# How to encode arbitrarily complex morphology in word embeddings, no corpus needed

#### Lane Schwartz

Department of Computer Science University of Alaska Fairbanks lane.schwartz@alaska.edu

## **Coleman Haley**

Institute for Language, Cognition and Computation University of Edinburgh Coleman. Haley@ed.ac.uk

## **Francis Tyers**

Department of Linguistics Indiana University ftyers@iu.edu

#### **Abstract**

In this paper, we present a straightforward technique for constructing interpretable word embeddings from morphologically analyzed examples (such as interlinear glosses) for all of the world's languages. Currently, fewer than 300–400 languages out of approximately 7000 have have more than a trivial amount of digitized texts; of those, between 100-200 languages (most in the Indo-European language family) have enough text data for BERT embeddings of reasonable quality to be trained. The word embeddings in this paper are explicitly designed to be both linguistically interpretable and fully capable of handling the broad variety found in the world's diverse set of 7000 languages, regardless of corpus size or morphological characteristics. We demonstrate the applicability of our representation through examples drawn from a typologically diverse set of languages whose morphology includes prefixes, suffixes, infixes, circumfixes, templatic morphemes, derivational morphemes, inflectional morphemes, and reduplication.

## 1 Better representations are needed

The past several years have seen the development of neural techniques capable of creating extremely high quality word embeddings, most notably BERT (Devlin et al., 2019) and its many variants. In total, however, fewer than 300–400 languages have have more than a trivial amount of digitized text data, thus rendering data-driven NLP approaches including BERT futile for more than 6000 remaining languages (representing over 1.2 billion people; Vannini and Crosnier, 2012; Joshi et al., 2020), even with aggressive multilingual models, transfer learning, bilingual anchoring, and typologically-aware modelling (Ponti et al., 2019; Michel et al., 2020; Eder et al., 2021; Hedderich et al., 2021).

Somewhere between 100–200 languages (most in the Indo-European language family) have enough digitized text data (Joshi et al., 2020; Conneau et al., 2020) for BERT embeddings of reasonable quality to be trained using a combination of techniques including unsupervised sub-word segmentation methods, multilingual bootstrapping, and transfer learning. Quality of word embeddings is substantially lower when corpus sizes are insufficiently large; Alabi et al. (2020), for example, constructed word embeddings using approximately 10 million tokens for Yorùbá¹ and Twi,² and found that the resulting embeddings are substantially poorer in quality those for high-resource languages.

## 1.1 Complex morphology is the norm

The issue of insufficient training data is exacerbated even more when productive derivational and inflectional morphology plays a significant role in word formation in a language. The average number of morphemes per word is medium or high for the vast majority of the world's approximately 7000 languages (see *World Atlas of Language Structures*, including Bickel and Nichols, 2013; Dryer, 2013). Despite this fact, since at least Oettinger (1954), the primary meaning-bearing unit used to represent language in natural language models has been the word.

While many modern NLP models can and sometimes do represent higher-level linguistics units (representing phrases, clauses, or sentences) and lower-level linguistic units (such as morphemes, sub-word chunks, or characters), and notwithstanding the widespread use of unsupervised subword

<sup>&</sup>lt;sup>1</sup>ISO 639-3: *yor*, an analytic language in the Yoruboid branch of the Niger-Congo language family

<sup>&</sup>lt;sup>2</sup>ISO 639-3: *twi*, an analytic language in the Tano branch of the Niger-Congo language family

segmentation methods (BPE, SentencePiece, etc), there remains a very common yet rarely stated assumption that the word should be treated as the primary meaning-bearing unit of language. This assumption likely stems from the historical and current dominance of English<sup>3</sup> as the language of study in NLP (Bender, 2011; Joshi et al., 2020), and the fact that in English, many words do in fact consist of only a single morpheme. English and Standard Mandarin Chinese<sup>4</sup> are prime examples of analytical languages where the average number of morphemes per word is low and for which existing neural representations such as BERT work very well (Peters et al., 2018; Devlin et al., 2019; Zhang et al., 2019).

## 1.2 Novel Contributions

Existing neural representations are insufficient (§1) for the thousands of languages which lack corpora. In this work, we take up this challenge,<sup>5</sup> surveying existing NLP methods for representing words (§2) and presenting a robust technique (§3) for constructing interpretable word embeddings from morphologically analyzed examples (such as interlinear glosses) for all of the world's languages, even when no corpus exists, and show how linguistic information encoded in these vectors can be successfully recovered.

As the primary contribution of this work, we present extensive proof-of-concept of our model gracefully handling immense morphological variety and hierarchical linguistic structures using complex examples that include concatenation and zero inflection (§4.1), circumfixation (§4.2), fusion (§4.3), polysynthesis (§4.4), agglutination (§4.5), infixation (§4.6), reduplication (§4.7), and templatic morphology (§4.8).

## 2 Existing Word Representations are Insufficient for Most Languages

Computational processing of natural language requires practical digital representations of the words of a language. We survey existing methods for representing words, arguing that while existing word representations work well for high resource ana-

lytic languages like English, existing representations are insufficient for effectively representing morphologically complex words in thousands of languages for which large corpora do not exist.

#### 2.1 Representing characters as integers

Oettinger (1954, ch. 2, p. 11), in the very first Ph.D. granted in the field of NLP, defined a word as "any string of letters preceded and followed by a space or a punctuation mark," and stored each word in an electronic dictionary as a sequence of characters, with each character represented digitally as a 5-bit integer. Nearly seventy years later, with relatively minor variations, this definition is still widely used in the NLP research community. Most digital word representations incorporate this technique, storing each character (or Unicode codepoint, as Clark et al., 2022, do) in a word as a multi-bit integer.

## 2.2 Representing words as feature bundles

During the 1960s through the early 1990s, most NLP systems utilized a knowledge-based paradigm in which words were represented as complex bundles of linguistic features, which were subsequently processed using linguistically-motivated rules (Hutchins, 1986). Finite-state morphological analyzers (Beesley and Karttunen, 2003) can be used to segment words into sequences of component morphemes; such segmentations can include explicit linguistic features such as case, number, and mood in addition to morpheme identity. Another modern example of this type of linguistically feature-rich word representation can be seen in the attribute-value matrices (AVMs) of Head-driven Phrase Structure Grammars (HPSG; Pollard and Sag, 1994). Such linguistically-based feature bundle representations can in principle work with any language, regardless of corpus size or morphological characteristics, but must be constructed by an expert linguist for each language, and do not naturally fit with many existing neural techniques.

## 2.3 Representing words as integers

The development of large digital corpora (primarily in English) and the rise of empirical approaches to NLP in the late 1980s and early 1990s, led to widespread use of statistical language models and translation models (see Church and Mercer, 1993; Manning and Schütze, 1999; Koehn, 2010). When implementing these statistical models, it is often convenient to map each word type to an integer,

<sup>&</sup>lt;sup>3</sup>ISO 639-3: *eng*, an analytic language in the Germanic branch of the Indo-European language family

<sup>&</sup>lt;sup>4</sup>ISO 639-3: *cmn*, an analytic language in the Sinitic branch of the Sino-Tibetan language family

<sup>&</sup>lt;sup>5</sup>"It is better to address the core scientific challenges than to continue to look for easy pickings that are no longer there." (Church, 2011)

allowing these integer word representations to directly serve as indices into probability tables (see for example §5 of Brown et al., 1993). A special integer value (often zero) is typically reserved to represent all words not seen during training.

While representing words as integers is efficient in its use of RAM, it suffers from a serious short-coming first observed by Bull et al. (1955), namely that no semantic, syntactic, or morphological information is encoded in the word representation (for example, *dog* and *dogs* are treated as completely unrelated word types). This problem is seriously exacerbated in languages with rich morphology, as productive derivational and inflectional morphology may result in extremely large numbers of closely-related word types, few of which are likely to appear in corpora. Schwartz et al. (2020a), for example, found that in one polysynthetic language, approximately every other word in running text will have never been previously seen.

## 2.4 Representing subwords as integers

Unsupervised techniques can be used to automatically segment words into sequences of shorter subword tokens generally longer than the character but shorter than the word. These techniques include approaches such as Morfessor (Creutz and Lagus, 2002; Smit et al., 2014) designed to segment words into units approximating morphemes, and compression-based subword segmentation techniques such as BPE (Sennrich et al., 2016; Wu et al., 2016; Kudo and Richardson, 2018). Most neural NLP systems in broad use today utilize integer representations of unsupervised subword tokens for both input and output.

This approach is more successful at representing words in languages with highly productive morphology than the integer word representations described in §2.3. When corpus sizes are small or nonexistent, however, as is the case for most of the world's languages, insufficient training signal exists to reliably train high-quality unsupervised subword segmentation. This problem can be mitigated through the use of a linguistically-based finite-state morphological analyzer (§2.2) for word segmentation instead of unsupervised segmentation methods (Park et al., 2021).

## 2.5 Representing (word or subword) types as embeddings

Distributed representations (Hinton et al., 1986), also called continuous representations and word

embeddings, represent each word as a point embedded in a high-dimensional vector space. When feed-forward or recurrent neural networks are trained as language models with the task of predicting the next element in a word sequence or a subword sequence, a side effect of the training process is a table of embeddings which can be indexed by the integer representation corresponding to each word (§2.3) or subword (§2.4) type. Other techniques for learning context-independent vector representations for each type include word2vec (Mikolov et al., 2013a) and GloVe (Pennington et al., 2014).

## 2.6 Representing (word or subword) tokens as embeddings

More recent neural techniques such as ELMo (Peters et al., 2018), BERT (Devlin et al., 2019), and Canine (Clark et al., 2022) can be used to obtain a context-dependent vector representation for each word or subword token. ELMo uses convolutional techniques to generalize over character sequences within the word in conjunction with deep bidirectional recurrent neural networks, while BERT utilizes unsupervised subword tokenization techniques (§2.4) in conjunction with a transformer architecture (Vaswani et al., 2017). Canine treats Unicode codepoints as the subword unit.

Learned context-free word embeddings empirically appear to implicitly encode at least some syntactic and semantic information (Mikolov et al., 2013b). Substantial recent work, summarized by Rogers et al. (2020) indicates that contextualized word embeddings learned by BERT are even more successful at implicitly encoding syntactic, semantic, and possibly morphological information. Interpretability of these embeddings is a challenging problem which is far from solved.

While multilingual training, transfer, and anchoring methods have been shown in some cases to somewhat improve the quality of very low-resource word embeddings over monolingually-trained low-resource word embeddings (see, for example, Eder et al., 2021), such methods rely on digitized monolingual and bilingual resources that exist for only a few hundred languages. It remains the case that at present, training high quality word embeddings is dependent on the availability of large corpora (Alabi et al., 2020; Joshi et al., 2020; Wu and Dredze, 2020; Budur et al., 2020; Michel et al., 2020) consisting of tens or hundreds of millions of

tokens, which are available for at most a few hundred languages (see §1).

## 2.7 Linguistically-informed word embeddings

No existing word representation is capable of robustly representing words in all of the world's languages regardless of corpus size and morphological characteristics. The existing representation that comes closest to meeting these needs is Linguistically Informed Multi-Task BERT (LIMIT-BERT Zhou et al., 2020b), a semi-supervised approach in which a trained parser (Zhou et al., 2020a) is used to annotate large unlabelled corpora. During LIMIT-BERT training, these silver linguistic annotations (part-of-speech tags, constituency trees, and dependency trees) are used along with the words themselves to train contextualized embeddings on five parsing-related tasks.

Unlike the embeddings learned by LIMIT-BERT, the representations we propose are explicitly interpretable by design, allowing for direct recovery of any linguistic features encoded in our word embeddings. Unlike LIMIT-BERT, our approach can produce high-quality word embeddings in the presence of arbitrarily complex morphology and in the absence of a corpus.

## 3 Embedding and retrieving rich linguistic information

As established in §1, there are thousands of languages which lack the large corpora needed for reliably training neural language models such as BERT. For many of these cases, the size of corpora may be very small or even nonexistent. While multilingual and bootstrapping approaches certainly have a role to play, we ought not ignore the rich linguistic information embedded in morphological analyses.

Essentially every language that is even partly documented has numerous such analyses in the form of interlinear glossed text (ILGs) created by expert linguists. Instead of relying on neural networks to induce linguistic patterns by processing massive corpora, we argue that for more than 6000 so-called "low-resource" languages, a more fruitful method for initializing meaningful word and subword embeddings is by directly embedding the rich linguistic information included in the morphological analyses found in ILGs and (when they exist) other morphologically analyzed corpora.

#### 3.1 Word Embedding Desiderata

We argue that the following desiderata are necessary in order to fulfill the use case of establishing meaningful word embeddings for all languages, even in the absence of any corpus. The representation must easily model words from polysynthetic languages, agglutinative languages, fusional languages, and isolating languages equally well, naturally incorporating any and all linguistic features which may be present in an interlinear gloss or available from other external resources. The representation must model words in ultra-low-resource settings where corpus sizes are very small or even non-existent just as well as it handles words in highresource settings with very large corpora. Finally, the representation must be interpretable; all linguistic features encoded in the resulting word embeddings should easily retrievable from the word embeddings.

## 3.2 Tensor Product Representation

To satisfy the word representation desiderata specified in §3.1, we utilize the Tensor Product Representation (TPR) proposed by Smolensky (1990). The use of TPRs provides a principled way of representing hierarchical symbolic information from external resources such as interlinear glosses or morphological analyzers into vector spaces, such as those used as the input and output domains of neural networks. The nature of TPRs enable simple linear algebra operations to be used to easily and fully recover this symbolic structure, including its compositional structure.

Constructing a TPR for a linguistic unit (such as a morpheme or a word) begins by decomposing the symbolic structure of that unit into *roles* and *fillers*. Each role represents a linguistic feature, while each filler represents the actual value of that feature.

The symbolic structure of a word is then represented as the *bindings* of fillers to roles for all feature-value pairs associated with that unit. Once decomposed, both roles and fillers are embedded into a vector space such that all roles are linearly independent from one another. Let b be a list of ordered pairs (i,j) representing filler i (with embedding vector  $\hat{\mathbf{f}}_i$ ) being bound to role j (with embedding vector  $\hat{\mathbf{r}}_j$ ). The *tensor product representation*  $\mathbf{T}$  of the information is then given by

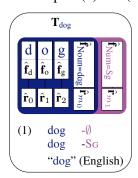
$$\mathbf{T} = \sum_{(i,j)\in b} \hat{\mathbf{f}}_i \otimes \hat{\mathbf{r}}_j \in \mathbb{R}^d \otimes \mathbb{R}^n.$$
 (1)

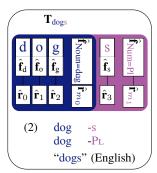
### 3.3 Constructing a TPR from an ILG

Our use of TPRs to represent ILGs is meant to be agnostic to linguistic theory. Considerable flexibility is available to the computational linguist in determining exactly how to map linguistic features from an ILG into the structure of a TPR. For example, one TPR design choice might involve linguistic features such as noun case or verb mood serving as roles, while the corresponding fillers represent actual values of those features, such as associative case or indicative mood.

For the sake of expositional simplicity in presenting a multilingual and typologically diverse set of linguistic examples (and without loss of generality), in Examples (1) and (2) below and in §4 we opt for a simplistic linguistic mapping where each TPR role represents a (grapheme or morpheme) position within the word and where the corresponding TPR fillers represent (grapheme or morpheme) identity at that position. Concretely, given a word comprised of  $\ell$  graphemes and m morphemes,  $\hat{\mathbf{r}}_i$ and  $\hat{\mathbf{r}}_{m_i}$  are one-hot<sup>6</sup> vectors respectively representing grapheme position i (where  $0 \le i < \ell$ ) and morpheme position j (where  $0 \le j < m$ ) within the word. For each linguistic element (grapheme or morpheme)  $\gamma$  in the language,  $\hat{\mathbf{f}}_{\gamma}$  is a vector<sup>7</sup> representing that element.

We now illustrate how morpheme and word embeddings can be constructed from interlinear glosses, using the English words 'dog' and 'dogs as Examples (1) and (2), respectively.





Each example is shown within a rounded rectangle; the example number and interlinear gloss are found at the bottom of the rounded rectangle, while a visualization of the TPR is shown at the top of the rectangle. At the top of each example is a label for the resulting word embedding. Colors are used to differentiate morpheme positions within the word.

In Example (1),  $\hat{\mathbf{r}}_0$  is a one-hot vector representing the initial grapheme position within the word, and  $\hat{\mathbf{f}}_d$  is a one-hot vector representing the English letter 'd'. The outer product  $\hat{\mathbf{r}}_0 \otimes \hat{\mathbf{f}}_d$  now represents a one-hot matrix encoding that the grapheme at position 0 is the English letter 'd'. Applying Equation (1), we add together three one-hot matrices  $(\hat{\mathbf{r}}_0 \otimes \hat{\mathbf{f}}_d + \hat{\mathbf{r}}_1 \otimes \hat{\mathbf{f}}_o + \hat{\mathbf{r}}_2 \otimes \hat{\mathbf{f}}_g)$ , to obtain a sparse matrix that encodes the surface form of the morpheme 'dog.' Similarly,  $\hat{\mathbf{r}}_{m_0} \otimes \hat{\mathbf{f}}_{\text{Noun=dog}}$  encodes that the identity of the initial morpheme in Example (1) is the noun 'dog.'

Recursive applications of Equation (1) result in multi-dimensional tensors  $T_{dog}$  (encoding the surface form and morpheme identity of each morpheme in the word 'dog') and  $T_{dogs}$  (encoding the surface form and morpheme identity of each morpheme in the word ''dogs').

## 3.4 Dense vectors from TPRs

Depending on how much linguistic information is encoded, each TPRs may consist of approximately  $10^3$  to  $10^9$  floating point values per tensor. Tensors of this size are far too large to be directly usable as neural word representations. It is therefore necessary to map each sparse TPR into an equivalent dense vector representation. Any of several existing techniques may be used to achieve this task; for simplicity in our work to date, we make use of an autoencoder. The autoencoder is trained using a dictionary of word or morpheme TPRs. The trained autoencoder can be used to encode a lowdimensional vector from a high-dimensional tensor by running the tensor through the first half of the autoencoder, and can be used to reconstitute the high-dimensional tensor from a vector by running the vector though the latter half of the autoencoder. For additional details, see Appendix A.

## 4 Supporting full linguistic diversity

We now demonstrate the broad applicability of our technique for encoding rich linguistic information from morphologically analyses such as ILGs using examples drawn from a typologically diverse set of polysynthetic, agglutinative, fusional, and analytic languages. The following examples include prefixes, suffixes, infixes, circumfixes, templatic morphemes, derivational morphemes, inflectional morphemes, and reduplication. The notation in the

<sup>&</sup>lt;sup>6</sup>In the general case, role vectors need not necessarily be one-hot.

<sup>&</sup>lt;sup>7</sup>For simplicity in our case, these filler vectors are one-hot. In the general case, filler vectors need not necessarily be one-hot, and may be separately pre-trained grapheme or morpheme embeddings if desired.

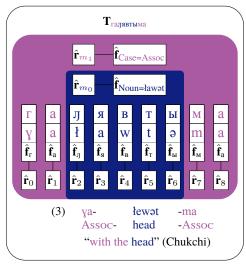
following examples follows the conventions established in §3.3.

## 4.1 Concatenative morphology and zero inflection in English

Concatenative morphology is extremely common cross-linguistically. Examples (1) and (2) in §3.3 demonstrate basic concatenative morphology in the English words 'dog' and 'dogs'. Example (1) illustrates that linguistic features of a word can be encoded even when those features are not explicitly marked in the surface form of the word. In Example (1), the tensor  $\mathbf{T}_{dog}$  explicitly encodes the null singular morpheme - $\emptyset$  marking number as singular in the word 'dog,' just as the morpheme -s marks number as plural in the word 'dogs in Example (2).' Unlike existing representations discussed in §2,  $\mathbf{T}_{dog}$  and  $\mathbf{T}_{dogs}$  are clearly distinguishable as variant inflections of the same root word.

#### 4.2 Circumfixes in Chukchi

The Chukchi<sup>8</sup> word галявтыма is composed of a noun root morpheme ławət and an inflectional circumfix ya...ma. The tensor  $T_{\text{галявтыма}}$  is a TPR that represents this word, *explicitly including* all information shown in Example (3):

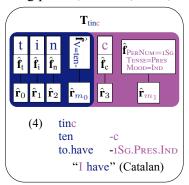


The individual characters positions in the word comprise roles  $\hat{\mathbf{r}}_0$  through  $\hat{\mathbf{r}}_8$ , while the characters (and respective phonemes) at those respective positions comprise fillers  $\hat{\mathbf{f}}_\Gamma$ ,  $\hat{\mathbf{f}}_a$ ,  $\hat{\mathbf{f}}_J$ ,  $\hat{\mathbf{f}}_B$ ,  $\hat{\mathbf{f}}_T$ ,  $\hat{\mathbf{f}}_{bI}$ , and  $\hat{\mathbf{f}}_M$  that encode character and phoneme identity. Roles  $\hat{\mathbf{r}}_{m_0}$  and  $\hat{\mathbf{r}}_{m_1}$  represent morpheme positions within the word, and are respectively filled by  $\hat{\mathbf{f}}_{Noun=lawet}$  (denoting the identity of the root morpheme) and

 $\hat{\mathbf{f}}_{\text{Case=Assoc}}$  (denoting the identity of the circumfix morpheme marking associative case).

## 4.3 Fusional suffixes in Catalan

Fusional morphology is also common crosslinguistics, as we can see in the Catalan<sup>9</sup> word tinc in Example (4), which is comprised only of only a verb root ten- 'to have' and a single inflectional suffix marking person, number, tense, and mood.



## 4.4 Polysynthesis with derivational and inflectional suffixes in Akuzipik

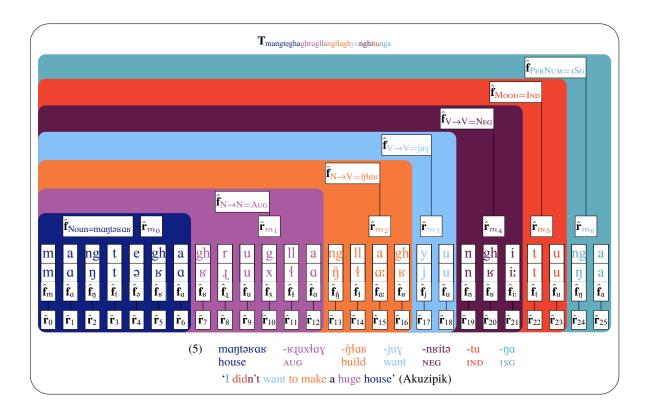
Productive derivational and inflectional suffixes are pervasive in the polysynthetic languages of the Inuit-Yupik language family. Words with 2-5 derivational morphemes are very common, often representing in a single word what in English would be represented by an entire clause or sentence.

The Akuzipik<sup>10</sup> word mangteghaghrugllangllaghyunghitunga shown in Example (5) can be translated into English as the sentence 'I didn't want to make a huge house' (Jacobson, 2001, pg. 43). The tensor  $T_{mangteghaghrugllangllaghyunghitunga}$  encodes the hierarchical structure of this word. grapheme position within the word is assigned a role  $(\hat{\mathbf{r}}_0 \dots \hat{\mathbf{r}}_{25})$ . For each of these grapheme position roles, a filler vector encodes the identity of the grapheme and corresponding phoneme at that position in the word  $(\hat{\mathbf{f}}_0 \dots \hat{\mathbf{f}}_{25})$ . The binding of grapheme position roles to grapheme filler vectors represents the first level of hierarchy in the TPR. The word is composed of 7 morphemes: a noun root mantakak, four derivational morphemes (-kjuxfay, -nglak, -juy, -nkitə) and two inflectional morphemes (-tu and -na). The subsequent levels of the TPR encode the identity, underlying form, surface form, and hierarchical scope of each

<sup>&</sup>lt;sup>8</sup>ISO 639-3: *ckt*, a polysynthetic language in the Chukotkan branch of the Chukotko–Kamchatkan language family

<sup>&</sup>lt;sup>9</sup>ISO 639-3: *cat*, a fusional language in the Romance branch of the Indo-European language family

<sup>&</sup>lt;sup>10</sup>ISO 639-3: *ess*, a polysynthetic language in the Yupik branch of the Inuit-Yupik-Unangan language family



morpheme. The resulting word representation is compositional and easily interpretable.

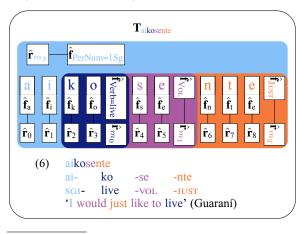
By inspecting the resulting tensor, the following structure of the word can be clearly observed:

- The noun root for 'house' mantakak is modified by the augmentatitive derivational morpheme -ktuxtay, resulting in an extended noun stem meaning 'big house' spanning grapheme positions 0 through 12.
- The resulting extended noun stem (mantəkak.juxtay) is verbalized by the derivational morpheme -ŋtak, resulting in an extended verb stem meaning 'to build a big house' spanning grapheme positions 0 through 16.
- The resulting extended verb stem (maŋtəʁaʁ.tuxɬaṅ̞ɬaʁ) is modified by the derivational
  morpheme -juy, resulting in an extended verb
  stem meaning 'to want to build a big house'
  spanning grapheme positions 0 through 18.
- The resulting extended verb stem (maŋtəʁaʁ-nuxtaŋtaʁjuɣ) is modified by the negating
  derivational morpheme -nʁitə), resulting in
  an extended verb stem meaning 'to not want
  to build a big house' spanning grapheme positions 0 through 21.
- The resulting extended verb stem (mantəkakıuxlanlar) is marked as being in

the indicative mood by the inflectional morpheme -tu and as having a first person singular subject by the inflectional morpheme -ŋa, resulting in the fully inflected word spanning grapheme positions 0 through 25.

## 4.5 Agglutination in Guaraní

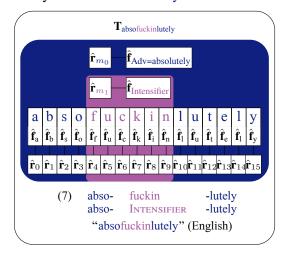
In the Guaraní<sup>11</sup> word aikosente shown in Example (6), the verb root ko 'to live' is modified in agglutinative manner by two suffixes (-se and -nte) and one inflectional prefix (ai-) which indicates a first person singular subject. Note that unlike the preceding example, which also encoded phoneme identity, in this example character fillers encode only character identity.



<sup>&</sup>lt;sup>11</sup>ISO 639-3: *gug*, an agglutinative language in the Tupian language family

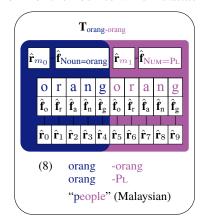
### 4.6 Infixation in English

Linguistic features such as infixes that are attested but relatively rare can also be included with no difficulty. Infixes are morphemes that break a given stem and appear inside it. In Seri, <sup>12</sup> for example, infixation after the first vowel in the root is used to mark number agreement. In Example (7), we observe an example of expletive infixation in English (McCarthy, 1982) with the infix fuckin serving to intensify the adverb absolutely.



## 4.7 Reduplication in Malaysian

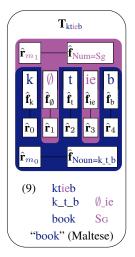
The Malaysian<sup>13</sup> word orang-orang 'people', is formed through reduplication of the noun root orang 'person'. Unlike in previous examples, in which morpheme fillers encoded underlying lexical form in addition to morpheme surface form and identity, in Example (8), the plural morpheme has no inherent underlying lexical form separate from the morpheme identity (Num=Pl). Instead the surface form of the plural morpheme (here, orang) is formed through reduplication, duplicating the form of the noun to which it attaches.

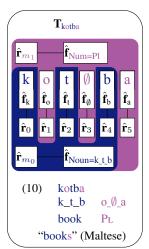


<sup>12</sup>ISO 639-3: *sei* a language isolate in north-west Mexico <sup>13</sup>ISO 639-3: *zsm*, a language in the Malayo-Polynesian branch of the Austronesian language family

### 4.8 Templatic morphology in Maltese

Our representation can easily encode nonconcatenative morphology such as that seen in the Maltese<sup>14</sup> words ktieb 'book' and kotba 'books.'





The noun root  $k_t_b$  acts as a template whose slots are filled by the vowels in the inflectional singular morpheme  $\emptyset$ \_ie (in Example (9)) or plural morpheme  $o_0\emptyset$ \_a (in Example (10)).

## 5 Conclusion

While corpora of anything greater than trivial size exist only for a few hundred languages (§1), morphologically analyzed examples in the form of interlinear glosses exist for essentially every human language. The vast array of human languages include a rich variety of morphological phenomenon that are not easily handled by existing word embedding methods (§2). This work presents a straightforward mechanism whereby meaningful, linguistically interpretable word and morpheme embeddings can be created for any word in any language (§3–§4). We have demonstrated the applicability of our method using linguistic examples of concatenation and zero inflection (§4.1), circumfixation (§4.2), fusion (§4.3), polysynthesis (§4.4), agglutination (§4.5), infixation (§4.6), reduplication (§4.7), and templatic morphology (§4.8).

In addition to their direct use in future research involving language documentation and revitalization, we anticipate that embeddings created using the methods described in this work may provide an important initial step in bootstrapping vastly multilingual models capable of embedding words from thousands of languages.

 $<sup>^{14}</sup>$ ISO 639-3: mlt, a templatic language in the Semitic language family

## Acknowledgements

This work was initially developed during the 2019 JSALT workshop on Neural Polysynthetic Language Modelling (Schwartz et al., 2020b) in Montréal, Canada. We wish to express our appreciation to the organizers, sponsors, and hosts of the 2019 JSALT workshop. We wish to express our deep respect and thanks to the many peoples whose languages we present in the examples in this paper. We wish to acknowledge and honor the Indigenous peoples on whose lands we live and work, both at Montréal and at our individual universities.

Our code is at https://github.com/ neural-polysynthetic-language-modelling/ iiksiin and the scripts we used to run our code are at https://github.com/ neural-polysynthetic-language-modelling/ iiksiin.experiment

#### References

- Jesujoba Alabi, Kwabena Amponsah-Kaakyire, David Adelani, and Cristina España-Bonet. 2020. Massive vs. curated embeddings for low-resourced languages: the case of Yorùbá and Twi. In Proceedings of the 12th Language Resources and Evaluation Conference, pages 2754–2762, Marseille, France. European Language Resources Association.
- Kenneth R. Beesley and Lauri Karttunen. 2003. *Finite State Morphology*. CSLI Studies in Computational Linguistics. CSLI Publications, Stanford, California.
- Emily M. Bender. 2011. On achieving and evaluating language-independence in NLP. *Linguistic Issues in Language Technology*, 6(3):1–26.
- Balthasar Bickel and Johanna Nichols. 2013. Inflectional synthesis of the verb. In Matthew S. Dryer and Martin Haspelmath, editors, *The World Atlas of Language Structures Online*. Max Planck Institute for Evolutionary Anthropology, Leipzig.
- Peter F. Brown, Stephen A. Della Pietra, Vincent J. Della Pietra, and Robert L. Mercer. 1993. The mathematics of statistical machine translation: Parameter estimation. *Computational Linguistics*, 19(2):263–311.
- Emrah Budur, Rıza Özçelik, Tunga Gungor, and Christopher Potts. 2020. Data and Representation for Turkish Natural Language Inference. In *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP)*, pages 8253–8267, Online. Association for Computational Linguistics.
- William E. Bull, Charles Africa, and Daniel Teichroew. 1955. Some problems of the "word". In William N.

- Locke and A. Donald Booth, editors, *Machine Translations of Languages*. Greenwood Press, Westport, Connecticut.
- Emily Chen and Lane Schwartz. 2018. A morphological analyzer for St. Lawrence Island / Central Siberian Yupik. In *Proceedings of the 11th Language Resources and Evaluation Conference (LREC 2018)*, Miyazaki, Japan. European Languages Resources Association (ELRA).
- Kenneth Church. 2011. A pendulum swung too far. Linguistic Issues in Language Technology, 6(3):1–27.
- Kenneth W. Church and Robert L. Mercer. 1993. Introduction to the special issue on computational linguistics using large corpora. *Computational Linguistics*, 19(1):1–24.
- Jonathan H. Clark, Dan Garrette, Iulia Turc, and John Wieting. 2022. Canine: Pre-training an efficient tokenization-free encoder for language representation. Transactions of the Association for Computational Linguistics, 10:73–91.
- Alexis Conneau, Kartikay Khandelwal, Naman Goyal, Vishrav Chaudhary, Guillaume Wenzek, Francisco Guzmán, Edouard Grave, Myle Ott, Luke Zettlemoyer, and Veselin Stoyanov. 2020. Unsupervised cross-lingual representation learning at scale. In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, pages 8440–8451, Online. Association for Computational Linguistics.
- Mathias Creutz and Krista Lagus. 2002. Unsupervised discovery of morphemes. In *Proceedings of the ACL-02 Workshop on Morphological and Phonological Learning*, pages 21–30. Association for Computational Linguistics.
- Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. 2019. BERT: Pre-training of deep bidirectional transformers for language understanding. In Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long and Short Papers), pages 4171–4186, Minneapolis, Minnesota. Association for Computational Linguistics.
- Matthew S. Dryer. 2013. Prefixing vs. suffixing in inflectional morphology. In Matthew S. Dryer and Martin Haspelmath, editors, *The World Atlas of Language Structures Online*. Max Planck Institute for Evolutionary Anthropology, Leipzig.
- Matthew S. Dryer and Martin Haspelmath, editors. 2013. *World Atlas of Language Structures*. Max Planck Institute for Evolutionary Anthropology, Leipzig.
- Tobias Eder, Viktor Hangya, and Alexander Fraser. 2021. Anchor-based bilingual word embeddings for low-resource languages. In *Proceedings of the*

- 59th Annual Meeting of the Association for Computational Linguistics and the 11th International Joint Conference on Natural Language Processing (Volume 2: Short Papers), pages 227–232, Online. Association for Computational Linguistics.
- Coleman Haley and Paul Smolensky. 2020. Invertible tree embeddings using a cryptographic role embedding scheme. In *Proceedings of the 28th International Conference on Computational Linguistics*, pages 3671–3683, Barcelona, Spain (Online). International Committee on Computational Linguistics.
- Michael A. Hedderich, Lukas Lange, Heike Adel, Jannik Strötgen, and Dietrich Klakow. 2021. A survey on recent approaches for natural language processing in low-resource scenarios. In *Proceedings of the 2021 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*, pages 2545–2568, Online. Association for Computational Linguistics.
- G.E. Hinton, J.L. McClelland, and D.E. Rumelhart. 1986. Distributed representations. In *Parallel distributed processing: Explorations in the microstructure of cognition*, volume 1: Foundations. MIT Press.
- W. John Hutchins. 1986. *Machine Translation: Past, Present, Future*. Computers and Their Applications. Ellis Horwood.
- Steven A. Jacobson. 2001. A Practical Grammar of the St. Lawrence Island/Siberian Yupik Eskimo Language, Preliminary Edition, 2nd edition. Alaska Native Language Center, Fairbanks, Alaska.
- Pratik Joshi, Sebastin Santy, Amar Budhiraja, Kalika Bali, and Monojit Choudhury. 2020. The state and fate of linguistic diversity and inclusion in the NLP world. In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, pages 6282–6293, Online. Association for Computational Linguistics.
- Philipp Koehn. 2010. *Statistical Machine Translation*. Cambridge University Press, Cambridge, UK.
- Taku Kudo and John Richardson. 2018. SentencePiece: A simple and language independent subword tokenizer and detokenizer for neural text processing. In Proceedings of the 2018 Conference on Empirical Methods in Natural Language Processing: System Demonstrations, pages 66–71, Brussels, Belgium. Association for Computational Linguistics.
- Christopher D. Manning and Hinrich Schütze. 1999. Foundations of Statistical Natural Language Processing. MIT Press, Cambridge, Massachusetts.
- John J. McCarthy. 1982. Prosodic structure and expletive infixation. *Language*, 58(3):574–590.
- Leah Michel, Viktor Hangya, and Alexander Fraser. 2020. Exploring bilingual word embeddings for Hiligaynon, a low-resource language. In *Proceedings of*

- the 12th Language Resources and Evaluation Conference, pages 2573–2580, Marseille, France. European Language Resources Association.
- Tomás Mikolov, Kai Chen, Greg Corrado, and Jeffrey Dean. 2013a. Efficient estimation of word representations in vector space. In 1st International Conference on Learning Representations, ICLR 2013, Scottsdale, Arizona, USA, May 2-4, 2013, Workshop Track Proceedings.
- Tomas Mikolov, Wen-tau Yih, and Geoffrey Zweig. 2013b. Linguistic regularities in continuous space word representations. In *Proceedings of the 2013 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*, pages 746–751, Atlanta, Georgia. Association for Computational Linguistics.
- Anthony Oettinger. 1954. A Study for the Design of an Automatic Dictionary. Ph.D. thesis, Harvard University, Cambridge, Massachusetts.
- Hyunji Hayley Park, Katherine J. Zhang, Coleman Haley, Kenneth Steimel, Han Liu, and Lane Schwartz. 2021. Morphology Matters: A Multilingual Language Modeling Analysis. *Transactions of the Association for Computational Linguistics*, 9:261–276.
- Jeffrey Pennington, Richard Socher, and Christopher Manning. 2014. GloVe: Global vectors for word representation. In *Proceedings of the 2014 Conference on Empirical Methods in Natural Language Processing (EMNLP)*, pages 1532–1543, Doha, Qatar. Association for Computational Linguistics.
- Matthew E. Peters, Mark Neumann, Mohit Iyyer, Matt Gardner, Christopher Clark, Kenton Lee, and Luke Zettlemoyer. 2018. Deep contextualized word representations. In *Proceedings of the 2018 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long Papers)*, pages 2227–2237, New Orleans, Louisiana. Association for Computational Linguistics.
- Carl Pollard and Ivan A. Sag. 1994. Head-Driven Phrase Structure Grammar. University of Chicago Press.
- Edoardo Maria Ponti, Helen O'Horan, Yevgeni Berzak, Ivan Vulić, Roi Reichart, Thierry Poibeau, Ekaterina Shutova, and Anna Korhonen. 2019. Modeling language variation and universals: A survey on typological linguistics for natural language processing. *Computational Linguistics*, 45(3):559–601.
- Anna Rogers, Olga Kovaleva, and Anna Rumshisky. 2020. A primer in BERTology: What we know about how BERT works. *Transactions of the Association for Computational Linguistics*, 8:842–866.
- Lane Schwartz, Sylvia Schreiner, and Emily Chen. 2020a. Community-focused language documentation in support of language education and revitalization for St. Lawrence Island Yupik. *Études Inuit Studies*, 43(1–2):291–312.

Lane Schwartz, Francis Tyers, Lori Levin, Christo Kirov, Patrick Littell, Chi-kiu Lo, Emily Prud'hommeaux, Hyunji Hayley Park, Kenneth Steimel, Rebecca Knowles, Jeffrey Micher, Lonny Strunk, Han Liu, Coleman Haley, Katherine J. Zhang, Robbie Jimmerson, Vasilisa Andriyanets, Aldrian Obaja Muis, Naoki Otani, Jong Hyuk Park, and Zhisong Zhang. 2020b. Neural polysynthetic language modelling. Final Report of the Neural Polysynthetic Language Modelling Team at the 2019 Frederick Jelinek Memorial Summer Workshop.

Rico Sennrich, Barry Haddow, and Alexandra Birch. 2016. Neural machine translation of rare words with subword units. In *Proceedings of the 54th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 1715–1725, Berlin, Germany. Association for Computational Linguistics.

Peter Smit, Sami Virpioja, Stig-Arne Grönroos, and Mikko Kurimo. 2014. Morfessor 2.0: Toolkit for statistical morphological segmentation. In *Proceedings of the Demonstrations at the 14th Conference of the European Chapter of the Association for Computational Linguistics*, pages 21–24, Gothenburg, Sweden. Association for Computational Linguistics.

Paul Smolensky. 1990. Tensor product variable binding and the representation of symbolic structures in connectionist systems. *Artificial Intelligence*, 46:159–216.

Laurent Vannini and Hervé Le Crosnier, editors. 2012. Net.lang: Towards the Multilingual Cyberspace. C&F éditions.

Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Ł ukasz Kaiser, and Illia Polosukhin. 2017. Attention is all you need. In *Advances in Neural Information Processing Systems*, volume 30. Curran Associates, Inc.

Shijie Wu and Mark Dredze. 2020. Are all languages created equal in multilingual BERT? In *Proceedings* of the 5th Workshop on Representation Learning for NLP, pages 120–130, Online. Association for Computational Linguistics.

Yonghui Wu, Mike Schuster, Zhifeng Chen, Quoc V. Le, Mohammad Norouzi, Wolfgang Macherey, Maxim Krikun, Yuan Cao, Qin Gao, Klaus Macherey, Jeff Klingner, Apurva Shah, Melvin Johnson, Xiaobing Liu, Lukasz Kaiser, Stephan Gouws, Yoshikiyo Kato, Taku Kudo, Hideto Kazawa, Keith Stevens, George Kurian, Nishant Patil, Wei Wang, Cliff Young, Jason Smith, Jason Riesa, Alex Rudnick, Oriol Vinyals, Greg Corrado, Macduff Hughes, and Jeffrey Dean. 2016. Google's neural machine translation system: Bridging the gap between human and machine translation. *CoRR*, abs/1609.08144.

Zhengyan Zhang, Xu Han, Zhiyuan Liu, Xin Jiang, Maosong Sun, and Qun Liu. 2019. ERNIE: Enhanced language representation with informative entities. In *Proceedings of the 57th Annual Meeting of* 

*the Association for Computational Linguistics*, pages 1441–1451, Florence, Italy. Association for Computational Linguistics.

Junru Zhou, Zuchao Li, and Hai Zhao. 2020a. Parsing all: Syntax and semantics, dependencies and spans. In *Findings of the Association for Computational Linguistics: EMNLP 2020*, pages 4438–4449, Online. Association for Computational Linguistics.

Junru Zhou, Zhuosheng Zhang, Hai Zhao, and Shuailiang Zhang. 2020b. LIMIT-BERT: Linguistics informed multi-task BERT. In Findings of the Association for Computational Linguistics: EMNLP 2020, pages 4450–4461, Online. Association for Computational Linguistics.

## A Unbinding

The core operation in retrieving structure from a TPR is called *unbinding*. Exact unbinding requires linear independence of the roles; however, Haley and Smolensky (2020) present an accurate approximate unbinding strategy for even densely packed TPRs. In this work, we use self-addressing unbinding, as it is quick to compute and proved sufficiently accurate for our purposes. Self-addressing unbinding retrieves the filler  $\tilde{\mathbf{f}}_i$  for the role  $\hat{\mathbf{r}}_i$  by simply computing the inner product between the role vector and the TPR:

$$\tilde{\mathbf{f}}_i = \mathbf{T} \cdot \hat{\mathbf{r}}_i \tag{2}$$

This unbinding is exact if the role vectors are orthogonal to one another. In our case, since we have a fixed filler vocabulary, we were able to snap our unbindings to the filler with the highest cosine similarity to the unbound vector with sufficient accuracy to render this intrusion irrelevant. Other unbinding strategies involve computing an inverse or pseudoinverse of a matrix of role vectors to perform a change of basis and decrease the intrusion.

## A.1 Unbinding loss

In order to effectively train the autoencoder in §3.4, gold standard TPRs must be compared against predicted tensors reconstituted by the autoencoder. However, these tensors are very high dimensional. In initial experiments, we used mean squared error as a loss function, but we found this was unable to converge for auto-encoding sparse TPRs.

To enable effective training of the autoencoder, we therefore define a novel loss function that makes use of the information encoded in the TPR. We define a loss function called *unbinding loss* that examines the unbinding properties of a predicted

morpheme tensor to answer the question, "What filler is closest to the unbinding of each role in the TPR?"

Given a predicted tensor, the unbinding loss is computed by recursively unbinding roles until the leaves of the structure are reached – that is, unbind each role until the result of unbinding is a single vector (rather than a higher-order tensor). When this point is reached, we compute the cosine similarity between the result of unbinding and all the fillers in the vocabulary.

This similarity vector can be used to define a probability distribution over possible fillers through the use of a softmax. We take the logarithm of the result of this computation to obtain log-probabilities. We call this distribution P. We then treat each filler (in this case, each character) as a class, and compute the negative log-likelihood loss over this probability distribution.

As we consider tree-structured representations, the number of fillers needing to be checked is exponential with the depth of our representation. This difficulty could be overcome by parallelizing the independent matrix computations for the loss of all the position roles for a given morpheme, trading space for time. For more complex TPRs, a potential avenue would be to exploit the fact that most roles will be empty (and their unbindings thus a matrix of zeros) by replacing the loss computations for unbound roles with mean squared error (which need only push that part of the representation to 0).

## A.2 Unbinding loss example

Given a predicted tensor, the first step to computing the unbinding loss is recursively unbinding roles until the leaves of the structure are reached that is, unbind each role until the result of unbinding is a single vector (rather than a higher-order tensor). When this point is reached, we compute the cosine similarity between the result of unbinding and all the fillers in the vocabulary. For example, assume a depth-4 structure is encoded in a morpheme TPR T, where the fillers are character embeddings, the second level is left-to-right positional roles, the third level is morpheme identity, and the fourth level is left-to-right morpheme position in the word. If we want to see what is bound to the first position of the English dog morpheme in T, we would first unbind from T as follows (assuming self-addressing unbinding):

$$\mathbf{f}_{dog,1} = \mathbf{T} \cdot \hat{\mathbf{r}}_{m0} \cdot \hat{\mathbf{f}}_{Noun=dog} \cdot \hat{\mathbf{r}}_{1}$$
 (3)

We then get the vector of similarities  $\hat{\mathbf{s}}_{dog,1}$  between this filler and the each of character embedding vectors in the vocabulary matrix V as follows:

$$\hat{\mathbf{s}}_{dog,1} = \frac{\mathbf{f}_{dog,1} \cdot \mathbf{V}}{||\mathbf{f}_{dog,1}||\mathbf{V}^i \mathbf{V}^i} \tag{4}$$

where  $\mathbf{V}^{i}\mathbf{V}^{i}$  denotes the column-wise vector norm of the vocabulary matrix (using Einstein summation notation).

This similarity vector can be used to define a probability distribution over possible fillers through the use of a softmax. We take the logarithm of the result of this computation to obtain log-probabilities. We call this distribution P.

$$P = \log\left(\frac{e^{\hat{\mathbf{s}}_{dog,1}}}{\sum e^{\hat{\mathbf{s}}_{dog,1}}}\right) \tag{5}$$

We then treat each filler (in this case, each character) as a class, and compute the negative log-likelihood loss over this probability distribution. The resulting loss for the first character of *dog* being "d" is then

$$loss(\hat{\mathbf{s}}_{dog,1}, d) = -\hat{\mathbf{s}}_{dog,1,d} + \log(\sum_{j} e^{\hat{\mathbf{s}}_{dog,1,j}}).$$

$$(6)$$

If the Tensor this loss is computed over is exactly  $T_{\rm dog}$  or  $T_{\rm dogs}$ , then this loss term would be 0. If we instead considered the loss for the fourth character of the word being "s" in the Num=Pl morpheme, This would be 0 only for  $T_{\rm dogs}$ .

## A.3 Successfully recovering surface forms from vectors

To demonstrate the successful recovery of linguistic data from embeddings, we construct TPRs for a dictionary of 6372 unique Akuzipik morpheme surface forms obtained by applying the finite-state morphological analyzer of Chen and Schwartz (2018) on a selection of Akuzupik New Testament data from https://github.com/SaintLawrenceIslandYupik/digital\_corpus. Using TPRs constructed from these morphemes, we trained a 3-layer autoencoder with vector sizes of 64, 128, 256, and 512 using unbinding loss (§A.1) as the loss function. We then reconstructed the morpheme surface forms from the trained morpheme vectors. For

vector size of 64, the reconstructed morpheme surface form exactly matched the original morpheme surface form for 97.8% of the morphemes. For vector sizes of 128, 256, and 512, the morpheme surface form reconstruction accuracy was 100%.