indicated by triangles.

As expected, Figure 8 shows a consistent tendency across the three parts of speech for more frequent markers to be shorter. For less frequent markers, however, there does not seem to be a strong tendency to be longer; we also find many infrequent markers that are short. As for zero markers, Figure 8 does not show clear tendencies either. For adjectives and nouns, they appear to have comparatively high frequencies, whereas no such trend is apparent for verbs.

To test the association shown in Figure 8, I fitted a Bayesian hurdle Poisson model, predicting the marker length from their frequencies. Similarly to the zero-one-inflated beta models, a hurdle poisson

Figure 8:
Association
between marker
token frequency
and length



model consists of two components. The poisson component predicts count data, and the hurdle consists of a logistic regression component that predicts the probability of zero. This allows us to compare the effect of frequency on marker length between zero and overt markers.

In order to determine which predictors other than token frequency should be included, I fitted a series of 9 models that included different combinations of token frequency with part of speech, affix position and the number of cells. The performance of these models was then compared to select the final model. I used approximated leave-one-out cross-validation for the comparison following the method described by Vehtari *et al.* (2017).²⁷ The final model includes token frequency and affix position as well as their interaction and the phylogenetic control.

Figure 9 shows the conditional effects for the Poisson components, i.e. the part of the model that predicts the marker length. We find a clear negative effect of marker frequency, confirming previous results from the literature. On average, low frequency markers are predicted to be about 0.15 phonological segments shorter than high frequency markers. The position of the affix also proves relevant for marker length. Despite the effect being smaller, the model predicts a

²⁷ See code-ud.R in the supplementary materials for details.

substantial difference in marker length between systems only using suffixes and all other systems. This becomes more evident when considering the interaction between token frequency and affix position. The effect of frequency is greater for systems using only suffixes than for all other systems, reaching an average difference of 0.25 phonological segments between low and high frequency markers. We can thus conclude that suffixes are more sensitive to the effect of marker frequency than the other types of affixes.

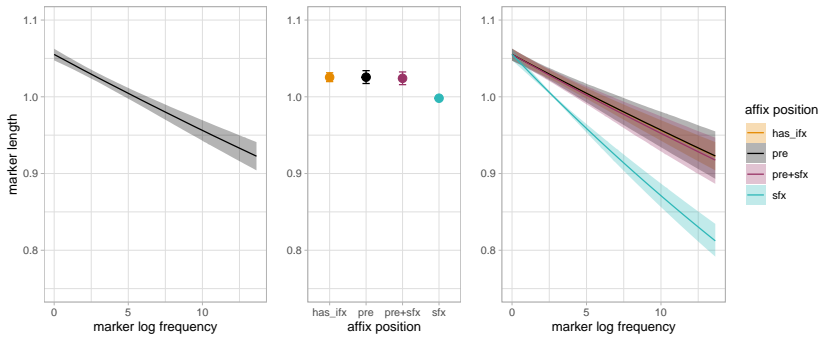


Figure 9: Conditional effects for the Poisson component

We see the conditional effects for the hurdle component in Figure 10. They represent the effect that the predictors have on the probability of a zero marker to occur.

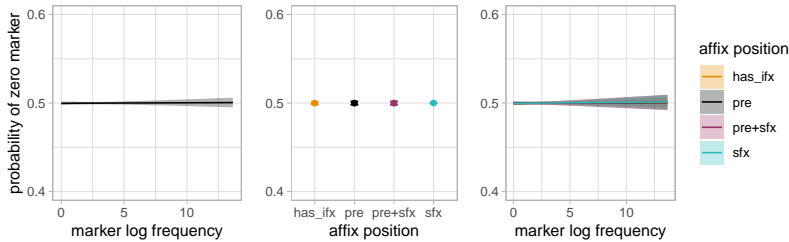


Figure 10: Conditional effects for the hurdle component

In stark contrast to the effects predicted for the phonological length of markers, neither the token frequency nor the affix position affect the probability of a zero marker. The small credible intervals show that this is not an issue of uncertainty or too few observations. We can be very confident in the model results that, given the data, the probability of zero marking to occur is not associated with the token frequency of that marker or the affix position the system uses. This means that

there is indeed a clear difference between the effect of frequency on marker length in general and the occurrence of zero marking. Zero marking does not simply follow the general trend of marker length being associated with marker frequency.

7

DISCUSSION

7.1

The probability of zero marking

The results of this study allow for a number of important insights into crosslinguistic trends of zero marking in inflection. The model results predicting the probability of zero marking in inflectional paradigms (Section 4) showed three important points. First, zero marking generally affects adjectives, nouns and verbs fairly equally and is not sensitive to the affix position(s) used for inflection. The overall probability of zero marking being rather low (0.1-0.3), zero marking is not a default strategy for inflection. This finding provides quantitative support for the proposal by Stolz and Levkovych (2019, 373), who argue that zero marking in inflection should be treated as a “morphological mismatch on a par with established categories such as suppletion and syncretism”. Zero marking is generally not a common strategy to encode inflection.

Second, we saw an effect of part of speech and affix position when analyzing zero marking in morphomic paradigms. Based on forms only with no information about cells, zero marking was more likely to be absent altogether in adjectives and verbs as opposed to nouns. The same applied to systems with prefixes and suffixes as opposed to suffixes only. This does not mean that nouns and systems with suffixes have a stronger preference for zero marking. It rather suggests that zero marking is less systematically excluded in those cases.

Third, an increasing number of values per cell was shown to be a strong predictor for a high probability of zero marking being avoided altogether. The predictor number of values per cell quantifies how functionally complex a marker is. The fact that more complex cells strictly avoid zero marking is reminiscent of what has been discussed

as isomorphism or iconicity in the literature (cf. Haspelmath 2008b; Lehmann 1974; Downing and Stiebels 2012; Givón 1991). While approaches differ in their details, the general idea is that the complexity or amount of linguistic structure reflects the complexity or amount of functional structure (meaning). It remains an open question, however, what the functional motivation for this effect is, if there is one. It is likely that usage distributions and frequencies are a confounding factor in that cells expressing more values may also be cells that are used less frequently. Their preference for longer markers could thus be a consequence of frequency rather than some iconicity principle.

Cells and values associated with zero marking

7.2

Sections 5.1 and 5.2 focused on a selection of cells and grammatical values and their association with zero marking. The results showed that even though zero marking exhibits a high degree of variation across lemmas and languages, it is not distributed randomly across inflectional paradigms. Some cells and values are comparatively likely to be zero marked across languages. For adjectives and nouns, indefinite, nominative and singular (and cell combinations thereof) were the values with the highest predicted probability of zero marking. For verbs, the probabilities of zero marking tended to be generally lower. The values of imperative, singular, third person and present (and cell combinations thereof) stood out as the ones with the highest probability of zero marking. The NOM;SG cell for adjectives was the only cell for which the probability of zero marking was predicted to be above 0.5. In other words, this is the only cell for which we can expect zero marking to be more likely than overt marking. In all other cases, predicted probabilities lied well below 0.5. This means that the vast majority of inflectional marking is in fact overt, and zero marking is more of an exception than the rule.

The values of nominative and singular, as well as their combination have long been associated with zero marking in the typological literature (e.g. Greenberg 1963, 1966; Jakobson 1983; Koch 1995; Croft 2003; Haspelmath and Karjus 2017; Haspelmath 2021). Interestingly, there is less discussion in the literature about zero marking of the indefinite value, which showed the strongest trend towards zero

marking in this study. Two verbal values that have been related to zero marking in the literature are third person (Bickel *et al.* 2015; Cysouw 2003; Siewierska 2010) and present tense (Bybee and Dahl 1989, 55; Bybee 1994, 248). The results of this study confirm the association. Although neither values show a crosslinguistic preference towards being zero marked, they are part of the values with the highest probabilities of zero marking.

Imperatives, especially second person singular forms, have also been mentioned in the literature as being prone to zero marking (e.g. Greenberg 1966; Koch 1995; Croft 2003; Haspelmath 2021; Aikhenvald 2010; Siewierska 2010). The results of the present study thus fit in well with the expectations from the literature. Instead of phonetic reduction, previous studies have argued for a functionally motivated non-development scenario for zero marking in (second person singular) imperatives. The idea is that the second person is highly recoverable in imperative contexts, e.g. as opposed to contexts of indicative verb forms. Thus, on the level of syntax, many languages allow or require the use of imperatives with no overt second person subject pronoun. This in turn means that the source construction of a verbal person marker is often not available for imperative forms (Aikhenvald 2010, 147; Nikolaeva 2007, 163; Sadock and Zwicky 1985, 173). The crosslinguistically common absence of a suitable source construction for person markers in imperative contexts may thus ultimately account for the high probability of zero marking especially for person-number agreement values. In addition, the use of bare verb forms for imperatives has been motivated by iconicity (Aikhenvald 2010, 46). According to her, using the shortest verb form makes imperatives very direct and abrupt. This can convey urgency and reflect that imperatives usually call for an immediate reaction.

7.3

Frequency effects and affix position

Section 6 examined the association between the token frequency of inflection markers and their length, including zero marking. For overt markers, the present study could confirm Zipfian effects for inflectional markers. Markers with a higher log frequency were predicted to

have longer forms (in number of phonological segments). This corroborates previous findings about form-frequency effects for inflectional markers by Haspelmath and Karjus (2017) and Stave *et al.* (2021).

An aspect that has not been addressed in quantitative corpus studies so far is the effect that the position of the inflection marker has. The results from this study showed a clear difference between inflectional systems using only suffixes and those that use different combinations of prefixes, suffixes and infixes. If inflectional markers are strictly suffixes, their length is predicted to be shorter than if the system uses a combination of affix positions. The effect of token frequency on marker length was also shown to be stronger for suffixes than for other combinations of marker positions. This means that suffixes are more susceptible to frequency effects on marker length than other affix positions.

A potential explanation for this difference across affix positions is phonetic reduction over time. We know from the literature that phonetic material at the end of words is reduced at higher rates than material at the beginning of words (Bybee *et al.* 1990, 19, Hall 1988). There is also evidence for word-initial (or domain-initial) syllables to be more prominent than other syllables (e.g. Beckman 1998; Smith 2005; Cho *et al.* 2007; Kim 2004; Keating *et al.* 2004). Especially word-initial consonants tend to be strengthened and lengthened (e.g. White *et al.* 2020; Cho and Keating 2009; Fougeron 2001; Cho and Keating 2001). This is relevant, since Bybee *et al.* (1990, 26) find that inflectional prefixes are crosslinguistically significantly more likely to have initial consonants than inflectional suffixes. Taken together, it is plausible that these properties contribute to suffixes being more likely candidates for phonetic reduction over time than affixes in other positions.

Support for the non-development scenario of zero markers

7.4

The other major finding from Section 6 is that the association between token frequency and marker length did not hold for zero markers. Their distributions in the Universal Dependencies treebanks showed that neither token frequency nor affix position were associated with

the occurrence of a zero marker. This is evidence against the traditional (implicit) assumption in typology that zero markers behave like short markers in terms of their distribution in language use (e.g. Bybee 2011; Croft 2003; Greenberg 1966; Haspelmath 2021). At the same time, the results from this study confirm previous studies, arguing that coding efficiency and frequency may not be suitable or a sufficient explanation for zero marking in inflectional morphology (Stolz and Levkovych 2019; Guzmán Naranjo and Becker 2021; Bickel *et al.* 2015; Cysouw 2003; Siewierska 2010; Seržant and Moroz 2022).

The difference between overt and zero markers in terms of their association with token frequencies also provides evidence for the non-development scenario leading to zero markers. The other, potential mechanism leading to zero marking is phonetic reduction. Phonetic reduction is commonly invoked as the mechanism responsible for the shortening of forms and the development of zero forms (Bybee 2003, 2007, 2015; Givón 2018; Haspelmath 2008a; Lehmann 2015). Bybee (2003, 2015) in particular has argued for phonetic reduction being a consequence of the repetition and automatization in production in the course of grammaticalization.

The main alternative to phonetic reduction is the differential non-development of a marker (cf. Bybee 1994; Cristofaro 2019, 2021; Haspelmath 2008a). For instance, we can imagine a scenario in which number is not marked on nouns at a given point in time. For independent reasons, plural marking could be developed. At the same time that the plural marker develops into an inflectional exponent, the absence of it becomes more systematically associated with the singular. Then, at some point, the singular is expressed by a zero form. In such a scenario, the zero marker results from the opposition to another new exponent of a different cell of the paradigm.

We can assume that phonetic reduction is at least in part responsible for the patterns found with overt markers, since we found a strong association between token frequency and marker length. Given that such an effect was not found for zero markers, the role of phonetic reduction as the main factor driving their development is questionable. As was mentioned above, the other main mechanism that can lead to the development of zero marking is the differential non-development of an inflection marker. For such a scenario, usage token frequencies may still play a role, but it would be much more indirect. In a non-

development scenario, the zero marker is only a consequence of the development of a different marker. The development process thus depends on a number of factors that are not directly related to the zero marker itself. The results from Section 6 cannot offer direct evidence in favor of the non-development scenario, but they are more compatible with this scenario than with the phonetic reduction scenario. There is certainly no single answer as to which mechanism leads to zero marking; it is likely that both and yet other mechanisms are involved, although likely to differing degrees. Diachronic corpus work is needed to shed more lights on the development of zero marking and its crosslinguistic tendencies.

CONCLUSION

8

This study offered a first token-based overview of zero marking in adjectival, nominal and verbal inflectional morphology across languages. Using the UniMorph dataset, it took into account the behavior of single lemmas to capture the variation across inflection classes and irregular forms. As for the probability of zero marking in inflection, the results showed that zero marking is generally not a preferred marking strategy, as it is predicted to only occur in 10-30% of inflected forms. No single cells or values showed a strong association with zero marking. Nevertheless, the values with the highest probability of zero marking (nominative, singular, indefinite, third person, present, imperative) confirmed observations from the typological literature. The findings further evidenced a high degree of idiosyncratic variation across languages and lemmas in the distribution of zero markers.

In addition, the study analyzed the token frequencies of zero markers together with those of overt markers in several corpora from the Universal Dependencies treebanks. For overt markers, the results showed that the token frequency has a stronger effect on the phonological length of suffixes compared to other affixes. This fits into a wider picture of phonetic differences between suffixes and other positions. For the probability of zero markers, however, no association with their frequency was found. This is new evidence for a fundamental difference between the distribution of overt and zero markers. Zero markers

do not simply follow the distributional patterns of short markers. This difference was argued to support a differential non-development scenario of zero marking rather than a phonetic reduction scenario.

REFERENCES

- Alexandra AIKHENVALD (2010), *Imperatives and Commands*, Oxford University Press.
- Sergiu AL-GEORGE (1967), The Semiosis of Zero According to Pāṇini, *East and West*, 17(1):115–124.
- Stephen R. ANDERSON (1992), *A-Morphous Morphology*, Cambridge University Press.
- Peter ARKADIEV (2016), Возможны Ли Однопадежные Системы? [Are Monocausal Systems Possible?], in Józefina PIĄTKOWSKA and Gennadij ZELDOWICZ, editors, *Znaki Czy Nie Znaki? - II. Zbiór Prac Lingwistycznych*, pp. 9–37, Wydawnictwa Uniwersytetu Warszawskiego.
- Matthew BAERMAN, Dunstan BROWN, and Greville G. CORBETT (2017), *Morphological Complexity*, Cambridge University Press.
- Matthew BAERMAN and Greville CORBETT (2012), Stem Alternations and Multiple Exponence, *Word Structure*, 5(1):52–68, doi:10.3366/word.2012.0019.
- Marianne BAKRÓ-NAGY (2022), Consonant Gradation, in Marianne BAKRÓ-NAGY, Johanna LAAKSO, and Elena SKRIBNIK, editors, *The Oxford Guide to Uralic Languages*, pp. 859–867, Oxford University Press.
- Sebastian BANK and Jochen TROMMER (2015), Learning and the Complexity of Ø-marking, in Matthew BAERMAN, Dunstan BROWN, and Greville CORBETT, editors, *Understanding and Measuring Morphological Complexity*, Oxford University Press.
- Jill BECKMAN (1998), Positional Faithfulness.
- Sacha BENIAMINE and Matías GUZMÁN NARANJO (2021), Multiple Alignments of Inflectional Paradigms, *Proceedings of the Society for Computation in Linguistics*, 4:216–227.
- Balthasar BICKEL, Alena WITZLACK-MAKAREVICH, Taras ZAKHARKO, and Giorgio IEMMOLO (2015), Exploring Diachronic Universals of Agreement: Alignment Patterns and Zero Marking across Person Categories, in Jürg FLEISCHER, Elisabeth RIEKEN, and Paul WIDMER, editors, *Agreement from a Diachronic Perspective*, pp. 29–52, De Gruyter.
- James BLEVINS (2003), Stems and Paradigms, *Language*, 79(4):737–767.

Zero marking in inflection

- James BLEVINS (2005), Word-Based Declensions in Estonian, in Geert BOOIJ and Jaap VAN MARLE, editors, *Yearbook of Morphology 2005*, pp. 1–25, Springer, doi:10.1007/1-4020-4066-0_1.
- James BLEVINS (2006), Word-Based Morphology, *Journal of Linguistics*, 42(3):531–573.
- James BLEVINS (2016), *Word and Paradigm Morphology*, Oxford University Press.
- Bernard BLOCH (1947), English Verb Inflection, *Language*, 23(4):399–418, doi:10.2307/410300.
- Leonard BLOOMFIELD (1933), *Language*, Holt.
- Olivier BONAMI (2012), Stems in Inflection and Lexeme Formation, *Word Structure*, 5(1).
- Olivier BONAMI and Sacha BENIAMINE (2021), Leaving the Stem by Itself, in Sedigheh MORADI, Marcia HAAG, Janie REES-MILLER, and Andrija PETROVIC, editors, *All Things Morphology: Its Independence and Its Interfaces*, pp. 81–98, Benjamins, doi:10.1075/cilt.353.05bon.
- Sami BOUDELAA and William D MARSLEN-WILSON (2001), Morphological Units in the Arabic Mental Lexicon, *Cognition*, 81(1):65–92, doi:10.1016/S0010-0277(01)00119-6.
- Gilles BOYÉ and Gauvin SCHALCHI (2016), The Status of Paradigms, in Andrew HIPPISEY and Gregory STUMP, editors, *The Cambridge Handbook of Morphology*, pp. 206–234, Cambridge University Press.
- Dunstan BROWN (1998), Stem Indexing and Morphological Selection in the Russian Verb: A Network Morphology Account, in Ray FABRI, Albert ORTMANN, and Teresa PARODI, editors, *Models of Inflection*, pp. 196–224, Niemeyer.
- Joan BYBEE (1994), The Grammaticization of Zero: Asymmetries in Tense and Aspect Systems, in William PAGLIUCA, editor, *Perspectives on Grammaticalization*, pp. 235–254, Benjamins.
- Joan BYBEE (2003), Mechanisms of Change in Grammaticization: The Role of Frequency, in Brian JOSEPH and Richard JANDA, editors, *Handbook of Historical Linguistics*, pp. 602–623, Blackwell.
- Joan BYBEE (2007), *Frequency of Use and the Organization of Language*, Oxford University Press.
- Joan BYBEE (2011), Markedness, in Jae Jung SONG, editor, *The Oxford Handbook of Typology*, pp. 1–11, Oxford University Press.
- Joan BYBEE (2015), *Language Change*, Cambridge University Press.
- Joan BYBEE and Östen DAHL (1989), The Creation of Tense and Aspect Systems in the Languages of the World, *Studies in Language*, 13(1):51–103.

- Joan BYBEE, William PAGLIUCA, and Revere PERKINS (1990), On the Asymmetries in the Affixation of Grammatical Material, in William CROFT, Suzanne KEMMER, and Keith DENNING, editors, *Studies in Typology and Diachrony. Papers Presented to Joseph H. Greenberg on His 75th Birthday*, pp. 1–42, Benjamins.
- Paul-Christian BÜRKNER (2017), Brms: An R Package for Bayesian Multilevel Models Using Stan, *Journal of Statistical Software*, 80(1):1–28, doi:10.18637/jss.v080.i01.
- Bob CARPENTER, Andrew GELMAN, Matthew HOFFMAN, Daniel LEE, Ben GOODRICH, Michael BETANCOURT, Marcus BRUBAKER, Jiqiang GUO, Peter LI, and Allen RIDDELL (2017), Stan: A Probabilistic Programming Language, *Journal of Statistical Software*, 76(1):1–32, doi:10.18637/jss.v076.i01, <https://www.jstatsoft.org/v076/i01>.
- Taehong CHO and Patricia KEATING (2001), Articulatory and Acoustic Studies on Domain-Initial Strengthening in Korean, *Journal of Phonetics*, 29(2):155–190, doi:10.1006/jpho.2001.0131.
- Taehong CHO and Patricia KEATING (2009), Effects of Initial Position versus Prominence in English, *Journal of Phonetics*, 37(4):466–485, doi:10.1016/j.wocn.2009.08.001.
- Taehong CHO, James MCQUEEN, and Ethan COX (2007), Prosodically Driven Phonetic Detail in Speech Processing: The Case of Domain-Initial Strengthening in English, *Journal of Phonetics*, 35(2):210–243, doi:10.1016/j.wocn.2006.03.003.
- Matt COLER (2015), Aymara Inflection, in Matthew BAERMAN, editor, *The Oxford Handbook of Inflection*, pp. 1–30, Oxford University Press.
- Matt COLER (2018), Subtractive Morphology & Disfixation in Aymara Case, pp. 1–9, doi:10.13140/RG.2.2.26153.03682.
- Ellen CONTINI-MORAVA (2006), The Difference between Zero and Nothing: Swahili Noun Class Prefixes 5 and 9/10, in Joseph DAVIS, Radmila GORUP, and Nancy STERN, editors, *Advances in Functional Linguistics*, pp. 211–222, Benjamins.
- Greville CORBETT (2007), Canonical Typology, Suppletion, and Possible Words, *Language*, 83(1):8–42, doi:10.1353/lan.2007.0006.
- Sonia CRISTOFARO (2019), Taking Diachronic Evidence Seriously: Result-oriented vs. Source-Oriented Explanations of Typological Universals, in Karsten SCHMIDTKE-BODE, Natalia LEVSHINA, Susanne Maria MICHAELIS, and Ilya SERŽANT, editors, *Explanation in Typology: Diachronic Sources, Functional Motivations and the Nature of the Evidence*, pp. 25–46, Language Science Press.
- Sonia CRISTOFARO (2021), Typological Explanations in Synchrony and Diachrony: On the Origins of Third Person Zeroes in Bound Person Paradigms, *Folia Linguistica*, 55(s42-s1):25–48, doi:10.1515/flin-2021-2013.

- William CROFT (2003), *Typology and Universals*, Cambridge University Press, 2nd edition.
- Michael CYSOUW (2003), *The Paradigmatic Structure of Person Marking*, Oxford University Press.
- Eystein DAHL and Antonio FÁBREGAS (2018), Zero Morphemes, in Rochelle LIEBER, editor, *Oxford Research Encyclopedia of Linguistics*, pp. 1–30, Oxford University Press, doi:10.1093/acrefore/9780199384655.013.592, <https://oxfordre.com/linguistics/view/10.1093/acrefore/9780199384655.001.0001/acrefore-9780199384655-e-592>.
- Catharine DIEHL (2008), The Empty Space in Structure: Theories of the Zero from Gauthiot to Deleuze, *Diacritics*, 38(3):93–119, ISSN 0300-7162, <https://www.jstor.org/stable/20616535>.
- Holger DIESSEL (2019), *The Grammar Network: How Linguistic Structure Is Shaped by Language Use*, Cambridge University Press.
- Laura J. DOWNING and Barbara STIEBELS (2012), Iconicity, in Jochen TROMMER, editor, *The Morphology and Phonology of Exponence*, Oxford Studies in Theoretical Linguistics, Oxford University Press.
- Cécile FOUGERON (2001), Articulatory Properties of Initial Segments in Several Prosodic Constituents in French, *Journal of Phonetics*, 29(2):109–135, doi:10.1006/jpho.2000.0114.
- Talmy GIVÓN (1991), Isomorphism in the Grammatical Code: Cognitive and Biological Considerations, *Studies in Language*, 15(1):85–114, doi:10.1075/sl.15.1.04giv.
- Talmy GIVÓN (2018), *On Understanding Grammar*, Benjamins.
- Joseph GREENBERG (1963), Some Universals of Grammar with Particular Reference to the Order of Meaningful Elements, in Joseph GREENBERG, editor, *Universals of Language*, pp. 73–113, MIT Press.
- Joseph GREENBERG (1966), *Language Universals: With Special Reference to Feature Hierarchies*, Mouton.
- Matías GUZMÁN NARANJO and Laura BECKER (2021), Coding Efficiency in Nominal Inflection: Expectedness and Type Frequency Effects, *Linguistics Vanguard*, 7(s3):20190075, doi:10.1515/lingvan-2019-0075.
- Matías GUZMÁN NARANJO and Laura BECKER (2022), Statistical Bias Control in Typology, *Linguistic Typology*, 26(3):605–670, doi:10.1515/lingty-2021-0002.
- William HAAS (1957), Zero in Linguistics Description, in John Rupert FIRTH, editor, *Studies in Linguistic Analysis*, pp. 33–53, Blackwell.
- Christopher HALL (1988), Integrating Diachronic and Processing Principles in Explaining the Suffixing Preference, in John HAWKINS, editor, *Explaining Language Universals*, pp. 321–349, Basil Blackwell.

Harald HAMMARSTRÖM, Robert FORKEL, Martin HASPELMATH, and Sebastian BANK (2021), *Glottolog 4.4*, Max Planck Institute for the Science of Human History, <http://glottolog.org>.

Martin HASPELMATH (2008a), Creating Economical Morphosyntactic Patterns in Language Change, in Jeff GOOD, editor, *Linguistic Universals and Language Change*, pp. 185–214, Oxford University Press.

Martin HASPELMATH (2008b), Frequency vs. Iconicity in Explaining Grammatical Asymmetries, *Cognitive Linguistics*, 19(1):1–33, doi:10.1515/COG.2008.001.

Martin HASPELMATH (2008c), A Frequentist Explanation of Some Universals of Reflexive Marking, *Linguistic Discovery*, 6(1):40–63, doi:10.1349/PS1.1537-0852.A.331.

Martin HASPELMATH (2018), How Comparative Concepts and Descriptive Linguistic Categories Are Different, in *Aspects of Linguistic Variation*, pp. 83–114, De Gruyter, doi:10.1515/9783110607963-004.

Martin HASPELMATH (2021), Explaining Grammatical Coding Asymmetries: Form–Frequency Correspondences and Predictability, *Journal of Linguistics*, pp. 1–29, doi:10.1017/S0022226720000535.

Martin HASPELMATH, Andreea CALUDE, Michael SPAGNOL, Heiko NARROG, and Elif BAMYACI (2014), Coding Causal–Noncausal Verb Alternations: A Form–Frequency Correspondence Explanation, *Journal of Linguistics*, 50(3):587–625.

Martin HASPELMATH and Andres KARJUS (2017), Explaining Asymmetries in Number Marking: Singulatives, Pluratives, and Usage Frequency, *Linguistics*, 55(6):1213–1235, doi:10.1515/ling-2017-0026.

George HEWITT (1995), *Georgian: A Structural Reference Grammar*, Benjamins.

Jane HILL and Ofelia ZEPEDA (1998), Tohono O'odham (Papago) Plurals, *Anthropological Linguistics*, 40(1):1–42.

Roman JAKOBSON (1983), Zero Sign, in Linda WAUGH and Morris HALLE, editors, *Russian and Slavic Grammar: Studies 1931-1981*, pp. 1–14, De Gruyter.

Patricia KEATING, Taehong CHO, Fougeron CECILE, and Chai-Shune HSU (2004), Domain-Initial Strengthening in Four Languages, in John LOCAL, Richard ODGEN, and Rosalind TEMPLE, editors, *Phonetic Interpretation. Papers in Laboratory Phonology VI*, pp. 145–163, Cambridge University Press.

Sahyang KIM (2004), The Role of Prosodic Phrasing in Korean Word Segmentation.

Harold KOCH (1995), The Creation of Morphological Zeros, in Geert BOOIJ and Jaap VAN MARLE, editors, *Yearbook of Morphology 1994*, pp. 31–731, Springer.

- Christian LEHMANN (1974), Isomorphismus Im Sprachlichen Zeichen [Isomorphism in the linguistic sign], in *Linguistic Workshop II: Arbeiten Des Kölner Universalienprojekts 1973/4*, pp. 98–123, Fink.
- Christian LEHMANN (2015), *Thoughts on Grammaticalization*, Language Science Press.
- Natalia LEVSHINA (2022), *Communicative Efficiency: Language Structure and Use*, Cambridge University Press.
- Martin MAIDEN (1992), Irregularity as a Determinant of Morphological Change, *Journal of Linguistics*, 28(2):285–312.
- Peter Hugoe MATTHEWS (1972), *Inflectional Morphology: A Theoretical Study Based on Aspects of Latin Verb Conjugation*, Cambridge University Press.
- Arya D. MCCARTHY, Christo KIROV, Matteo GRELLA, Amrit NIDHI, Patrick XIA, Kyle GORMAN, Ekaterina VYLOMOVA, Sabrina J. MIELKE, Garrett NICOLAI, Miikka SILFVERBERG, Timofey ARKHANGELSKIY, Nataly KRIZHANOVSKY, Andrew KRIZHANOVSKY, Elena KLYACHKO, Alexey SOROKIN, John MANSFIELD, Valts ERNŠTREITS, Yuval PINTER, Cassandra L. JACOBS, Ryan COTTERELL, Mans HULDEN, and David YAROWSKY (2020), UniMorph 3.0: Universal Morphology, in *Proceedings of the 12th Language Resources and Evaluation Conference*, pp. 3922–3931, European Language Resources Association, <https://aclanthology.org/2020.lrec-1.483>.
- William MCGREGOR (2003), The Nothing That Is, the Zero That Isn't, *Studia Linguistica*, 57(2):75–119, doi:10.1111/1467-9582.00100.
- Georg Friedrich MEIER (1961), *Das Zéro-Problem in Der Linguistik. Kritische Untersuchungen Zur Strukturalistischen Analyse Der Relevanz Sprachlicher Form*, Akademie Verlag.
- Igor MEL'ČUK (1994), Suppletion: Toward a Logical Analysis of the Concept, *Studies in Language*, 18(2):339–410, doi:10.1075/sl.18.2.03mel.
- Igor MEL'ČUK (2002), Towards a Formal Concept Zero Linguistic Sign: Applications in Typology, in Sabrina BENDJABALLAH, Wolfgang DRESSLER, Oskar PFEIFFER, and Maria VOEIKOVA, editors, *Morphology 2000: Selected Papers from the 9th Morphology Meeting, Vienna, 24–28 February 2000*, pp. 241–258, Benjamins.
- Marianne MITHUN (1986), When Zero Isn't There, *Annual Meeting of the Berkeley Linguistics Society*, 12(0):195–211, doi:10.3765/bls.v12i0.1882.
- Fabio MONTERMINI and Olivier BONAMI (2013), Stem Spaces and Predictability in Verbal Inflection, *Lingue e linguaggio*, 2:171–190, doi:10.1418/75040.
- David R MORTENSEN, Siddharth DALMIA, and Patrick LITTELL (2018), Epitran: Precision G2P for Many Languages, in *Proceedings of the Eleventh International Conference on Language Resources and Evaluation (LREC 2018)*.

- Irina NIKOLAEVA (2007), *Finiteness: Theoretical and Empirical Foundations*, Oxford University Press.
- Mary PASTER (2016), Alternations: Stems and Allomorphy, in Andrew HIPPISEY and Gregory STUMP, editors, *The Cambridge Handbook of Morphology*, pp. 93–116, Cambridge University Press.
- Vito PIRELLI and Marco BATTISTA (2000), The Paradigmatic Dimension of Stem Allomorphy in Italian Verb Inflection, *Rivista di Linguistica*, 12(2):307–380.
- Geoffrey PULLUM and Arnold ZWICKY (1991), A Misconceived Approach to Morphology, *Proceedings of the West Coast Conference on Formal Linguistics*, 10.
- R CORE TEAM (2021), *R: A Language and Environment for Statistical Computing*, R Foundation for Statistical Computing, <https://www.R-project.org/>.
- Robert RATCLIFFE (1998), *The “Broken” Plural Problem in Arabic and Comparative Semitic*, Benjamins.
- R. H. ROBINS (1997), *A Short History of Linguistics*, Routledge, 4th edition.
- Jerrold SADOCK and Arnold ZWICKY (1985), Speech Act Distinctions in Grammar, in Timothy SHOPEN, editor, *Language Typology and Syntactic Description. Volume 1*, pp. 155–196, Cambridge University Press.
- Gerald SANDERS (1988), Zero Derivation and the Overt Analogue Criterion, in Michael HAMMOND and Michael NOONAN, editors, *Theoretical Morphology*, pp. 155–175, Academic Press.
- Ferdinand SAUSSURE (1916), *Cours de Linguistique Générale*, Payot.
- Gene SCHRAMM (1962), An Outline of Classical Arabic Verb Structure, *Language*, 38(4):360–375, doi:10.2307/410672.
- Ilja SERŽANT and George MOROZ (2022), Universal Attractors in Language Evolution Provide Evidence for the Kinds of Efficiency Pressures Involved, *Humanities and Social Sciences Communications*, 9(1):1–9, doi:10.1057/s41599-022-01072-0.
- Anna SIEWIERSKA (2010), Person Asymmetries in Zero Expression and Grammatical Functions, in Franck FLORICIC, editor, *Essais de Linguistique Generale et de Typologie Linguistique Offerts Au Professeur Denis Creissels à l’occasion de Ses 65 Ans*, pp. 425–438, Presses de l’École Normale Supérieure.
- Jennifer SMITH (2005), *Phonological Augmentation in Prominent Positions*, Taylor & Francis, doi:10.4324/9780203506394.
- Jae Jung SONG (2018), *Linguistic Typology*, Oxford University Press.
- Andrew SPENCER (2012), Identifying Stems, *Word Structure*, 5(1):88–108, doi:10.3366/word.2012.0021.
- Matthew STAVE, Ludger PASCHEN, François PELLEGRINO, and Frank SEIFART (2021), Optimization of Morpheme Length: A Cross-Linguistic Assessment of Zipf’s and Menzerath’s Laws, *Linguistics Vanguard*, 7(s3), doi:10.1515/lingvan-2019-0076.

- Thomas STOLZ and Nataliya LEVKOVYCH (2019), Absence of Material Exponence, *Language Typology and Universals*, 72(3):373–400, doi:10.1515/stuf-2019-0015.
- Gregory STUMP (2001), *Inflectional Morphology: A Theory of Paradigm Structure*, Cambridge University Press.
- Gregory STUMP and Rafael FINKEL (2013), *Morphological Typology: From Word to Paradigm*, volume 138 of *Cambridge Studies in Linguistics*, Cambridge University Press.
- John SYLAK-GLASSMAN (2016), The Composition and Use of the Universal Morphological Feature Schema (UniMorph Schema), <https://unimorph.github.io/doc/unimorph-schema.pdf>.
- Anna THORNTON (2012), Reduction and Maintenance of Overabundance. A Case Study on Italian Verb Paradigms, *Word Structure*, 5(2):183–207, doi:10.3366/word.2012.0026.
- Jochen TROMMER (2012), Ø-Exponence, in Jochen TROMMER, editor, *The Morphology and Phonology of Exponence*, pp. 326–354, Oxford University Press.
- Aki VEHTARI, Andrew GELMAN, and Jonah GABRY (2017), Practical Bayesian Model Evaluation Using Leave-One-out Cross-Validation and WAIC, *Statistics and Computing*, 27(5):1413–1432, doi:10.1007/s11222-016-9696-4.
- Laurence WHITE, Silvia BENAVIDES-VARELA, and Katalin MÁDY (2020), Are Initial-Consonant Lengthening and Final-Vowel Lengthening Both Universal Word Segmentation Cues?, *Journal of Phonetics*, 81:100982, doi:10.1016/J.WOCN.2020.100982.
- Jingting YE (2020), Independent and Dependent Possessive Person Forms: Three Universals, *Studies in Language*, 44(2):363–406, doi:10.1075/sl.19020.ye.
- Moira YIP (1988), Template Morphology and the Direction of Association, *Natural Language & Linguistic Theory*, 6(4):551–577.
- Daniel ZEMAN, Joakim NIVRE, Mitchell ABRAMS, Elia ACKERMANN, Noëmi AEPLI, Hamid AGHAEI, Željko AGIĆ, Amir AHMADI, Lars AHRENBORG, Chika Kennedy AJEDE, Gabrielé ALEKSANDRAVIČIŪTĖ, Ika ALFINA, Lene ANTONSEN, Katya APLONOVA, Angelina AQUINO, Carolina ARAGON, Maria Jesus ARANZABE, Bilge Nas ARICAN, Órunn ARNARDÓTTIR, Gashaw ARUTIE, Jessica Naraiswari ARWIDARASTI, Masayuki ASAHARA, Deniz Baran ASLAN, Luma ATEYAH, Furkan ATMACA, Mohammed ATTIA, Aitziber ATUTXA, Liesbeth AUGUSTINUS, Elena BADMAEVA, Keerthana BALASUBRAMANI, Miguel BALLESTEROS, Esha BANERJEE, Sebastian BANK, Verginica BARBU MITITELU, Starkaður BARKARSON, Rodolfo BASILE, Victoria BASMOV, Colin BATCHELOR, John BAUER, Seyyit Talha BEDIR, Kepa BENGOTXEA, Gözde BERK, Yevgeni BERZAK, Irshad Ahmad BHAT, Riyaz Ahmad BHAT, Erica BIAGETTI, Eckhard BICK, Agnė BIELINSKIENĖ, Kristín BJARNADÓTTIR, Rogier

BLOKLAND, Victoria BOBICEV, Loïc BOIZOU, Emanuel BORGES VÖLKER, Carl BÖRSTELL, Cristina BOSCO, Gosse BOUMA, Sam BOWMAN, Adriane BOYD, Anouck BRAGGAAR, Kristina BROKAITĖ, Aljoscha BURCHARDT, Marie CANDITO, Bernard CARON, Gauthier CARON, Lauren CASSIDY, Tatiana CAVALCANTI, Gülşen CEBIROĞLU ERYİĞİT, Flavio Massimiliano CECCHINI, Giuseppe G. A. CELANO, Slavomír ČÉPLÖ, Neslihan CESUR, Savas CETIN, Özlem ÇETİNOĞLU, Fabricio CHALUB, Shweta CHAUHAN, Ethan CHI, Taishi CHIKA, Yongseok CHO, Jinho CHOI, Jayeol CHUN, Juyeon CHUNG, Alessandra T. CIGNARELLA, Silvie CINKOVÁ, Aurélie COLLOMB, Çağrı ÇÖLTEKİN, Miriam CONNOR, Marine COURTIN, Mihaela CRISTESCU, Philemon DANIEL, Elizabeth DAVIDSON, Marie-Catherine DE MARNEFFE, Valeria DE PAIVA, Mehmet Oguz DERIN, Elvis DE SOUZA, Arantza DIAZ DE ILARRAZA, Carly DICKERSON, Arawinda DINAKARAMANI, Elisa DI NUOVO, Bamba DIONE, Peter DIRIX, Kaja DOBROVOLJC, Timothy DOZAT, Kira DROGANOVA, Puneet DWIVEDI, Hanne ECKHOFF, Sandra EICHE, Marhaba ELI, Ali ELKAHKY, Binyam EPHREM, Olga ERINA, Tomaz ERJAVEC, Aline ETIENNE, Wograine EVELYN, Sidney FACUNDES, Richárd FARKAS, Jannatul FERDAOUSI, Marília FERNANDA, Hector FERNANDEZ ALCALDE, Jennifer FOSTER, Cláudia FREITAS, Kazunori FUJITA, Katarína GAJDOŠOVÁ, Daniel GALBRAITH, Marcos GARCIA, Moa GÄRDENFORS, Sebastian GARZA, Fabrício Ferraz GERARDI, Kim GERDES, Filip GINTER, Gustavo GODOY, Iakes GOENAGA, Koldo GOJENOLA, Memduh GÖKIRMAK, Yoav GOLDBERG, Xavier GÓMEZ GUINOVART, Berta GONZÁLEZ SAAVEDRA, Bernadeta GRICIŪTĖ, Matias GRIONI, Loïc GROBOL, Normunds GRŪZTIS, Bruno GUILLAUME, Céline GUILLOT-BARBANCE, Tunga GÜNGÖR, Nizar HABASH, Hinrik HAFSTEINSSON, Jan HAJIČ, Jan HAJIČ JR., Mika HÄMÄLÄINEN, Linh HÀ MỸ, Na-Rae HAN, Muhammad Yudistira HANIFMUTI, Sam HARDWICK, Kim HARRIS, Dag HAUG, Johannes HEINECKE, Oliver HELLWIG, Felix HENNIG, Barbora HLADKÁ, Jaroslava HLAVÁČOVÁ, Florinel HOCIUNG, Petter HOHLE, Eva HUBER, Jena HWANG, Takumi IKEDA, Anton Karl INGASON, Radu ION, Elena IRIMIA, Olájidé ISHOLA, Kaoru ITO, Siratun JANNAT, Tomáš JELÍNEK, Apoorva JHA, Anders JOHANNSEN, Hildur JÓNSDÓTTIR, Fredrik JØRGENSEN, Markus JUUTINEN, Sarveswaran K, Hüner KAŞIKARA, Andre KAASEN, Nadezhda KABAEVA, Sylvain KAHANE, Hiroshi KANAYAMA, Jenna KANERVA, Neslihan KARA, Boris KATZ, Tolga KAYADELEN, Jessica KENNEY, Václava KETTNEROVÁ, Jesse KIRCHNER, Elena KLEMENTIEVA, Elena KLYACHKO, Arne KÖHN, Abdullatif KÖKSAL, Kamil KOPACEWICZ, Timo KORKIAKANGAS, Mehmet KÖSE, Natalia KOTSYBA, Jolanta KOVALEVSKAITĖ, Simon KREK, Parameswari KRISHNAMURTHY, Sandra KÜBLER, Oğuzhan KUYRUKÇU, Asli KUZGUN, Sookyoung KWAK, Veronika LAIPPALA, Lucia LAM, Lorenzo LAMBERTINO, Tatiana LANDO, Septina Dian LARASATI, Alexei LAVRENTIEV, John LEE, Phương LÊ HỒNG, Alessandro LENCI, Saran LERTPRADIT, Herman LEUNG, Maria LEVINA, Cheuk Ying LI, Josie LI, Keying LI, Yuan LI, KyungTae LIM, Bruna

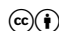
LIMA PADOVANI, Krister LINDÉN, Nikola LJUBEŠIĆ, Olga LOGINOVA, Stefano LUSITO, Andry LUTHFI, Mikko LUUKKO, Olga LYASHEVSKAYA, Teresa LYNN, Vivien MACKETANZ, Menel MAHAMDI, Jean MAILLARD, Aibek MAKAZHANOV, Michael MANDL, Christopher MANNING, Ruli MANURUNG, Büşra MARŞAN, Cătălina MĂRĂNDUC, David MAREČEK, Katrin MARHEINECKE, Héctor MARTÍNEZ ALONSO, Lorena MARTÍN-RODRÍGUEZ, André MARTINS, Jan MAŠEK, Hiroshi MATSUDA, Yuji MATSUMOTO, Alessandro MAZZEI, Ryan McDONALD, Sarah MCGUINNESS, Gustavo MENDONÇA, Tatiana MERZHEVICH, Niko MIEKKA, Karina MISCHENKOVA, Margarita MISIRPASHAYEVA, Anna MISSILÄ, Cătălin MITITELU, Maria MITROFAN, Yusuke MIYAO, AmirHossein MOJIRI FOROUSHANI, Judit MOLNÁR, Amirsaeid MOLOODI, Simonetta MONTEMAGNI, Amir MORE, Laura MORENO ROMERO, Giovanni MORETTI, Keiko Sophie MORI, Shinsuke MORI, Tomohiko MORIOKA, Shigeki MORO, Bjartur MORTENSEN, Bohdan MOSKALEVSKYI, Kadri MUISCHNEK, Robert MUNRO, Yugo MURAWAKI, Kaili MÜÜRISep, Pinkey NAINWANI, Mariam NAKHLÉ, Juan Ignacio NAVARRO HORÑIACEK, Anna NEDOLUZHKO, Gunta NEŠPORE-BĚRZKALNE, Manuela NEVACI, Lư ỡng NGUYỄN THỊ, Huy ẽn NGUYỄN THỊ MINH, Yoshihiro NIKAIDO, Vitaly NIKOLAEV, Rattima NITISAROJ, Alireza NOURIAN, Hanna NURMI, Stina OJALA, Atul Kr. OJHA, Adédayọ OLÚÒKUN, Mai OMURA, Emeka ONWUEGBUZIA, Petya OSENOVA, Robert ÖSTLING, Lilja ØVRELID, Şaziye Betül ÖZATEŞ, Merve ÖZÇELİK, Arzucan ÖZGÜR, Balkız ÖZTÜRK BAŞARAN, Hyunji Hayley PARK, Niko PARTANEN, Elena PASCUAL, Marco PASSAROTTI, Agnieszka PATEJUK, Guilherme PAULINO-PASSOS, Angelika PELJAK-ŁAPIŃSKA, Siyao PENG, Genel-Augusto PEREZ, Natalia PERKOVA, Guy PERRIER, Slav PETROV, Daria PETROVA, Jason PHELAN, Jussi PIITULAINEN, Tommi A PIRINEN, Emily PITLER, Barbara PLANK, Thierry POIBEAU, Larisa PONOMAREVA, Martin POPEL, Lauma PRETKALNIŃA, Sophie PRÉVOST, Prokopis PROKOPIDIS, Adam PRZEPIÓRKOWSKI, Tiina PUOLAKAINEN, Sampo PYYSALO, Peng QI, Andriela RÄÄBIS, Alexandre RADEMAKER, Mizanur RAHOMAN, Taraka RAMA, Loganathan RAMASAMY, Carlos RAMISCH, Fam RASHEL, Mohammad Sadegh RASOOLI, Vinit RAVISHANKAR, Livy REAL, Petru REBEJA, Siva REDDY, Mathilde REGNAULT, Georg REHM, Ivan RIABOV, Michael RIESSLER, Erika RIMKUTĚ, Larissa RINALDI, Laura RITUMA, Putri RIZQIYAH, Luisa ROCHA, Eiríkur RÖGNVALDSSON, Mykhailo ROMANENKO, Rudolf ROSA, Valentin ROŞCA, Davide ROVATI, Olga RUDINA, Jack RUETER, Kristján RÚNARSSON, Shoval SADDE, Pegah SAFARI, Benoît SAGOT, Aleksí SAHALA, Shadi SALEH, Alessio SALOMONI, Tanja SAMARDŽIĆ, Stephanie SAMSON, Manuela SANGUINETTI, Ezgi SANIYAR, Dage SÄRG, Baiba SAULTE, Yanin SAWANAKUNANON, Shefali SAXENA, Kevin SCANNELL, Salvatore SCARLATA, Nathan SCHNEIDER, Sebastian SCHUSTER, Lane SCHWARTZ, Djamé SEDDAH, Wolfgang SEEKER, Mojgan SERAJI, Syeda SHAHZADI, Mo SHEN, Atsuko SHIMADA, Hiroyuki SHIRASU, Yana SHISHKINA, Muh SHOHIBUSSIRRI,

Dmitry SICHINA, Janine SIEWERT, Einar Freyr SIGURÐSSON, Aline SILVEIRA, Natalia SILVEIRA, Maria SIMI, Radu SIMIONESCU, Katalin SIMKÓ, Mária ŠIMKOVÁ, Kiril SIMOV, Maria SKACHEDUBOVA, Aaron SMITH, Isabela SOARES-BASTOS, Shafi SOUROV, Carolyn SPADINE, Rachele SPRUGNOLI, Steinór STEINGRÍMSSON, Antonio STELLA, Milan STRAKA, Emmett STRICKLAND, Jana STRNADOVÁ, Alane SUHR, Yogi Lesmana SULESTIO, Umur SULUBACAK, Shingo SUZUKI, Zsolt SZÁNTÓ, Chihiro TAGUCHI, Dima TAJI, Yuta TAKAHASHI, Fabio TAMBURINI, Mary Ann C. TAN, Takaaki TANAKA, Dipta TANAYA, Samson TELLA, Isabelle TELLIER, Marinella TESTORI, Guillaume THOMAS, Liisi TORGA, Marsida TOSKA, Trond TROSTERUD, Anna TRUKHINA, Reut TSARFATY, Utku TÜRK, Francis TYERS, Sumire UEMATSU, Roman UNTILOV, Zdeňka UREŠOVÁ, Larraitz URÍA, Hans USZKOREIT, Andrius UTKA, Sowmya VAJJALA, Rob VAN DER GOOT, Martine VANHOVE, Daniel VAN NIEKERK, Gertjan VAN NOORD, Viktor VARGA, Eric VILLEMONT DE LA CLERGERIE, Veronika VINCZE, Natalia VLASOVA, Aya WAKASA, Joel C. WALLENBERG, Lars WALLIN, Abigail WALSH, Jing Xian WANG, Jonathan North WASHINGTON, Maximilian WENDT, Paul WIDMER, Sri Hartati WIJONO, Seyi WILLIAMS, Mats WIRÉN, Christian WITTERN, Tsegay WOLDEMARIAM, Tak-sum WONG, Alina WRÓBLEWSKA, Mary YAKO, Kayo YAMASHITA, Naoki YAMAZAKI, Chunxiao YAN, Koichi YASUOKA, Marat M. YAVRUMYAN, Arife Betül YENICE, Olcay Taner YILDIZ, Zhuoran YU, Arlisa YULIAWATI, Zdeněk ŽABOKRTSKÝ, Shorouq ZAHRA, Amir ZELDES, He ZHOU, Hanzhi ZHU, Anna ZHURAVLEVA, and Rayan ZIANE (2023), *Universal Dependencies 2.13*, <http://hdl.handle.net/11234/1-5287>.

George Kingsley ZIPF (1935), *The Psychobiology of Language: An Introduction to Dynamic Philology*, MIT Press.

Arnold ZWICKY (1985), How to Describe Inflection, in Mary NIEPOKUIJ, Mary VAN CLAY, Vassiliki NIKIFORIDOU, and Deborah FEDER, editors, *Proceedings of the Eleventh Annual Meeting of the Berkeley Linguistics Society*, pp. 372–386, Berkeley Linguistics Society.

This work is licensed under the *Creative Commons Attribution 4.0 Public License*.

 <http://creativecommons.org/licenses/by/4.0/>