

# How to Compute the Probability of a Word



Tiago Pimentel, Clara Meister

ETH

{tiago.pimentel, clara.meister}@inf.ethz.ch

## Abstract

Language models (LMs) estimate a probability distribution over strings in a natural language; these distributions are crucial for computing perplexity and surprisal in linguistics research. While we are usually concerned with measuring these values for words, most LMs operate over subwords. Despite seemingly straightforward, accurately computing probabilities over one unit given probabilities over the other requires care. Indeed, we show here that many recent linguistic studies have been incorrectly computing these values. This paper derives the correct methods for computing word probabilities, highlighting issues when relying on language models that use beginning-of-word (bow)-marking tokenisers, e.g., the GPT family. Empirically, we show that correcting the widespread bug in probability computations affects measured outcomes in sentence comprehension and lexical optimisation analyses.

 [tpimentelms/probability-of-a-word](https://github.com/tpimentelms/probability-of-a-word)  
 `pip install wordsprobability`

## 1 Introduction

Language models (LMs) define probability distributions. After being trained on language data, these models can be used to compute estimates of the probability of a sequence of characters  $\mathbf{c} \in \mathcal{C}^*$ , or of a word  $w_t \in \mathcal{W}$  in context  $\mathbf{w}_{<t} \in \mathcal{W}^*$ . While deriving such estimates is now rarely the explicit goal of training such models,<sup>1</sup> this use case is still critical in several fields. Estimating the probability of a sequence of characters, for instance, is necessary to compute a model’s perplexity; a core evaluation metric in LM training. Estimating the probability of a word in context is necessary to compute a **word’s surprisal**:  $-\log p(w_t | \mathbf{w}_{<t})$ , an important value in both psycho- and computational

<sup>1</sup>Rather, LMs have become known for their high performance on downstream natural language processing (NLP) tasks (Radford et al., 2019; Touvron et al., 2023).

TL;DR: How to correctly compute word probabilities

Given a word  $w$  in context  $\mathbf{w}_{<t}$ , let  $\mathbf{s}^w$  and  $\mathbf{s}^{\mathbf{w}_{<t}}$  be their respective subword sequences output by a tokeniser. Further, let:

$$p(\mathbf{s}^w | \mathbf{s}^{\mathbf{w}_{<t}}) = \prod_{i=1}^{|\mathbf{s}^w|} p(s_i^w | \mathbf{s}^{\mathbf{w}_{<t}} \circ \mathbf{s}_{<i}^w)$$

- LM with end-of-word marking tokeniser

$$p(w | \mathbf{w}_{<t}) = p(\mathbf{s}^w | \mathbf{s}^{\mathbf{w}_{<t}})$$

- LM with beginning-of-word marking tokeniser

$$p(w | \mathbf{w}_{<t}) = p(\mathbf{s}^w | \mathbf{s}^{\mathbf{w}_{<t}}) \frac{\sum_{\{s \in \mathcal{S}_{\text{bow}}\}} p(s | \mathbf{s}^{\mathbf{w}_{<t}} \circ \mathbf{s}^w)}{\underbrace{\sum_{\{s \in \mathcal{S}_{\text{bow}}\}} p(s | \mathbf{s}^{\mathbf{w}_{<t}})}_{\text{"bug" fix}}}$$

Figure 1: Equations for computing a word’s contextual probability  $p(w | \mathbf{w}_{<t})$  using a subword-based LM  $p(s_t | \mathbf{s}_{<t})$ .  $\mathcal{S}_{\text{bow}}$  is a subset of the tokeniser’s vocabulary marking beginnings of words. The “bug” fix can be computed for “free”, i.e., within a single model pass.

linguistics (Hale, 2001; Levy and Jaeger, 2007; Piantadosi et al., 2011; Pimentel et al., 2023a).

Notably, most recent LMs operate over **subwords** (Sennrich et al., 2016; Kudo and Richardson, 2018): sequences of characters that frequently occur together. This is done for both optimisation and efficiency reasons (Gallé, 2019; Mielke et al., 2021; Zouhar et al., 2023). Subwords, however, do not necessarily constitute actual words, as defined by a language’s lexicon.<sup>2</sup> At least superficially, converting from a probability distribution over subwords  $p(\mathbf{s})$  into one over characters  $p(\mathbf{c})$  or words  $p(\mathbf{w})$  appears straightforward. However, some technical details are easy to overlook. For example, several sequences of subwords  $\mathbf{s}$  can map to a single sequence of characters  $\mathbf{c}$ , implying an accurate computation of  $p(\mathbf{c})$  should marginalise over these options (Cao and Rimell, 2021).

<sup>2</sup>Despite the name, which we use out of convention, a subword need not strictly be a subunit of a word. For example, subwords can span multiple words, containing the markers used to delineate words, e.g., white spaces.

In this work, we discuss how to correctly compute a word’s contextual probability:  $p(w_t | \mathbf{w}_{<t})$ . This value’s computation depends on the choice of tokeniser used to define an LM’s vocabulary. When using an end-of-word (eow)-marking tokeniser, computing  $p(w_t | \mathbf{w}_{<t})$  is simple. However, when using a beginning-of-word (bow)-marking tokeniser, correctly computing this value is not as straightforward. We derive methods for these tokenisation schemes, which we present in Fig. 1. Since many widely-used LMs employ bow-marking tokenisers (e.g., the GPT models, Pythia, Mistral), this highlights a wide-spread “bug” in how most recent psycholinguistics and computational linguistics works compute word probabilities (present in, e.g., Oh and Schuler, 2023b; Wilcox et al., 2023a; Pimentel et al., 2023a; Shain et al., 2024).<sup>3</sup>

Empirically, we evaluate how correcting this computation affects the results of two prior empirical analyses: one on sentence comprehension and another on the lexicon’s communicative efficiency. While these studies’ conclusions do not change, we do observe statistically significant differences between the measured quantities when using the correct vs. buggy methods for computing word probabilities. We conclude this methodological choice may impact empirical analyses, and that future work should adopt these proposed corrections.

## 2 What is a Word?

Despite decades of discussion and debate, there is no single, widely accepted definition of what constitutes a word (Haspelmath, 2023). Typically, definitions are made with respect to some system within the language, such as its orthography, phonology, or grammar. As a concrete example, one can delineate words using the sound system of a language: if we assume words define the domain over which certain phonological processes operate (e.g., vowel harmony), we can delineate words based on those processes’ boundaries (Hall and Kleinhenz, 1999; Nespov and Vogel, 2007). Alternatively, one could define words as grammatical elements (e.g., a root plus affixes) that are cohesive, occur in a fixed order, and have a coherent meaning (Dixon and Aikhenvald, 2003). Notably grammatical and phonological words are non-isomorphic. For example, English hyphenated elements like *editor-in-chief* or *mother-in-law* are typically analysed as

<sup>3</sup>Concurrent work by Oh and Schuler (2024) points out this same issue and proposes a solution similar to **Bug Fix 1** (in our Theorem 2).

a single grammatical word that contains multiple phonological words (Dixon and Aikhenvald, 2003).

We abstain from this broader discussion here. While we use the definition common to natural language processing applications—where words are defined orthographically<sup>4</sup>—our methods only assume the existence of a deterministic set of rules for segmenting a string of characters into words.

## 3 Words and Distributions Over Them

Let  $\mathcal{W}$  be a lexicon—the (potentially infinite) set of all words in a language—and  $w \in \mathcal{W}$  a word in this lexicon. Further, let  $\mathbf{w} \in \mathcal{W}^*$  be a sequence of words;  $\mathcal{W}^*$  denotes the set of all finite-length word sequences. Now, assume distribution  $p$  describes the probability with which users of this language produce sequences  $\mathbf{w}$ . We can decompose these probabilities autoregressively as:

$$p(\mathbf{w}) = p(\text{eos} | \mathbf{w}) \prod_{t=1}^{|\mathbf{w}|} p(w_t | \mathbf{w}_{<t}) \quad (1)$$

where eos is a special end-of-sequence symbol that makes this probability distribution over  $\mathcal{W}^*$  valid.<sup>5</sup>

This paper is concerned with the proper method for computing the probability of a word in context, i.e.,  $p(w_t | \mathbf{w}_{<t})$ , using a pretrained language model. To this end, we first discuss its equivalence to other quantities, which will ultimately reveal a flaw in prior approaches to its computation. We start by defining a probability function  $\mathbb{P}_{\mathcal{W}}$ , which operates over sets of strings  $\Psi_{\mathcal{W}} \subseteq \mathcal{W}^*$ .

**Definition 1.** Given distribution  $p(\mathbf{w})$ , we define the probability function  $\mathbb{P}_{\mathcal{W}} : \mathcal{P}(\mathcal{W}^*) \rightarrow [0, 1]$ , which returns the probability of occurrence of any word sequence  $\mathbf{w} \in \Psi_{\mathcal{W}} \subseteq \mathcal{W}^*$ . As these events are disjoint,  $\mathbb{P}_{\mathcal{W}}(\Psi_{\mathcal{W}})$  can be defined as:

$$\mathbb{P}_{\mathcal{W}}(\Psi_{\mathcal{W}}) \stackrel{\text{def}}{=} \sum_{\mathbf{w} \in \Psi_{\mathcal{W}}} p(\mathbf{w}) \quad (2)$$

Now, let  $\circ$  denote concatenation (between either strings or sets of strings). For instance, we can write  $\mathbf{w} \circ \mathcal{W}^* = \{\mathbf{w} \circ \mathbf{w}' | \mathbf{w}' \in \mathcal{W}^*\}$  to represent the set of all strings with prefix  $\mathbf{w}$ . We can compute our desired conditional distribution as the quotient of two evaluations of  $\mathbb{P}_{\mathcal{W}}$ :

$$p(w | \mathbf{w}_{<t}) = \frac{\mathbb{P}_{\mathcal{W}}(\mathbf{w}_{<t} \circ w \circ \mathcal{W}^*)}{\mathbb{P}_{\mathcal{W}}(\mathbf{w}_{<t} \circ \mathcal{W}^*)} \quad (3)$$

<sup>4</sup>Orthographic words are defined as sequences of characters surrounded by white space or other special delimiters. One such delimiter is ‘, present in the English clitic ‘s.

<sup>5</sup>See Du et al. (2023) for a longer discussion on when probability distributions over  $\mathcal{W}^*$  are valid.

Note that this is a trivial invocation of the joint rule of probability: the conditional  $p(w \mid \mathbf{w}_{<t})$  is equal to the probability of observing prefix  $\mathbf{w}_{<t} \circ w$ —represented by  $\mathbb{P}_{\mathcal{W}}(\mathbf{w}_{<t} \circ w \circ \mathcal{W}^*)$ —divided by the probability of observing prefix  $\mathbf{w}_{<t}$ —represented by  $\mathbb{P}_{\mathcal{W}}(\mathbf{w}_{<t} \circ \mathcal{W}^*)$ . We call probabilities of the form  $\mathbb{P}_{\mathcal{W}}(\mathbf{w} \circ \mathcal{W}^*)$  the **prefix probability** of  $\mathbf{w}$ . As we will show, careful consideration of these prefix probabilities is critical for converting between our desired distributions (over words) and the ones provided by language models (over subwords).

**Orthography.** We assume here this language can be written, and that it has a standardised orthographic convention. Formally, given a language’s alphabet  $\mathcal{C}$ , each string  $\mathbf{w}$  can be mapped to a sequence of characters  $\mathbf{c} \in \mathcal{C}^*$  via function  $\mathbb{S}_{\mathcal{W}^* \rightarrow \mathcal{C}^*} : \mathcal{W}^* \rightarrow \mathcal{C}^*$ . Further, we assume this language allows for straightforward segmentation from orthography. Given a sequence of characters  $\mathbf{c}$ , we can thus extract a sequence of words as  $\mathbb{S}_{\mathcal{C}^* \rightarrow \mathcal{W}^*}(\mathbf{c}) = \mathbf{w}$ .

## 4 Subwords and Language Models

Most modern language models are not defined directly as distributions over words  $\mathbf{w}$ , but rather as distributions over *subwords*. These subwords are themselves defined by a choice of **tokenizer**.<sup>6</sup> In this section, we first introduce tokenisers, and how they map words to subwords (and back). We then use these building blocks to show how we can compute word probabilities from subword probabilities.

### 4.1 From Words to Subwords and Back

We define a tokenizer here as a tuple  $\langle \mathcal{S}, \mathbb{S}_{\mathcal{S}^* \rightarrow \mathcal{C}^*}, \mathbb{S}_{\mathcal{C}^* \rightarrow \mathcal{S}^*} \rangle$ . This tuple consists of: (i) a **vocabulary**  $\mathcal{S}$ , whose elements are subwords  $s \in \mathcal{S}$ , each of which represents a sequence of characters  $\mathbf{c} \in \mathcal{C}^*$ ,<sup>7</sup> (ii) a **detokenisation function**  $\mathbb{S}_{\mathcal{S}^* \rightarrow \mathcal{C}^*} : \mathcal{S}^* \rightarrow \mathcal{C}^*$ , which is simply a function that maps a sequence of subwords to the characters they represent and concatenates them together; (iii) a **tokenisation function**  $\mathbb{S}_{\mathcal{C}^* \rightarrow \mathcal{S}^*} : \mathcal{C}^* \rightarrow \mathcal{S}^*$ , which takes as input a character sequence and maps it to a subword sequence. Notably, multiple subword sequences may map to the same character sequence. However, most tokenisers specify one of

<sup>6</sup>We are not concerned with most aspects of individual tokenisers, and will focus on general considerations here. See Mielke et al. (2021) for a more comprehensive discussion.

<sup>7</sup>While subwords can be mapped back to a set of characters, they need not consist of only characters from the alphabet  $\mathcal{C}$ . Additional markers—such as bow—can be used.

these subword sequences as the canonical mapping and employ a deterministic tokenisation function.

Collectively, the mapping functions we have defined give us the ability to convert between words and subwords, which will be necessary when using subword distributions to compute word probabilities. We write word-to-subword mappings as:

$$\mathbb{S}_{\mathcal{W}^* \rightarrow \mathcal{S}^*} \stackrel{\text{def}}{=} \mathbb{S}_{\mathcal{W}^* \rightarrow \mathcal{C}^*} \bullet \mathbb{S}_{\mathcal{C}^* \rightarrow \mathcal{S}^*}, \quad \mathbb{S}_{\mathcal{S}^* \rightarrow \mathcal{W}^*} \stackrel{\text{def}}{=} \mathbb{S}_{\mathcal{S}^* \rightarrow \mathcal{C}^*} \bullet \mathbb{S}_{\mathcal{C}^* \rightarrow \mathcal{W}^*} \quad (4)$$

where  $\bullet$  represents function composition. Importantly, these functions reverse each other when applied as  $\mathbb{S}_{\mathcal{S}^* \rightarrow \mathcal{W}^*}(\mathbb{S}_{\mathcal{W}^* \rightarrow \mathcal{S}^*}(\mathbf{w})) = \mathbf{w}$ , but not necessarily when applied in the opposite order. The implication of this is that each  $\mathbf{w}$  maps to a unique  $\mathbf{s}$ , and every  $\mathbf{w}$  can be represented by some  $\mathbf{s}$ ; but there are subword sequences that will *not* be mapped to by our tokenisation function. For example, if a tokenizer maps word *probability* to subwords  $[_{prob}, ability]$ , then the subword sequence  $[_{p}, r, o, b, \dots]$  will never be mapped to. We denote **unmapped subword sequences** as:

$$\mathcal{S}_x \stackrel{\text{def}}{=} \mathcal{S}^* \setminus \left\{ \mathbb{S}_{\mathcal{W}^* \rightarrow \mathcal{S}^*}(\mathbf{w}) \mid \mathbf{w} \in \mathcal{W}^* \right\} \quad (5)$$

### 4.2 From Word to Subword Probabilities

Now let  $p_{\theta}$  be a language model with parameters  $\theta$  and a vocabulary  $\mathcal{S}$ . This model defines a probability distribution over the set of all finite subword sequences  $\mathbf{s} \in \mathcal{S}^*$  and its parameters are optimized to provide good estimates of the true distribution over subwords, given by:

$$p(\mathbf{s}) = \sum_{\mathbf{w} \in \mathcal{W}^*} p(\mathbf{w}) \mathbb{1} \left\{ \mathbf{s} = \mathbb{S}_{\mathcal{W}^* \rightarrow \mathcal{S}^*}(\mathbf{w}) \right\} \quad (6)$$

As not all subword sequences are mapped to, and because each mapping in  $\mathbb{S}_{\mathcal{W}^* \rightarrow \mathcal{S}^*}$  is unique, we can re-write this distribution as:

$$p(\mathbf{s}) = \begin{cases} p(\mathbf{w}) & \text{if } \mathbf{s} = \mathbb{S}_{\mathcal{W}^* \rightarrow \mathcal{S}^*}(\mathbf{w}) \\ 0 & \text{if } \mathbf{s} \in \mathcal{S}_x \end{cases} \quad (7)$$

### 4.3 From Subword to Word Probabilities

Eq. (7) suggests a way to extract probabilities over words from a language model; we can simply use the equivalence:<sup>8</sup>

$$p(\mathbf{w}) = p(\mathbf{s}), \quad \text{for } \mathbf{s} = \mathbb{S}_{\mathcal{W}^* \rightarrow \mathcal{S}^*}(\mathbf{w}) \quad (8)$$

<sup>8</sup>Notably, to apply this equivalence in practice, one needs an **exact language model**, which we define as a model  $p_{\theta}$  with the same support as  $p$ , i.e.,  $p_{\theta}(\mathbf{s}) = 0$  when  $p(\mathbf{s}) = 0$ . Note that most neural language models *cannot* assign zero probability to any subword sequence due to their use of a softmax projection in the final step of computing probabilities;

The implication of eq. (8) is that if we can create a subword set  $\Psi_s$  that is “equivalent” to a chosen word set  $\Psi_w$ , we would be able to compute  $\Psi_w$ ’s probability by summing over the subwords in  $\Psi_s$ . Formally, we define the set equivalence  $\triangleq$  between two sets of sequences as:

$$\Psi_w \triangleq \Psi_s \implies \left( w \in \Psi_w \iff \underset{w^* \rightarrow s^*}{\mathbb{S}}(w) \in \Psi_s \right) \quad (9)$$

Now let  $\mathbb{P}_s$  be a probability function defined analogously to  $\mathbb{P}_w$  (in Defn. 1). It then follows that:

$$\mathbb{P}_w(\Psi_w) = \mathbb{P}_s(\Psi_s), \quad \text{for } \Psi_w \triangleq \Psi_s \quad (10)$$

We are now in a position to define our quantity of interest  $p(w \mid \mathbf{w}_{<t})$  in terms of subword probabilities: it is simply the quotient of  $\mathbb{P}_s(\cdot)$  for two different sets  $\Psi_s$ .

**Lemma 1.** *The contextual probability of a word can be computed using probability distributions over subwords as:*

$$p(w \mid \mathbf{w}_{<t}) = \frac{\mathbb{P}_s(\Psi'_s)}{\mathbb{P}_s(\Psi''_s)} \quad (11)$$

where  $\Psi'_s \triangleq \mathbf{w}_{<t} \circ w \circ \mathcal{W}^*$  and  $\Psi''_s \triangleq \mathbf{w}_{<t} \circ \mathcal{W}^*$ .

*Proof.* This result follows from a simple application of the equivalence in eq. (10) to the definition of  $p(w \mid \mathbf{w}_{<t})$  in eq. (3).  $\square$

Luckily, it is straightforward to find the sets  $\Psi'_s$  and  $\Psi''_s$  required by Lemma 1. This is because, for a given word set  $\Psi_w$ , the subword set

$$\Psi_s = \left\{ \underset{w^* \rightarrow s^*}{\mathbb{S}}(w) \mid w \in \Psi_w \right\} \quad (12)$$

satisfies  $\Psi_w \triangleq \Psi_s$ : first, by construction, we have that  $w \in \Psi_w \implies \underset{w^* \rightarrow s^*}{\mathbb{S}}(w) \in \Psi_s$ ; second, due to the injectivity of  $\underset{w^* \rightarrow s^*}{\mathbb{S}}$ , it must be that  $\underset{w^* \rightarrow s^*}{\mathbb{S}}(w) \in \Psi_s \implies w \in \Psi_w$ . These sets thus meet the *iff* criteria required by our definition in eq. (9).

Before making use of eq. (11) for computing contextual probabilities, however, there is still one hurdle to overcome: the two sets  $\Psi'_w = (\mathbf{w}_{<t} \circ w \circ \mathcal{W}^*)$  and  $\Psi''_w = (\mathbf{w}_{<t} \circ \mathcal{W}^*)$  are infinite. We must thus find a more efficient strategy to compute these probabilities than summing over the (also infinite) sets  $\Psi'_s$  and  $\Psi''_s$ .

they will thus not be exact in this sense. While we focus on exact language models in this paper, we note that extending our results to inexact ones simply requires marginalising out potential ambiguities, i.e., computing  $p(w)$  for a given word requires summing over the (finite) set of subword sequences which map to it (Cao and Rimell, 2021).

#### 4.4 Leveraging LMs’ Autoregressiveness

We now discuss how we can leverage the fact that most LMs compute probabilities autoregressively to efficiently compute the probabilities in Lemma 1. In short, most LMs provide estimates of conditional probabilities:  $p(s \mid \mathbf{s}_{<t})$ . Given eq. (3) and the fact that  $\mathbb{P}_s(\mathcal{S}^*) = 1$ , we can use these conditionals to compute prefix probabilities efficiently.

**Lemma 2.** *Prefix probabilities can be computed using conditional probabilities as:*

$$\mathbb{P}_s(\mathbf{s} \circ \mathcal{S}^*) = \prod_{t=1}^{|\mathbf{s}|} \frac{\mathbb{P}_s(\mathbf{s}_{<t} \circ s_t \circ \mathcal{S}^*)}{\mathbb{P}_s(\mathbf{s}_{<t} \circ \mathcal{S}^*)} = \prod_{t=1}^{|\mathbf{s}|} p(s_t \mid \mathbf{s}_{<t}) \quad (13)$$

It follows that if we can find a set of subword sequences  $\Psi_s = \{\mathbf{s}^{(k)}\}_{k=1}^K$  for which we have the equivalence  $\mathbf{w} \circ \mathcal{W}^* \triangleq \bigcup_{\mathbf{s} \in \Psi_s} \mathbf{s} \circ \mathcal{S}^*$ , then we can compute prefix probabilities as:<sup>9</sup>

$$\mathbb{P}_s \left( \bigcup_{\mathbf{s} \in \Psi_s} \mathbf{s} \circ \mathcal{S}^* \right) = \sum_{\mathbf{s} \in \Psi_s} \mathbb{P}_s(\mathbf{s} \circ \mathcal{S}^*) \quad (14)$$

In turn, these let us compute  $p(w \mid \mathbf{w}_{<t})$  efficiently through eq. (11). For most tokenisers, finding a set  $\Psi_s$  for which the equivalence  $\mathbf{w} \circ \mathcal{W}^* \triangleq \bigcup_{\mathbf{s} \in \Psi_s} \mathbf{s} \circ \mathcal{S}^*$  holds is not actually possible due to the existence of unmapped sequences in  $\mathbf{s} \circ \mathcal{S}^*$ ; unmapped sequences, however, have zero probability and including them in  $\Psi'_s$  or  $\Psi''_s$  does not affect the equality in eq. (11). We thus ignore this issue in our exposition, while still considering it in our theorem proofs. We now outline tokeniser-specific considerations which influence how to choose these sets.

## 5 The Nuances of Mapping: Tokeniser-dependent Strategies

We are left with the task of finding a set of subword prefixes which will allow us to compute the probabilities of  $\Psi'_s \triangleq \Psi'_w$  and  $\Psi''_s \triangleq \Psi''_w$ . In this section, we discuss how our tokeniser—specifically whether it uses end- or beginning-of-word markings in its vocabulary—affects this task.

<sup>9</sup>In practice, we also need these prefix sets to be disjoint:  $(\mathbf{s} \circ \mathcal{S}^*) \cap (\mathbf{s}' \circ \mathcal{S}^*) = \emptyset$  for  $\mathbf{s}, \mathbf{s}' \in \Psi_s$ . This will be the case whenever no  $\mathbf{s} \in \Psi_s$  is a prefix of another  $\mathbf{s}' \in \Psi_s$  (i.e.,  $\mathbf{s}' \notin \mathbf{s} \circ \mathcal{S}^*$ ). If there is an  $\mathbf{s}$  which is a prefix of  $\mathbf{s}'$ , however, we can easily find a new set which still satisfies the equivalence above by dropping  $\mathbf{s}'$  from  $\Psi_s$ .

$$\begin{aligned}
& \mathbb{S}_{\mathcal{W}^* \rightarrow \mathcal{S}^*}(\text{How} \circ \text{do} \circ \text{you} \circ \text{compute} \circ \text{a} \circ \text{word} \circ \text{'s} \circ \text{probability} \circ ?) && \text{eow-marked, split punctuation} \\
& = \mathbb{S}_{\mathcal{W}^* \rightarrow \mathcal{S}^*}(\text{How}) \circ \mathbb{S}_{\mathcal{W}^* \rightarrow \mathcal{S}^*}(\text{do}) \circ \mathbb{S}_{\mathcal{W}^* \rightarrow \mathcal{S}^*}(\text{you}) \circ \mathbb{S}_{\mathcal{W}^* \rightarrow \mathcal{S}^*}(\text{compute}) \circ \mathbb{S}_{\mathcal{W}^* \rightarrow \mathcal{S}^*}(\text{a}) \circ \mathbb{S}_{\mathcal{W}^* \rightarrow \mathcal{S}^*}(\text{word}) \circ \mathbb{S}_{\mathcal{W}^* \rightarrow \mathcal{S}^*}(\text{'s}) \circ \mathbb{S}_{\mathcal{W}^* \rightarrow \mathcal{S}^*}(\text{probability}) \circ \mathbb{S}_{\mathcal{W}^* \rightarrow \mathcal{S}^*}(?) \\
& = \text{How}_\_ \circ \text{do}_\_ \circ \text{you}_\_ \circ \text{comp} \circ \text{ute}_\_ \circ \text{a}_\_ \circ \text{word} \circ \text{'s}_\_ \circ \text{prob} \circ \text{ability} \circ ?
\end{aligned}$$

$$\begin{aligned}
& \mathbb{S}_{\mathcal{W}^* \rightarrow \mathcal{S}^*}(\text{How} \circ \text{do} \circ \text{you} \circ \text{compute} \circ \text{a} \circ \text{word} \circ \text{'s} \circ \text{probability} \circ ?) && \text{bow-marked, split punctuation} \\
& = \mathbb{S}_{\mathcal{W}^* \rightarrow \mathcal{S}^*}(\text{How}) \circ \mathbb{S}_{\mathcal{W}^* \rightarrow \mathcal{S}^*}(\text{do}) \circ \mathbb{S}_{\mathcal{W}^* \rightarrow \mathcal{S}^*}(\text{you}) \circ \mathbb{S}_{\mathcal{W}^* \rightarrow \mathcal{S}^*}(\text{compute}) \circ \mathbb{S}_{\mathcal{W}^* \rightarrow \mathcal{S}^*}(\text{a}) \circ \mathbb{S}_{\mathcal{W}^* \rightarrow \mathcal{S}^*}(\text{word}) \circ \mathbb{S}_{\mathcal{W}^* \rightarrow \mathcal{S}^*}(\text{'s}) \circ \mathbb{S}_{\mathcal{W}^* \rightarrow \mathcal{S}^*}(\text{probability}) \circ \mathbb{S}_{\mathcal{W}^* \rightarrow \mathcal{S}^*}(?) \\
& = \text{How} \circ \_\text{do} \circ \_\text{you} \circ \_\text{comp} \circ \text{ute} \circ \_\text{a} \circ \_\text{word} \circ \text{'s} \circ \_\text{prob} \circ \text{ability} \circ ?
\end{aligned}$$

Figure 2: The output of tokenisers with different methods of handling word delineations.

## 5.1 Segmentation-compatible Tokenisers

In the following sections, we consider  $\mathbb{S}_{\mathcal{W}^* \rightarrow \mathcal{S}^*}$  that operate independently over words in a sequence  $\mathbf{w}$ . This is necessary for our methods below, and is a common practice in NLP (typically called pre-tokenisation) where a text is segmented according to some criterion (e.g., white space) before being converted into subwords by a tokeniser. Here, we consider pre-tokenisation to be one of the steps implemented by  $\mathbb{S}_{\mathcal{W}^* \rightarrow \mathcal{S}^*}$ . We formalise this in the following definition.

**Definition 2.** We define a *segmentation-compatible tokeniser* as one whose operations can be decomposed across words in a sequence, i.e.:

$$\begin{aligned}
\mathbb{S}_{\mathcal{W}^* \rightarrow \mathcal{S}^*}(\mathbf{w}) &= \mathbb{S}_{\mathcal{W}^* \rightarrow \mathcal{S}^*}(\mathbf{w}_{<t}) \circ \mathbb{S}_{\mathcal{W}^* \rightarrow \mathcal{S}^*}(w_t) \circ \mathbb{S}_{\mathcal{W}^* \rightarrow \mathcal{S}^*}(\mathbf{w}_{>t}) \quad (15) \\
&= \mathbb{S}_{\mathcal{W}^* \rightarrow \mathcal{S}^*}(w_1) \circ \mathbb{S}_{\mathcal{W}^* \rightarrow \mathcal{S}^*}(w_2) \circ \cdots \circ \mathbb{S}_{\mathcal{W}^* \rightarrow \mathcal{S}^*}(w_{|\mathbf{w}|})
\end{aligned}$$

While it is possible to create tokenisers with vocabularies in which subwords can cross word boundaries, the majority of them meet Defn. 2.<sup>10</sup>

The decomposition in Defn. 2 has an important implication. As discussed in §4.1, the (sequence-level) tokenisation function  $\mathbb{S}_{\mathcal{W}^* \rightarrow \mathcal{S}^*}$  must be injective, meaning that each word sequence must map to a unique subword sequence; this, in turn, implies that concatenating the outputs of  $\mathbb{S}_{\mathcal{W}^* \rightarrow \mathcal{S}^*}$  should always result in unique subword sequences. This property is known in the compression literature as unique decodability (Cover and Thomas, 2006, page 105). At an intuitive level, we can see why this is a desirable property of a tokenisation function: when working with NLP models, we want to be able to deterministically map a sequence of subwords to a sequence of words. A relatively simple strategy to ensure unique decodability, which is used by the

<sup>10</sup>E.g., the sentencepiece library (Kudo and Richardson, 2018) has an option which allows multi-word subwords to be added to a tokeniser’s vocabulary; by default, though, this option is disabled and it does not consider tokens of this format.

majority of tokenisers, is to mark either the ends or beginnings of words (eow or bow) using a subset of the subwords in  $\mathcal{S}$ . We discuss these strategies next.

## 5.2 End of Word Markings

We now consider eow-marking tokenisers. These tokenisers use a subset of their vocabulary  $\mathcal{S}_{\text{eow}} \subseteq \mathcal{S}$  to indicate the end of words,<sup>11</sup> with the rest of the vocabulary  $\mathcal{S}_{\text{mid}} \stackrel{\text{def}}{=} \mathcal{S} \setminus \mathcal{S}_{\text{eow}}$  mapping back to the beginning or middle of words.

**Definition 3.** An *eow-marking tokeniser* is a segmentation-compatible tokeniser which marks ends of words. Its word-level tokenisation function can be written as  $\mathbb{S}_{\mathcal{W}^* \rightarrow \mathcal{S}^*}^{\text{eow}} : \mathcal{W} \rightarrow \mathcal{S}_{\text{mid}}^* \circ \mathcal{S}_{\text{eow}}$ .<sup>12</sup>

Importantly, given the definition above, when a subword  $s_t \in \mathcal{S}_{\text{eow}}$  is observed, it means that the current subsequence  $\mathbf{s}_{t':t}$  (where  $t' \leq t$ ) can be mapped back to a word, and that a subsequence representing a new word will begin at  $s_{t+1}$ . (The current subsequence  $\mathbf{s}_{t':t}$  is thus determined by the smallest  $t'$  for which  $\mathbf{s}_{t':t-1} \in \mathcal{S}_{\text{mid}}^*$ ; note that this means either  $t' = 1$  or  $s_{t'-1} \in \mathcal{S}_{\text{eow}}$ .) This property implies that eow-marking tokenisers provide **instantaneous decodability** (Cover and Thomas, 2006, page 106): prefix  $\mathbf{s}_{\leq t}$  with  $s_t \in \mathcal{S}_{\text{eow}}$  is instantaneously decodable, as it always maps to the same words, regardless of its continuation  $\mathbf{s}_{>t}$ . Instantaneous decodability allows us to compute the contextual probability of a word as follows.

**Theorem 1.** Let  $\mathbb{S}_{\mathcal{W}^* \rightarrow \mathcal{S}^*}^{\text{eow}}$  be a eow-marking tokeniser. Further, let  $\mathbf{s}^{\mathbf{w}} \stackrel{\text{def}}{=} \mathbb{S}_{\mathcal{W}^* \rightarrow \mathcal{S}^*}^{\text{eow}}(\mathbf{w})$ . We can show the following equivalence:

$$\begin{aligned}
\mathbb{P}_{\mathcal{W}}(\mathbf{w}_{<t} \circ \mathcal{W}^*) &= \mathbb{P}_{\mathcal{S}}(\mathbf{s}^{\mathbf{w}_{<t}} \circ \mathcal{S}^*) \quad (16) \\
\mathbb{P}_{\mathcal{W}}(\mathbf{w}_{<t} \circ w \circ \mathcal{W}^*) &= \mathbb{P}_{\mathcal{S}}(\mathbf{s}^{\mathbf{w}_{<t}} \circ \mathbf{s}^w \circ \mathcal{S}^*)
\end{aligned}$$

<sup>11</sup>The case of  $\mathcal{S}_{\text{eow}} = \mathcal{S}$  or  $\mathcal{S}_{\text{bow}} = \mathcal{S}$  happens when  $\mathcal{S} = \mathcal{W}$ ; while possible in theory, it will not happen in practice since a language model cannot have an infinite vocabulary.

<sup>12</sup>Note that only subword sequences ending in  $s \in \mathcal{S}_{\text{eow}}$  or the empty sequence (i.e., “”) are valid under this tokeniser. This is because:  $\bigcup_{i=0}^{\infty} (\mathcal{S}_{\text{mid}}^* \circ \mathcal{S}_{\text{eow}})^i \subseteq \{\text{""}\} \cup (\mathcal{S}^* \circ \mathcal{S}_{\text{eow}})$ .

Further, we can compute a word’s probability as:

$$p(w \mid \mathbf{w}_{<t}) = \underbrace{\prod_{t'=1}^{|\mathbf{s}^w|} p(\mathbf{s}_{t'}^w \mid \mathbf{s}^{\mathbf{w}_{<t}} \circ \mathbf{s}_{<t'}^w)}_{p(\mathbf{s}^w \mid \mathbf{s}^{\mathbf{w}_{<t}})} \quad (17)$$

*Proof.* See App. D.1 for formal proof.  $\square$

Eq. (16) follows from instantaneous decodability, as every sequence  $\mathbf{s} \in \mathbf{s}^w \circ \mathcal{S}^*$  maps back to  $\mathbf{w} \circ \mathcal{W}^*$ . Eq. (17) then follows from a simple application of Lemmas 1 and 2:

$$p(\mathbf{s}^w \mid \mathbf{s}^{\mathbf{w}_{<t}}) = \frac{\prod_{t'=1}^{|\mathbf{s}^{\mathbf{w}_{<t} \circ \mathbf{w}}|} p(\mathbf{s}_{t'}^{\mathbf{w}_{<t} \circ \mathbf{w}} \mid \mathbf{s}_{<t'}^{\mathbf{w}_{<t} \circ \mathbf{w}})}{\prod_{t'=1}^{|\mathbf{s}^{\mathbf{w}_{<t}}|} p(\mathbf{s}_{t'}^{\mathbf{w}_{<t}} \mid \mathbf{s}_{<t'}^{\mathbf{w}_{<t}})} \quad (18)$$

Notably, eq. (17) is fairly straightforward and is how most NLP practitioners would compute a word’s probability. In the next section, however, we see that it would not compute the correct probabilities if using bow-marking tokenisers.

### 5.3 Beginning of Word Markings

We now consider bow-marking tokenisers. Analogously to the eow case, a subset of a bow-marking tokeniser’s vocabulary  $\mathcal{S}_{\text{bow}} \subseteq \mathcal{S}$  is used exclusively to indicate word beginnings. The rest of the vocabulary  $\mathcal{S}_{\text{mid}} \stackrel{\text{def}}{=} \mathcal{S} \setminus \mathcal{S}_{\text{bow}}$  then represents either the middle or end of words. We provide a formal definition of this tokeniser below.

**Definition 4.** A *bow-marking tokeniser* is a segmentation-compatible tokeniser which marks beginnings of words. Its word-level tokenisation function is written as  $\mathbb{S}_{\mathcal{W} \rightarrow \mathcal{S}^*}^{\text{bow}} : \mathcal{W} \rightarrow \mathcal{S}_{\text{bow}} \circ \mathcal{S}_{\text{mid}}^*$ .<sup>13</sup>

Given the definition above, when a subword  $s_t \in \mathcal{S}_{\text{bow}}$  is observed, it thus means that a *previous* subsequence  $\mathbf{s}_{t':t-1}$  can be mapped back to a word, and that a subsequence representing a new word begins at  $s_t$ . (The previous subsequence  $\mathbf{s}_{t':t-1}$  is determined by  $s_{t'} \in \mathcal{S}_{\text{bow}}$  and  $\mathbf{s}_{t'+1:t-1} \in \mathcal{S}_{\text{mid}}^*$ .) Such tokenisers are thus *not* instantaneously decodable. They only provide what we term **near-instantaneous decodability**: a prefix  $\mathbf{s}_{\leq t}$  does *not* always map to the same words, as its mapping depends on whether the following subword  $s_{t+1}$  is in  $\mathcal{S}_{\text{bow}} \cup \{\text{eos}\}$ .<sup>14</sup> Computing probabilities with near-instantaneous codes thus requires discounting the

<sup>13</sup>Similarly to with eow, not all subword sequences are valid under bow tokenisers, only sequences in  $\{\text{eos}\} \cup (\mathcal{S}_{\text{bow}} \circ \mathcal{S}^*)$ .

<sup>14</sup>Here, we define the concatenation of any sequence with eos to be itself, e.g.,  $\mathbf{s} \circ \text{eos} = \mathbf{s}$ .

probability of continuations  $s_{t+1} \notin \mathcal{S}_{\text{bow}} \cup \{\text{eos}\}$ ; we label this discount factor as **Bug Fix ①**.

**Theorem 2.** Let  $\mathbb{S}_{\mathcal{W} \rightarrow \mathcal{S}^*}$  be a bow-marking tokeniser. Further, let  $\bar{\cdot}$  represent the union of a set with eos, e.g.,  $\bar{\mathcal{S}}_{\text{bow}} = \mathcal{S}_{\text{bow}} \cup \{\text{eos}\}$ . We can show the following equivalence:

$$\begin{aligned} \mathbb{P}_{\mathcal{W}}(\mathbf{w}_{<t} \circ \mathcal{W}^*) &= \mathbb{P}_{\mathcal{S}}(\mathbf{s}^{\mathbf{w}_{<t}} \circ \overline{\mathcal{S}_{\text{bow}} \circ \mathcal{S}^*}) \quad (19) \\ \mathbb{P}_{\mathcal{W}}(\mathbf{w}_{<t} \circ w \circ \mathcal{W}^*) &= \mathbb{P}_{\mathcal{S}}(\mathbf{s}^{\mathbf{w}_{<t}} \circ \mathbf{s}^w \circ \overline{\mathcal{S}_{\text{bow}} \circ \mathcal{S}^*}) \end{aligned}$$

Further, we can compute a word’s probability as:

$$p(w \mid \mathbf{w}_{<t}) = \underbrace{\prod_{t'=1}^{|\mathbf{s}^w|} p(\mathbf{s}_{t'}^w \mid \mathbf{s}^{\mathbf{w}_{<t}} \circ \mathbf{s}_{<t'}^w)}_{p(\mathbf{s}^w \mid \mathbf{s}^{\mathbf{w}_{<t}})} \underbrace{\frac{\sum_{\{s \in \bar{\mathcal{S}}_{\text{bow}}\}} p(s \mid \mathbf{s}^{\mathbf{w}_{<t}} \circ \mathbf{s}^w)}{\sum_{\{s \in \bar{\mathcal{S}}_{\text{bow}}\}} p(s \mid \mathbf{s}^{\mathbf{w}_{<t}})}}_{\text{Bug Fix ①}} \quad (20)$$

*Proof.* See App. D.2 for formal proof.  $\square$

Eq. (19) follows from near-instantaneous decodability, as every sequence  $\mathbf{s}^w \circ \overline{\mathcal{S}_{\text{bow}} \circ \mathcal{S}^*}$  maps back to  $\mathbf{w} \circ \mathcal{W}^*$ , but sequences in  $\mathbf{s}^w \circ \mathcal{S}_{\text{mid}} \circ \mathcal{S}^*$  do not. Fig. 2 contains an example of a sequence tokenised using either eow- or bow-marking tokenisers; Fig. 3 contains an example motivating **Bug Fix ①**.

### 5.4 Practical Concerns and Corner Cases

In this section, we discuss corner cases that deserve special consideration. Many of these cases arise because of practical demands, e.g., ensuring the presence or absence of white space where appropriate. Notably, the need for these corner cases is often language-dependent, as they arise due to orthographic conventions. We discuss the implications of two tokeniser conventions that handle special cases: the treatment of the beginnings and ends of sequences.

**Non-eow-marked Final Words.** Several eow-marking tokenisers do not decompose exactly as in eq. (15), but treat the final word in a sequence differently. Specifically, they override the behaviour of  $\mathbb{S}_{\mathcal{W} \rightarrow \mathcal{S}^*}$  on these words and do *not* use subwords from  $\mathcal{S}_{\text{eow}}$  to mark its ends. This is also often the treatment applied to words followed immediately by punctuation. This mechanism allows tokenisers to avoid implying the existence of a white space that does not exist, e.g., at the end of a string. Notably, this breaks instantaneous decodability, making this code only near-instantaneous. A simple example demonstrates this fact: let  $\mathbb{S}_{\text{mid}}^w \stackrel{\text{def}}{=} \mathbb{S}_{\mathcal{W} \rightarrow \mathcal{S}^*}^{\text{mid}}(w)$ , where  $\mathbb{S}_{\mathcal{W} \rightarrow \mathcal{S}^*}^{\text{mid}} : \mathcal{W} \rightarrow \mathcal{S}_{\text{mid}}^*$ .

### Worked Example: Contextual probability computation for word `mark` using a bow-marking tokeniser

Let “`She saw the mark...`” be our context of interest; we thus have  $\mathbf{w} = \langle \text{She, saw, the, mark, ...} \rangle$ . Further, let  $p_\theta$  be our language model with vocabulary:

$$\mathcal{S} = \{ \_a, \_an, \_mark, \_saw, \_She, \_the, er, tion, ing, ed \}$$

Let’s assume that we are interested in estimating  $p(\text{mark} \mid \langle \text{She, saw, the} \rangle)$  using  $p_\theta$ . To employ eq. (11), we must compute  $\mathbb{P}_{\mathcal{S}}(\Psi_{\mathcal{S}})$  for  $\Psi'_{\mathcal{S}} \triangleq \langle \text{She, saw, the} \rangle \circ \text{mark} \circ \mathcal{W}^*$  and  $\Psi''_{\mathcal{S}} \triangleq \langle \text{She, saw, the} \rangle \circ \mathcal{W}^*$ . For vocabularies derived using bow-marking tokenisers, Theorem 2 states that we should use:

$$\begin{aligned} \Psi'_{\mathcal{S}} &= \_She \circ \_saw \circ \_the \circ \_mark \circ \overline{\mathcal{S}_{\text{bow}} \circ \mathcal{S}^*} \\ \Psi''_{\mathcal{S}} &= \_She \circ \_saw \circ \_the \circ \overline{\mathcal{S}_{\text{bow}} \circ \mathcal{S}^*} \end{aligned}$$

where  $\mathcal{S}_{\text{bow}} = \{ \_a, \_an, \_mark, \_saw, \_She, \_the \}$ . Using this theorem’s eq. (20) we arrive at:

$$p(\text{mark} \mid \langle \text{She, saw, the} \rangle) = p_\theta(\_mark \mid \langle \_She, \_saw, \_the \rangle) \cdot \frac{\sum_{\{s \in \overline{\mathcal{S}_{\text{bow}}}\}} p_\theta(s \mid \langle \_She, \_saw, \_the \rangle \circ \_mark)}{\sum_{\{s \in \overline{\mathcal{S}_{\text{bow}}}\}} p_\theta(s \mid \langle \_She, \_saw, \_the \rangle)}$$

Note that this computation specifically discounts the probabilities  $p(\text{marker} \mid \langle \text{She, saw, the} \rangle)$ ,  $p(\text{marktion} \mid \langle \text{She, saw, the} \rangle)$ ,  $p(\text{markerer} \mid \langle \text{She, saw, the} \rangle)$ , etc., which otherwise would have incorrectly counted towards our estimate of  $p(\text{mark} \mid \langle \text{She, saw, the} \rangle)$ .

Figure 3: Example for computing a word’s probability using a LM over subwords defined by a bow-marked tokeniser.

Upon observing subsequence  $\mathbf{s}_{\text{mid}}^w$ , we cannot instantaneously map it back to  $w$ , and must wait for the next symbol: if  $\mathbf{s}_{\text{mid}}^w$  is followed by either eos or punctuation, then it is mapped back to  $w$ ; if not, it is mapped to another word. Handling this thus requires the following fix (termed **Bug Fix ②** here):

$$p(w \mid \mathbf{w}_{<t}) = \underbrace{\left( p(\mathbf{s}_{\text{mid}}^w \mid \mathbf{s}^{\mathbf{w}_{<t}}) \sum_{s \in \overline{\mathcal{S}}_{! ?}} p(s \mid \mathbf{s}^{\mathbf{w}_{<t}} \circ \mathbf{s}_{\text{mid}}^w) \right) + p(\mathbf{s}^w \mid \mathbf{s}^{\mathbf{w}_{<t}})}_{\text{Bug Fix ②}} \quad (21)$$

**Non-bow-marked First Words.** Just as eow-marking tokenisers often treat final words differently, bow-marking tokenisers treat the first word in a sequence differently to handle white space appropriately. These tokenisers typically do *not* mark first words with bow, and instead apply  $\mathbb{S}_{\mathcal{W} \rightarrow \mathcal{S}^*}^{\text{mid}}$  to  $w_1$ . This affects the probability computation of the first word in a sequence. In such cases, the prefix  $\mathbf{w}_{<t}$  of the first word is empty (denoted here as “”). While computing a word’s contextual probability according to eq. (19) requires computing  $\mathbb{P}_{\mathcal{S}}(\overline{\mathcal{S}_{\text{bow}} \circ \mathcal{S}^*})$ , the first subword in a sequence will not be in  $\mathcal{S}_{\text{bow}}$ , but in  $\mathcal{S}_{\text{mid}}$  instead. The probability computation of such words thus requires the following correction (**Bug Fix ③**):

$$p(w \mid \text{“”}) = p(\mathbf{s}_{\text{mid}}^w \mid \text{“”}) \underbrace{\frac{\sum_{\{s \in \overline{\mathcal{S}_{\text{bow}}}\}} p(s \mid \mathbf{s}^w)}{\sum_{\{s \in \overline{\mathcal{S}_{\text{mid}}}\}} p(s \mid \text{“”})}}_{\text{Bug Fix ③}} \quad (22)$$

### 5.5 Defining $\mathcal{S}_{\text{bow}}$ and $\mathcal{S}_{\text{eow}}$

Defining sets  $\mathcal{S}_{\text{bow}}$  or  $\mathcal{S}_{\text{eow}}$  for a given tokeniser is not necessarily straightforward, as tokenisers do not explicitly mark bow or eow in their vocabularies.<sup>15</sup> Further, these sets’ definitions will depend on what a researcher considers a word to be. As an example, we use the sentence in Fig. 2: *How do you compute a word’s probability?*. One could define words to be the set of whitespace-separated character sequences: `How` `do` `you` `compute` `a` `word’s` `probability?`. However, one may also consider punctuation and clitics to impose word boundaries, meaning words would instead be delineated as: `How` `do` `you` `compute` `a` `word` `’s` `probability` `?`. In the former case, we would define  $\mathcal{S}_{\text{bow}}$  and  $\mathcal{S}_{\text{eow}}$  simply as the set of subwords with a leading or trailing white space (e.g., `_word`  $\in \mathcal{S}_{\text{bow}}$  or `’s_`  $\in \mathcal{S}_{\text{eow}}$ ). In the latter case though, subwords starting with punctuation or clitics should also be included in  $\mathcal{S}_{\text{bow}}$  (e.g., we require `’s`  $\in \mathcal{S}_{\text{bow}}$ ). This choice further impacts probability computations: computing eow-marking probabilities in the former case simply requires eq. (17), while in the latter case it requires **Bug Fix ②**.<sup>16</sup>

<sup>15</sup>They often mark white spaces instead (denoted here as `_`), but white space need not be the only word-boundary marker.

<sup>16</sup>A recent work (Giulianelli et al., 2024) proposes a method allowing the computation of the probability of any character span within a sequence. E.g., one can compute the probability of `word` or `ord’s_prob` in the example above. While computing the probability of arbitrary character spans can be valuable, we note that there is no single sequence of characters that is equivalent to a word. For example,  $p(\text{compute\_} \mid \mathbf{c}_{<t})$  is the probability of `compute` followed by `_`; the methods here, however,

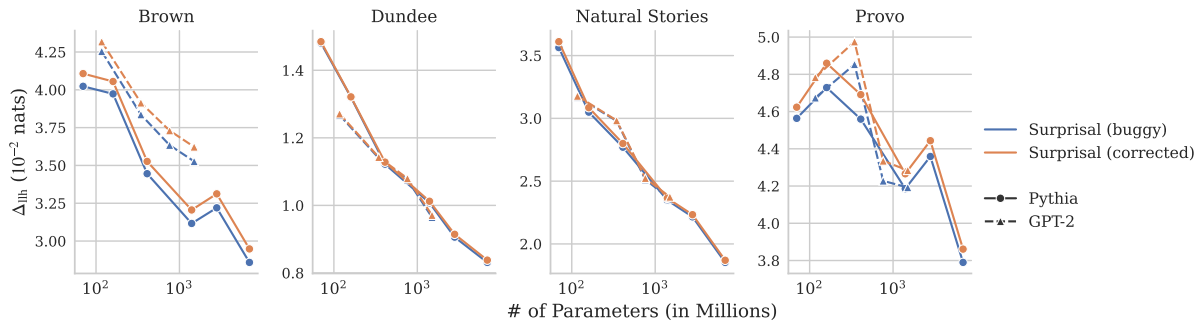


Figure 4:  $\Delta_{llh}$  between regressors with and without surprisal as a predictor. We include  $\Delta_{llh}$  when using surprisal estimates computed from language models across several sizes and families. Results are presented both when using the buggy and correct methods for surprisal estimation.

## 6 Experiments

We now investigate how correcting the computation of word probability estimates affects the results of prior studies. These works incorrectly computed probabilities as  $p(\mathbf{s}^w | \mathbf{s}^{w < t})$  (i.e., using eq. (19) without **Bug Fix** ①), which we term as **buggy** estimates here. We explore two settings: psycholinguistics experiments surrounding sentence comprehension (Hale, 2001; Levy, 2008) and computational linguistics experiments assessing the lexicon’s communicative efficiency (Piantadosi et al., 2011; Gibson et al., 2019). We follow these works’ experimental methodologies, observing how the use of corrected surprisal estimates impacts the conclusions that were originally drawn.

**Models.** In our first experiment, we estimate contextual probabilities using GPT-2 (Radford et al., 2019) and Pythia (Biderman et al., 2023); in the second, we focus only on Pythia. Both these suites contain language models of various sizes. We use these models’ open-source versions from the transformers library (Wolf et al., 2020). GPT-2 and Pythia use bow-marking tokenisers, meaning we employ the methods discussed in §5.3 to compute words’ contextual probabilities.

### 6.1 Sentence Comprehension

Surprisal theory (Hale, 2001; Levy, 2008) hypothesises that readers keep a belief distribution over meanings while reading; after observing each word in a sentence, they must thus update this distribution. Under some assumptions about how these

compute the probability of a word,  $p(\text{compute} | \mathbf{w}_{<t})$ , regardless of what follows it. We can combine our considerations and their method to recover the probability of a word by first defining a set of word-ending characters  $\mathcal{C}_{\text{bow}}$ , and then using it to marginalise a word’s probability over possible word-ending continuations:  $p(\text{compute} | \mathbf{w}_{<t}) = \sum_{c \in \bar{\mathcal{C}}_{\text{bow}}} p(\text{compute} \circ c | c_{<t})$ . We thus see our works as complementary.

belief updates are performed, surprisal theory then predicts that their cost is related to a word’s **surprisal**, defined as the negative log-probability:

$$h(w_t) \stackrel{\text{def}}{=} -\log p(w_t | \mathbf{w}_{<t}) \quad (23)$$

Surprisal theory is widely accepted as a model of comprehension effort, with numerous works empirically supporting it (Smith and Levy, 2008, 2013; Goodkind and Bicknell, 2018; Shain, 2019; Wilcox et al., 2020, 2023a; Oh et al., 2022; Shain et al., 2024, *inter alia*). Notably, the true contextual probabilities  $p(w_t | \mathbf{w}_{<t})$  required to compute surprisal are unknown, and must be approximated. All of the works above use language models to do so, with the most recent using LMs which operate on top of subwords produced by bow-marking tokenisers (e.g., Oh and Schuler, 2023b,a; Shain et al., 2024; Pimentel et al., 2023b). Notably, these works compute surprisal estimates using the aforementioned buggy estimates of  $p(w_t | \mathbf{w}_{<t})$ . In this section, we reproduce some of these prior works’ results, observing how this correction affects results.

**Setup Summary.** We run our analyses on 4 reading times datasets—Brown, Dundee, Natural Stories, and Provo. Further, following prior work (Wilcox et al., 2023a; Oh and Schuler, 2023b), we evaluate surprisal’s predictive power over reading times by measuring the change in data log-likelihood  $\Delta_{llh}$  when using linear regressors with and without surprisal as a predictor. More details about our experimental setup are in App. A.1.

**Results.** Fig. 4 shows the change in data log-likelihood under regressors with and without surprisal as a predictor; values are detailed in Tab. 1 (in the appendix). We first note that the predictive power of surprisal decreases as language model size increases, as observed in prior work (Oh and



Schuler, 2023b; Shain et al., 2024). Here however, we are more interested in the effect of our corrections on these results—labelled as buggy vs. corrected surprisal. Interestingly, we observe only small changes in predictive power due to our correction; individually, these changes are only significant for a few models (see Tab. 1 for detailed results). However, when analysed in aggregate for all models, we see this positive improvement is consistent and significant in the four dataset ( $\alpha < 0.01$  in our permutation tests). We also confirm the same patterns in seven other languages in App. C.

## 6.2 Communicative Efficiency

Languages’ lexicons have been studied for decades in an effort to gain better insights about the forces that shape natural languages (Zipf, 1935; Howes, 1968; Bentz and Ferrer-i-Cancho, 2016; Levshina, 2022). One characteristic of particular interest has been word lengths and how a tendency for communicative efficiency has influenced them. There are several hypotheses about the exact way in which this tendency takes effect. Zipf (1935) argues that speakers have a tendency towards minimising utterance lengths, and therefore that word lengths should correlate with frequencies. Piantadosi et al. (2011) argues that speakers maximise information transfer, and thus word lengths should correlate with a word’s expected surprisal instead:

$$\mathbb{E}[h(w_t)] \stackrel{\text{def}}{=} \mathbb{E}_{\mathbf{w}_{<t}} [-\log p(w_t | \mathbf{w}_{<t}) | w_t] \quad (24)$$

We follow Pimentel et al. (2023a) in calling this the **channel capacity hypothesis** (CCH). Finally, Pimentel et al. (2023a) point out an issue with Piantadosi et al.’s solution, and argue that to maximise information transfer, lengths should correlate with the following value instead:<sup>17</sup>

$$\frac{\mathbb{E}[h^2(w_t)]}{\mathbb{E}[h(w_t)]} \stackrel{\text{def}}{=} \frac{\mathbb{E}_{\mathbf{w}_{<t}} [(-\log p(w_t | \mathbf{w}_{<t}))^2 | w_t]}{\mathbb{E}_{\mathbf{w}_{<t}} [-\log p(w_t | \mathbf{w}_{<t}) | w_t]} \quad (25)$$

**Setup Summary.** We run our analysis using a subset of the English portion of Wiki-40B (Guo et al., 2020). We compare the three values above (unigram frequency, and eqs. (24) and (25)); evaluating them based on their correlation with words’ lengths. Two of these values depend on a word’s contextual probability, and we thus also compare their fixed vs. buggy versions.

<sup>17</sup>See their paper for a derivation for this fix.

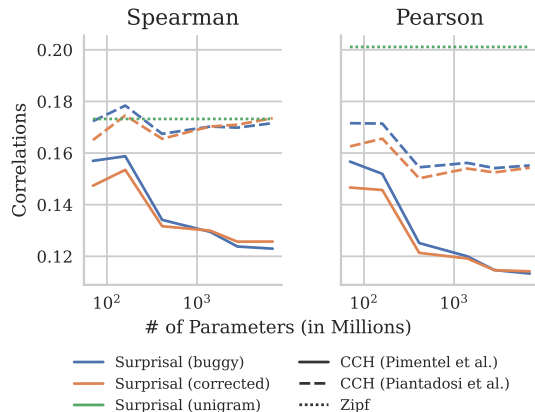


Figure 5: Correlation between English word lengths and the values predicted by either a Zipfian notion of optimality, or the channel capacity hypothesis. CCH (Pimentel et al.) and CCH (Piantadosi et al.) refer to eqs. (24) and (25).

**Results.** The results in Fig. 5 confirm the findings of Pimentel et al. (2023a): once larger (and better) language models are used to estimate words’ surprisals, the metrics under the CCH hypothesis (both Piantadosi et al.’s and Pimentel et al.’s versions) become weaker predictors of word lengths. Interestingly, correcting the computation of surprisals also leads to a drop in the correlations between CCH predictors and word lengths. Improving CCH’s predictors thus consistently hurts its predictive power over word lengths—either when using better models, Pimentel et al.’s fix to CCH’s optimal solution, or our fix to probability computations. We conclude, as Pimentel et al., that word lengths are best predicted by Zipf’s hypothesis.

## 7 Conclusion

This work expounds on the intricacies of accurately computing contextual word probabilities using language models. We focus on the challenges posed by the use of subword vocabularies. We show that subword vocabularies defined using beginning-of-word (bow) tokenisers—common in many modern LMs—introduce complexities that are often overlooked. We point out that this has led to potential inaccuracies in computing probability estimates of various prior empirical analyses. Our methodological corrections lead to significant differences in results, although the overarching conclusions of the previous studies that we explore remain the same. This finding underscores the importance of precise computational methods in linguistic research. Future work should ensure these corrections are adopted to enhance the reliability of their analyses.

## Limitations

The authors see limitations in both the theoretical and empirical aspects of this work. Perhaps the main theoretical limitation is the lack of consideration of all potential corner cases which tokenisers might implement (similar to, e.g., those discussed in §5.4). The use of white space differs from language to language, and many corner cases of tokeniser behaviour are designed specifically to handle this. There are likely other fixes to probability computations that would need to be derived to handle paradigms not discussed in §5.4. In Spanish, for instance, words following “¿” are usually not bow-marked, and might thus require the use of an approach similar to **Bug Fix** ③. Our theoretical results are also limited to autoregressive models. While the majority of today’s language models meet this criterion, it is feasible that future language models would be designed differently and consequently, our methods would no longer be necessarily applicable. On the empirical side, a large limitation of our work is the exploration of the impact of our methods in only two studies. Additional studies are thus needed to understand the full extent to which our corrections impact empirical results in other areas of computational linguistics (and of NLP, more broadly).

## Acknowledgments

We thank Ethan Wilcox for many discussions about this paper, and for helping to draft parts of it. We also thank Sotiris Anagnostidis and Pietro Lesci for feedback on earlier versions of this manuscript, and Yahya Emara and Mario Giulianelli for feedback on the final version.

## References

Christian Bentz and Ramon Ferrer-i-Cancho. 2016. [Zipf’s law of abbreviation as a language universal](#). In *Proceedings of the Leiden Workshop on Capturing Phylogenetic Algorithms for Linguistics*. Universität Tübingen.

Stella Biderman, Hailey Schoelkopf, Quentin Anthony, Herbie Bradley, Kyle O’Brien, Eric Hallahan, Mohammad Aflah Khan, Shivanshu Purohit, USVSN Sai Prashanth, Edward Raff, Aviya Skowron, Lintang Sutawika, and Oskar Van Der Wal. 2023. [Pythia: A suite for analyzing large language models across training and scaling](#). In *Proceedings of the 40th International Conference on Machine Learning*, ICML’23.

Kris Cao and Laura Rimell. 2021. [You should evaluate your language model on marginal likelihood over tokenisations](#). In *Proceedings of the 2021 Conference on Empirical Methods in Natural Language Processing*, pages 2104–2114, Online and Punta Cana, Dominican Republic. Association for Computational Linguistics.

Thomas M. Cover and Joy A. Thomas. 2006. *Elements of Information Theory*, second edition. Wiley-Interscience.

Robert M. W. Dixon and Alexandra Y. Aikhenvald. 2003. [Word: a typological framework](#). In *Word: A Cross-linguistic Typology*, page 1–41. Cambridge University Press.

Li Du, Lucas Torroba Hennigen, Tiago Pimentel, Clara Meister, Jason Eisner, and Ryan Cotterell. 2023. [A measure-theoretic characterization of tight language models](#). In *Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 9744–9770, Toronto, Canada. Association for Computational Linguistics.

Richard Futrell, Edward Gibson, Harry J. Tily, Idan Blank, Anastasia Vishnevetsky, Steven Piantadosi, and Evelina Fedorenko. 2018. [The natural stories corpus](#). In *Proceedings of the Eleventh International Conference on Language Resources and Evaluation*. European Language Resources Association.

Matthias Gallé. 2019. [Investigating the effectiveness of BPE: The power of shorter sequences](#). In *Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing (EMNLP-IJCNLP)*, pages 1375–1381, Hong Kong, China. Association for Computational Linguistics.

Edward Gibson, Richard Futrell, Steven T. Piantadosi, Isabelle Dautriche, Kyle Mahowald, Leon Bergen, and Roger Levy. 2019. [How efficiency shapes human language](#). *Trends in Cognitive Sciences*, 23(5):389–407.

Mario Giulianelli, Luca Malagutti, Juan Luis Gastaldi, Brian DuSell, Tim Vieira, and Ryan Cotterell. 2024. [On the proper treatment of tokenization in psycholinguistics](#). In *Proceedings of the 2024 Conference on Empirical Methods in Natural Language Processing*, Miami, Florida, USA. Association for Computational Linguistics.

Adam Goodkind and Klinton Bicknell. 2018. [Predictive power of word surprisal for reading times is a linear function of language model quality](#). In *Proceedings of the 8th Workshop on Cognitive Modeling and Computational Linguistics (CMCL 2018)*, pages 10–18, Salt Lake City, Utah. Association for Computational Linguistics.

Mandy Guo, Zihang Dai, Denny Vrandečić, and Rami Al-Rfou. 2020. [Wiki-40B: Multilingual language](#)

- model dataset. In *Proceedings of the 12th Language Resources and Evaluation Conference*, page 2440–2452, Marseille, France. European Language Resources Association.
- John Hale. 2001. A probabilistic Earley parser as a psycholinguistic model. In *Second Meeting of the North American Chapter of the Association for Computational Linguistics*, pages 1–8.
- Tracy Alan Hall and Ursula Kleinhenz. 1999. *Studies on the phonological word*. John Benjamins.
- Martin Haspelmath. 2023. Defining the word. *WORD*, 69(3):283–297.
- Davis Howes. 1968. Zipf’s law and Miller’s random-monkey model. *The American Journal of Psychology*, 81(2):269–272.
- Marcel Adam Just, Patricia A. Carpenter, and Jacqueline D. Woolley. 1982. Paradigms and processes in reading comprehension. *Journal of Experimental Psychology: General*, 111(2):228–238.
- Alan Kennedy, Robin Hill, and Joel Pynte. 2003. The Dundee corpus. In *Proceedings of the 12th European Conference on Eye Movements*.
- 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.
- Natalia Levshina. 2022. Frequency, informativity and word length: Insights from typologically diverse corpora. *Entropy*, 24(2).
- Roger Levy. 2008. Expectation-based syntactic comprehension. *Cognition*, 106(3):1126–1177.
- Roger Levy and T. Florian Jaeger. 2007. Speakers optimize information density through syntactic reduction. In *Advances in Neural Information Processing Systems*, volume 19. MIT Press.
- Steven G. Luke and Kiel Christianson. 2018. The Provo corpus: A large eye-tracking corpus with predictability norms. *Behavior Research Methods*, 50(2):826–833.
- Stephan C. Meylan and Thomas L. Griffiths. 2021. The challenges of large-scale, web-based language datasets: Word length and predictability revisited. *Cognitive Science*, 45(6):e12983.
- Sabrina J. Mielke, Zaid Alyafeai, Elizabeth Salesky, Colin Raffel, Manan Dey, Matthias Gallé, Arun Raja, Chenglei Si, Wilson Y. Lee, Benoît Sagot, and Samson Tan. 2021. Between words and characters: A brief history of open-vocabulary modeling and tokenization in NLP. *arXiv preprint arXiv:2112.10508*.
- Marina Nespor and Irene Vogel. 2007. *Prosodic Phonology: With a New Foreword*. De Gruyter Mouton, Berlin, Boston.
- Byung-Doh Oh, Christian Clark, and William Schuler. 2022. Comparison of structural parsers and neural language models as surprisal estimators. *Frontiers in Artificial Intelligence*, 5.
- Byung-Doh Oh and William Schuler. 2023a. Transformer-based language model surprisal predicts human reading times best with about two billion training tokens. In *Findings of the Association for Computational Linguistics: EMNLP 2023*, pages 1915–1921, Singapore. Association for Computational Linguistics.
- Byung-Doh Oh and William Schuler. 2023b. Why does surprisal from larger transformer-based language models provide a poorer fit to human reading times? *Transactions of the Association for Computational Linguistics*, 11:336–350.
- Byung-Doh Oh and William Schuler. 2024. Leading whitespaces of language models’ subword vocabulary poses a confound for calculating word probabilities. *arXiv preprint arXiv:2406.10851*.
- Steven T. Piantadosi, Harry Tily, and Edward Gibson. 2011. Word lengths are optimized for efficient communication. *Proceedings of the National Academy of Sciences*, 108(9):3526–3529.
- Tiago Pimentel, Clara Meister, Ethan Wilcox, Kyle Mahowald, and Ryan Cotterell. 2023a. Revisiting the optimality of word lengths. In *Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing*, pages 2240–2255, Singapore. Association for Computational Linguistics.
- Tiago Pimentel, Clara Meister, Ethan G. Wilcox, Roger P. Levy, and Ryan Cotterell. 2023b. On the effect of anticipation on reading times. *Transactions of the Association for Computational Linguistics*, 11:1624–1642.
- Alec Radford, Jeffrey Wu, Rewon Child, David Luan, Dario Amodei, and Ilya Sutskever. 2019. Language models are unsupervised multitask learners. *OpenAI Blog*, 1(8):9.
- 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.
- Cory Shain. 2019. A large-scale study of the effects of word frequency and predictability in naturalistic reading. 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 4086–4094, Minneapolis, Minnesota. Association for Computational Linguistics.

- Cory Shain, Clara Meister, Tiago Pimentel, Ryan Cotterell, and Roger Levy. 2024. [Large-scale evidence for logarithmic effects of word predictability on reading time](#). *Proceedings of the National Academy of Sciences*, 121(10):e2307876121.
- Oleh Shliachko, Alena Fenogenova, Maria Tikhonova, Vladislav Mikhailov, Anastasia Kozlova, and Tatiana Shavrina. 2022. [mGPT: Few-shot learners go multilingual](#). *arXiv preprint arXiv:2204.07580*.
- Noam Siegelman, Sascha Schroeder, Cengiz Acartürk, Hee-Don Ahn, Svetlana Alexeeva, Simona Amenta, Raymond Bertram, Rolando Bonandrini, Marc Brysbaert, Daria Chernova, Sara Maria Da Fonseca, Nicolas Dirix, Wouter Duyck, Argyro Fella, Ram Frost, Carolina A. Gattei, Areti Kalaitzi, Nayoung Kwon, Kaidi Lõo, Marco Marelli, Timothy C. Papadopoulos, Athanassios Protopapas, Satu Savo, Diego E. Shalom, Natalia Slioussar, Roni Stein, Longjiao Sui, Analí Taboh, Veronica Tønnesen, Kerem Alp Usal, and Victor Kuperman. 2022. [Expanding horizons of cross-linguistic research on reading: The multilingual eye-movement corpus \(MECO\)](#). *Behavior Research Methods*, 54:2843–2863.
- Nathaniel J. Smith and Roger Levy. 2008. [Optimal processing times in reading: a formal model and empirical investigation](#). In *Proceedings of the Cognitive Science Society*, volume 30, pages 595–600.
- Nathaniel J. Smith and Roger Levy. 2013. [The effect of word predictability on reading time is logarithmic](#). *Cognition*, 128(3):302–319.
- Hugo Touvron, Thibaut Lavril, Gautier Izacard, Xavier Martinet, Marie-Anne Lachaux, Timothée Lacroix, Baptiste Rozière, Naman Goyal, Eric Hambro, Faisal Azhar, Aurelien Rodriguez, Armand Joulin, Edouard Grave, and Guillaume Lample. 2023. [LLaMA: Open and efficient foundation language models](#). *arXiv preprint arXiv:2302.13971*.
- Ethan Wilcox, Jon Gauthier, Jennifer Hu, Peng Qian, and Roger Levy. 2020. [On the predictive power of neural language models for human real-time comprehension behavior](#). In *Proceedings of the Cognitive Science Society*.
- Ethan Wilcox, Clara Meister, Ryan Cotterell, and Tiago Pimentel. 2023a. [Language model quality correlates with psychometric predictive power in multiple languages](#). In *Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing*, pages 7503–7511, Singapore. Association for Computational Linguistics.
- Ethan Wilcox, Tiago Pimentel, Clara Meister, Ryan Cotterell, and Roger P. Levy. 2023b. [Testing the predictions of surprisal theory in 11 languages](#). *Transactions of the Association for Computational Linguistics*, 11:1451–1470.
- Thomas Wolf, Lysandre Debut, Victor Sanh, Julien Chaumond, Clement Delangue, Anthony Moi, Pierric Cistac, Tim Rault, Rémi Louf, Morgan Funtowicz, Joe Davison, Sam Shleifer, Patrick von Platen, Clara Ma, Yacine Jernite, Julien Plu, Canwen Xu, Teven Le Scao, Sylvain Gugger, Mariama Drame, Quentin Lhoest, and Alexander M. Rush. 2020. [Transformers: State-of-the-art natural language processing](#). In *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing: System Demonstrations*, pages 38–45, Online. Association for Computational Linguistics.
- George K. Zipf. 1935. *The Psychobiology of Language*. London: Routledge.
- Vilém Zouhar, Clara Meister, Juan Gastaldi, Li Du, Mrinmaya Sachan, and Ryan Cotterell. 2023. [Tokenization and the noiseless channel](#). In *Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 5184–5207, Toronto, Canada. Association for Computational Linguistics.

## A Experimental Setup

### A.1 Sentence Comprehension

**Data.** We use four well-established reading time datasets, in which participants were given text passages to read and their reading time was recorded. For two of these datasets—Natural Stories (Futrell et al., 2018) and Brown (Smith and Levy, 2013)—measurements were collected using the self-paced paradigm (Just et al., 1982). For the other two datasets—Provo (Luke and Christianson, 2018) and Dundee (Kennedy et al., 2003)—eye-tracking movements were recorded. Each of these datasets provides the reading time each participant spent on a word. Following the works whose experiments we aim to replicate, we aggregate reading times per word (i.e., across participants). We thus analyse the average reading time participants spent on a word.

**Evaluation.** Studies of sentence comprehension are often concerned with a variable’s **predictive power**: its ability to predict sentence comprehension data. Formally, let  $\mathcal{D} = \{\mathbf{x}_n, y_n\}_{n=1}^N$  be a reading times dataset, where  $y_n \in \mathbb{R}_+$  represents the average time participants spent reading a word  $w_n$ , and  $\mathbf{x}_n \in \mathbb{R}^d$  is a vector containing a number of measurements taken on that word. Among these quantities is a word’s length (in characters) and unigram frequency. Further, let  $f_\psi$  be a regressor that takes  $\mathbf{x}_n$  as input and predicts  $y_n$ . We use  $\psi$  to denote this regressor’s parameters. A variable’s predictive power is then the change in  $\mathcal{D}$ ’s log-likelihood (denoted as  $\Delta_{\text{llh}}$ ) under two regressors: one where  $\mathbf{x}$  includes this variable ( $f_{\psi_1}$ ), and one where it does not ( $f_{\psi_2}$ ):

$$\Delta_{\text{llh}} \stackrel{\text{def}}{=} \text{llh}(f_{\psi_1}, \mathcal{D}) - \text{llh}(f_{\psi_2}, \mathcal{D}) \quad (26)$$

Here, we use this equation to measure surprisal’s predictive power. Further, we estimate this change in data log-likelihood (denoted as  $\Delta_{\text{llh}}$ ) using 10-fold cross-validation, and we leverage these results to run paired permutation tests. Finally, we account for spillover effects by including features of word  $w_n$  as well as its three preceding words in  $\mathbf{x}$ .

### A.2 Communicative Efficiency

We largely follow the setup of Pimentel et al. (2023a). We highlight the points where our setups differ below.

**Data.** We use the publicly available Wiki40b dataset (Guo et al., 2020), a large text corpus derived from Wikipedia articles. We use only the English portion of this dataset because the language models that we consider were trained solely on English data. We randomly sample a subset of the data, of size  $\approx 20\text{M}$  tokens. We do not perform any pre-processing of the text, beyond that carried out by the native HuggingFace tokenisers for the respective language models. Unigram frequencies—which are used to estimate the unigram surprisals required by the Zipfian hypothesis—are computed on a separate subset of this same dataset.

**Evaluation.** We look at correlations between word lengths and the quantities put forward by various hypotheses about the influencing factors in a lexicon’s word lengths. We expect to see that the hypotheses offering more accurate accounts of such factors have higher correlations with word lengths. In line with prior work (Piantadosi et al., 2011; Meylan and Griffiths, 2021; Levshina, 2022; Pimentel et al., 2023a), we look at Spearman and Pearson correlations.

## B Detailed Surprisal Theory Results

Model	Brown			Natural Stories			Provo			Dundee		
	Improvement	Corrected	Buggy	Improvement	Corrected	Buggy	Improvement	Corrected	Buggy	Improvement	Corrected	Buggy
gpt2-small	0.02	5.25***	5.24***	0.02	4.35***	4.33***	0.16	3.63***	3.47***	0.01	1.07***	1.06***
gpt2-medium	0.03	4.51***	4.48***	0.02	4.08***	4.07***	0.20	3.47***	3.27***	0.01	1.01***	1.00***
gpt2-large	0.04	4.53***	4.49***	0.02	3.68***	3.65***	0.21	3.10***	2.89***	0.01	0.98***	0.97***
gpt2-xl	0.03	4.23***	4.20***	0.03	3.28***	3.25***	0.19	3.09**	2.90**	0.01	0.89***	0.88***
pythia-70m	0.01	4.70***	4.69***	0.07	4.86***	4.79***	0.10	3.69***	3.59***	0.01	1.19***	1.18***
pythia-160m	0.02	4.81***	4.78***	0.05*	4.27***	4.22***	0.15	3.61***	3.46***	0.01	1.14***	1.13***
pythia-410m	0.03	4.34***	4.31***	0.05**	3.86***	3.81***	0.20	3.24***	3.04***	0.01	1.05***	1.04***
pythia-14b	0.03	4.01***	3.98***	0.04*	3.24***	3.20***	0.16	2.80***	2.64***	0.01	0.96***	0.95***
pythia-28b	0.04	3.92***	3.89***	0.03	2.96***	2.94***	0.18	3.09***	2.91***	0.01	0.88***	0.87***
pythia-69b	0.04	3.59***	3.55***	0.03	2.55***	2.52***	0.16	2.61***	2.46***	0.01	0.81***	0.80***
pythia-120b	0.04	3.51***	3.46***	0.03	2.47***	2.45***	0.17	2.36***	2.19***	0.01	0.76***	0.75***

Table 1:  $\Delta_{lh}$  between regressors with and without surprisal as a predictor.

## C Multilingual Surprisal Theory Results

In this section, we expand our surprisal theory experiments (in §6.1) to multiple languages, following a similar experimental setup to Wilcox et al. (2023b). Specifically, we analyse the MECO dataset (Siegelman et al., 2022), running our experiments on seven of its languages: Finnish, German, Greek, Hebrew, Italian, Spanish, and Turkish. We estimate surprisals for the words in these languages using mGPT (Shliachko et al., 2022)—a language model defined over the output of a bow-marking tokeniser. We thus analyse the effect of our correction **Bug Fix** ① when using this model to estimate surprisals. The  $\Delta_{lh}$  when predicting reading times on this dataset are presented in Fig. 6.

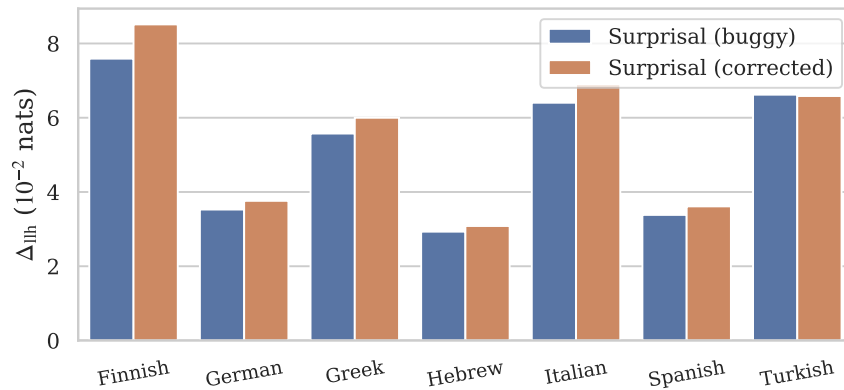


Figure 6:  $\Delta_{lh}$  between regressors with and without surprisal as a predictor in a subset of the languages in the MECO dataset (Siegelman et al., 2022).

## D Proofs of Lemmas and Theorems

### D.1 Proof of End-of-Word Tokeniser’s Theorem 1

**Lemma 3.** Let  $\mathbb{S}_{w^* \rightarrow S^*}$  be a eow-marking tokeniser. We can show the following equivalence:

$$\begin{aligned} \mathbb{P}_w(\mathbf{w}_{<t} \circ \mathcal{W}^*) &= \mathbb{P}_s(\mathbf{s}^{\mathbf{w}_{<t}} \circ \mathcal{S}^*) \\ \mathbb{P}_w(\mathbf{w}_{<t} \circ w \circ \mathcal{W}^*) &= \mathbb{P}_s(\mathbf{s}^{\mathbf{w}_{<t}} \circ \mathbf{s}^w \circ \mathcal{S}^*) \end{aligned} \quad (27)$$

*Proof.* This lemma assumes a segmentation-compatible tokeniser. Therefore, we can rely on Defn. 2, whose equation we rewrite here for convenience:

$$\mathbb{S}_{w^* \rightarrow S^*}(\mathbf{w}) = \mathbb{S}_{w \rightarrow S^*}(w_1) \circ \mathbb{S}_{w \rightarrow S^*}(w_2) \circ \dots \circ \mathbb{S}_{w \rightarrow S^*}(w_{|\mathbf{w}|}) \quad (28)$$

Further, as this tokeniser is eow-marking, we have that:  $\mathbb{S}_{\mathcal{W} \rightarrow \mathcal{S}^*} : \mathcal{W} \rightarrow \mathcal{S}_{\text{mid}} \circ \mathcal{S}_{\text{eow}}$ . We now prove the equivalences above. First, we show that  $\mathbf{w}' \in (\mathbf{w}_{<t} \circ \mathcal{W}^*) \implies \mathbb{S}_{\mathcal{W} \rightarrow \mathcal{S}^*}(\mathbf{w}') \in (\mathbf{s}^{\mathbf{w}_{<t}} \circ \mathcal{S}^*)$ ; this shows that the tokenised version of all strings  $\mathbf{w}' \in (\mathbf{w}_{<t} \circ \mathcal{W}^*)$  are present in the set  $(\mathbf{s}^{\mathbf{w}_{<t}} \circ \mathcal{S}^*)$ .

$$\mathbf{w}_{<t} \circ \mathcal{W}^* = \{ \mathbf{w}_{<t} \circ \mathbf{w}' \mid \mathbf{w}' \in \mathcal{W}^* \} \quad \text{definition of } \circ \quad (29a)$$

$$\stackrel{\triangle}{=} \left\{ \mathbb{S}_{\mathcal{W} \rightarrow \mathcal{S}^*}(\mathbf{w}_{<t} \circ \mathbf{w}') \mid \mathbf{w}' \in \mathcal{W}^* \right\} \quad \text{definition of } \stackrel{\triangle}{=} \quad (29b)$$

$$= \left\{ \mathbb{S}_{\mathcal{W} \rightarrow \mathcal{S}^*}(\mathbf{w}_{<t}) \circ \mathbb{S}_{\mathcal{W} \rightarrow \mathcal{S}^*}(\mathbf{w}') \mid \mathbf{w}' \in \mathcal{W}^* \right\} \quad \text{decomposition of } \mathbb{S}_{\mathcal{W} \rightarrow \mathcal{S}^*} \quad (29c)$$

$$= \mathbb{S}_{\mathcal{W} \rightarrow \mathcal{S}^*}(\mathbf{w}_{<t}) \circ \left\{ \mathbb{S}_{\mathcal{W} \rightarrow \mathcal{S}^*}(\mathbf{w}') \mid \mathbf{w}' \in \mathcal{W}^* \right\} \quad \text{definition of } \circ \text{ over sets} \quad (29d)$$

$$= \mathbf{s}^{\mathbf{w}_{<t}} \circ \left\{ \mathbb{S}_{\mathcal{W} \rightarrow \mathcal{S}^*}(\mathbf{w}') \mid \mathbf{w}' \in \mathcal{W}^* \right\} \quad \text{definition of } \mathbf{s}^{\mathbf{w}_{<t}} \quad (29e)$$

$$\subseteq \mathbf{s}^{\mathbf{w}_{<t}} \circ \mathcal{S}^* \quad (29f)$$

We now define the set  $\Psi_{\mathcal{S}}^{\mathbf{w}_{<t} \circ \mathcal{W}^*} \stackrel{\text{def}}{=} \{ \mathbb{S}_{\mathcal{W} \rightarrow \mathcal{S}^*}(\mathbf{w}') \mid \mathbf{w}' \in (\mathbf{w}_{<t} \circ \mathcal{W}^*) \}$ , and note that  $\mathbf{w}_{<t} \circ \mathcal{W}^* \stackrel{\triangle}{=} \Psi_{\mathcal{S}}^{\mathbf{w}_{<t} \circ \mathcal{W}^*}$ . We can thus split the probability we are computing into two parts:

$$\mathbb{P}_{\mathcal{S}}(\mathbf{s}^{\mathbf{w}_{<t}} \circ \mathcal{S}^*) = \mathbb{P}_{\mathcal{S}}(\Psi_{\mathcal{S}}^{\mathbf{w}_{<t} \circ \mathcal{W}^*}) + \mathbb{P}_{\mathcal{S}}((\mathbf{s}^{\mathbf{w}_{<t}} \circ \mathcal{S}^*) \setminus \Psi_{\mathcal{S}}^{\mathbf{w}_{<t} \circ \mathcal{W}^*}) \quad (30)$$

If we prove that  $\mathbb{P}_{\mathcal{S}}((\mathbf{s}^{\mathbf{w}_{<t}} \circ \mathcal{S}^*) \setminus \Psi_{\mathcal{S}}^{\mathbf{w}_{<t} \circ \mathcal{W}^*}) = 0$ , then it must be that  $\mathbb{P}_{\mathcal{W}}(\mathbf{w}_{<t} \circ \mathcal{W}^*) = \mathbb{P}_{\mathcal{S}}(\mathbf{s}^{\mathbf{w}_{<t}} \circ \mathcal{S}^*)$ , which completes our proof. Towards this end, we show that  $\mathbf{s}' \in (\mathbf{s}^{\mathbf{w}_{<t}} \circ \mathcal{S}^*) \implies \mathbb{S}_{\mathcal{W} \rightarrow \mathcal{S}^*}(\mathbf{s}') \in (\mathbf{w}_{<t} \circ \mathcal{W}^*)$ . For the reader's convenience, we first rewrite eq. (7) here:

$$p(\mathbf{s}) = \begin{cases} p(\mathbf{w}) & \text{if } \mathbf{s} = \mathbb{S}_{\mathcal{W} \rightarrow \mathcal{S}^*}(\mathbf{w}) \\ 0 & \text{if } \mathbf{s} \in \mathcal{S}_x \end{cases} \quad (31)$$

We now proceed with our proof.

$$\mathbf{s}^{\mathbf{w}_{<t}} \circ \mathcal{S}^* = \{ \mathbf{s}^{\mathbf{w}_{<t}} \circ \mathbf{s}' \mid \mathbf{s}' \in \mathcal{S}^* \} \quad \text{definition of } \circ \quad (32a)$$

$$\stackrel{\implies}{=} \left\{ \mathbb{S}_{\mathcal{S}^* \rightarrow \mathcal{W}^*}(\mathbf{s}^{\mathbf{w}_{<t}} \circ \mathbf{s}') \mid \mathbf{s}' \in \mathcal{S}^* \right\} \quad \text{definition of } \stackrel{\implies}{=} \quad (32b)$$

$$= \left\{ \mathbb{S}_{\mathcal{S}^* \rightarrow \mathcal{W}^*}(\mathbf{s}^{\mathbf{w}_{<t}}) \circ \mathbb{S}_{\mathcal{S}^* \rightarrow \mathcal{W}^*}(\mathbf{s}') \mid \mathbf{s}' \in \mathcal{S}^* \right\} \quad \mathbf{s}^{\mathbf{w}_{<t}} \text{ ends in } \mathcal{S}_{\text{eow}}, \text{ decomposition of } \mathbb{S}_{\mathcal{S}^* \rightarrow \mathcal{W}^*} \quad (32c)$$

$$= \mathbb{S}_{\mathcal{S}^* \rightarrow \mathcal{W}^*}(\mathbf{s}^{\mathbf{w}_{<t}}) \circ \left\{ \mathbb{S}_{\mathcal{S}^* \rightarrow \mathcal{W}^*}(\mathbf{s}') \mid \mathbf{s}' \in \mathcal{S}^* \right\} \quad \text{definition of } \circ \quad (32d)$$

$$= \mathbf{w}_{<t} \circ \left\{ \mathbb{S}_{\mathcal{S}^* \rightarrow \mathcal{W}^*}(\mathbf{s}') \mid \mathbf{s}' \in \mathcal{S}^* \right\} \quad \text{definition of } \mathbf{s}^{\mathbf{w}_{<t}} \quad (32e)$$

$$= \mathbf{w}_{<t} \circ \mathcal{W}^* \quad \text{co-domain of } \mathbb{S}_{\mathcal{S}^* \rightarrow \mathcal{W}^*} \quad (32f)$$

This result implies that any string  $\mathbf{s} \in (\mathbf{s}^{\mathbf{w}_{<t}} \circ \mathcal{S}^*)$  is either mapped to from a string  $\mathbf{w}' \in (\mathbf{w}_{<t} \circ \mathcal{W}^*)$ , or not mapped to at all by the tokenisation function  $\mathbb{S}_{\mathcal{W} \rightarrow \mathcal{S}^*}$ . (We note again that  $\mathbb{S}_{\mathcal{W} \rightarrow \mathcal{S}^*}$  only maps each  $\mathbf{w}$  to a single subword sequence  $\mathbf{s}$ , even if multiple subword sequences would be detokenised to the same  $\mathbf{w}$ .) As  $\Psi_{\mathcal{S}}^{\mathbf{w}_{<t} \circ \mathcal{W}^*}$  is defined as the set of all subword sequences mapped to from  $\mathbf{w}' \in (\mathbf{w}_{<t} \circ \mathcal{W}^*)$ , we have that  $((\mathbf{s}^{\mathbf{w}_{<t}} \circ \mathcal{S}^*) \setminus \Psi_{\mathcal{S}}^{\mathbf{w}_{<t} \circ \mathcal{W}^*}) \subseteq \mathcal{S}_x$ . It follows that the probability of the set  $(\mathbf{s}^{\mathbf{w}_{<t}} \circ \mathcal{S}^*)$  includes the probability of no other string  $\mathbf{w}' \notin (\mathbf{w}_{<t} \circ \mathcal{W}^*)$ . By the property  $((\mathbf{s}^{\mathbf{w}_{<t}} \circ \mathcal{S}^*) \setminus \Psi_{\mathcal{S}}^{\mathbf{w}_{<t} \circ \mathcal{W}^*}) \subseteq \mathcal{S}_x$ , we have that  $\mathbb{P}_{\mathcal{S}}((\mathbf{s}^{\mathbf{w}_{<t}} \circ \mathcal{S}^*) \setminus \Psi_{\mathcal{S}}^{\mathbf{w}_{<t} \circ \mathcal{W}^*}) = 0$ , which completes the proof.  $\square$

**Theorem 1.** Let  $\mathbb{S}_{\mathcal{W} \rightarrow \mathcal{S}^*}$  be a eow-marking tokeniser. Further, let  $\mathbf{s}^{\mathbf{w}} \stackrel{\text{def}}{=} \mathbb{S}_{\mathcal{W} \rightarrow \mathcal{S}^*}(\mathbf{w})$ . We can show the following equivalence:

$$\mathbb{P}_{\mathcal{W}}(\mathbf{w}_{<t} \circ \mathcal{W}^*) = \mathbb{P}_{\mathcal{S}}(\mathbf{s}^{\mathbf{w}_{<t}} \circ \mathcal{S}^*) \quad (16)$$

$$\mathbb{P}_{\mathcal{W}}(\mathbf{w}_{<t} \circ \mathbf{w} \circ \mathcal{W}^*) = \mathbb{P}_{\mathcal{S}}(\mathbf{s}^{\mathbf{w}_{<t}} \circ \mathbf{s}^{\mathbf{w}} \circ \mathcal{S}^*)$$

Further, we can compute a word's probability as:

$$p(w | \mathbf{w}_{<t}) = \underbrace{\prod_{t'=1}^{|\mathbf{s}^w|} p(\mathbf{s}_{t'}^w | \mathbf{s}^{\mathbf{w}_{<t}} \circ \mathbf{s}_{<t'}^w)}_{p(\mathbf{s}^w | \mathbf{s}^{\mathbf{w}_{<t}})} \quad (17)$$

*Proof.* The first part of this theorem utilizes Lemma 3. We can then derive the probabilities in eq. (17) as:

$$p(w | \mathbf{w}_{<t}) = \frac{\mathbb{P}(\mathbf{w}_{<t} \circ w \circ \mathcal{W}^*)}{\mathbb{P}(\mathbf{w}_{<t} \circ \mathcal{W}^*)} \quad (33a)$$

$$= \frac{\mathbb{P}(\mathbf{s}^{\mathbf{w}_{<t} \circ w} \circ \mathcal{S}^*)}{\mathbb{P}(\mathbf{s}^{\mathbf{w}_{<t}} \circ \mathcal{S}^*)} \quad (33b)$$

$$= \frac{\prod_{t'=1}^{|\mathbf{s}^{\mathbf{w}_{<t} \circ w}|} p(\mathbf{s}_{t'}^{\mathbf{w}_{<t} \circ w} | \mathbf{s}_{<t'}^{\mathbf{w}_{<t} \circ w})}{\prod_{t'=1}^{|\mathbf{s}^{\mathbf{w}_{<t}}|} p(\mathbf{s}_{t'}^{\mathbf{w}_{<t}} | \mathbf{s}_{<t'}^{\mathbf{w}_{<t}})} \quad (33c)$$

$$= \prod_{t'=|\mathbf{s}^{\mathbf{w}_{<t}}|+1}^{|\mathbf{s}^{\mathbf{w}_{<t} \circ w}|} p(\mathbf{s}_{t'}^{\mathbf{w}_{<t} \circ w} | \mathbf{s}_{<t'}^{\mathbf{w}_{<t} \circ w}) \quad (33d)$$

$$= \prod_{t'=1}^{|\mathbf{s}^w|} p(\mathbf{s}_{t'}^w | \mathbf{s}^{\mathbf{w}_{<t}} \circ \mathbf{s}_{<t'}^w) \quad (33e)$$

This completes the proof.  $\square$

## D.2 Proof of Beginning-of-Word Tokeniser's Theorem 2

**Lemma 4.** Let  $\mathbb{S}_{\mathcal{W}^* \rightarrow \mathcal{S}^*}$  be a bow-marking tokeniser. We can show the following equivalence:

$$\mathbb{P}_{\mathcal{W}}(\mathbf{w}_{<t} \circ \mathcal{W}^*) = \mathbb{P}_{\mathcal{S}}(\mathbf{s}^{\mathbf{w}_{<t}} \circ \overline{\mathcal{S}_{\text{bow}} \circ \mathcal{S}^*}) \quad (34)$$

$$\mathbb{P}_{\mathcal{W}}(\mathbf{w}_{<t} \circ w \circ \mathcal{W}^*) = \mathbb{P}_{\mathcal{S}}(\mathbf{s}^{\mathbf{w}_{<t}} \circ \mathbf{s}^w \circ \overline{\mathcal{S}_{\text{bow}} \circ \mathcal{S}^*})$$

*Proof.* This lemma assumes a segmentation-compatible tokeniser. Therefore, we can rely on Defn. 2, whose mathematical formulation we rewrite here for convenience:

$$\mathbb{S}_{\mathcal{W}^* \rightarrow \mathcal{S}^*}(\mathbf{w}) = \mathbb{S}_{\mathcal{W}^* \rightarrow \mathcal{S}^*}(w_1) \circ \mathbb{S}_{\mathcal{W}^* \rightarrow \mathcal{S}^*}(w_2) \circ \dots \circ \mathbb{S}_{\mathcal{W}^* \rightarrow \mathcal{S}^*}(w_{|\mathbf{w}|}) \quad (35)$$

Further, as this tokeniser is bow-marking, we have that:  $\mathbb{S}_{\mathcal{W}^* \rightarrow \mathcal{S}^*} : \mathcal{W} \rightarrow \mathcal{S}_{\text{bow}} \circ \mathcal{S}_{\text{mid}}^*$ . We now prove the equivalences above. First, we show that  $\mathbf{w}' \in (\mathbf{w}_{<t} \circ \mathcal{W}^*) \implies \mathbb{S}_{\mathcal{W}^* \rightarrow \mathcal{S}^*}(\mathbf{w}') \in (\mathbf{s}^{\mathbf{w}_{<t}} \circ \overline{\mathcal{S}_{\text{bow}} \circ \mathcal{S}^*})$ ; this shows that the tokenised version of all strings  $\mathbf{w}' \in (\mathbf{w}_{<t} \circ \mathcal{W}^*)$  are present in the set  $(\mathbf{s}^{\mathbf{w}_{<t}} \circ \overline{\mathcal{S}_{\text{bow}} \circ \mathcal{S}^*})$ .

$$\mathbf{w}_{<t} \circ \mathcal{W}^* = \{ \mathbf{w}_{<t} \circ \mathbf{w}' \mid \mathbf{w}' \in \mathcal{W}^* \} \quad \text{definition of } \circ \quad (36a)$$

$$\stackrel{\triangle}{=} \left\{ \mathbb{S}_{\mathcal{W}^* \rightarrow \mathcal{S}^*}(\mathbf{w}_{<t} \circ \mathbf{w}') \mid \mathbf{w}' \in \mathcal{W}^* \right\} \quad \text{definition of } \xrightarrow{\mathbb{S}_{\mathcal{W}^* \rightarrow \mathcal{S}^*}} \quad (36b)$$

$$= \left\{ \mathbb{S}_{\mathcal{W}^* \rightarrow \mathcal{S}^*}(\mathbf{w}_{<t}) \circ \mathbb{S}_{\mathcal{W}^* \rightarrow \mathcal{S}^*}(\mathbf{w}') \mid \mathbf{w}' \in \mathcal{W}^* \right\} \quad \text{decomposition of } \mathbb{S}_{\mathcal{W}^* \rightarrow \mathcal{S}^*} \quad (36c)$$

$$= \mathbb{S}_{\mathcal{W}^* \rightarrow \mathcal{S}^*}(\mathbf{w}_{<t}) \circ \left\{ \mathbb{S}_{\mathcal{W}^* \rightarrow \mathcal{S}^*}(\mathbf{w}') \mid \mathbf{w}' \in \mathcal{W}^* \right\} \quad \text{definition of } \circ \text{ over sets} \quad (36d)$$

$$= \mathbf{s}^{\mathbf{w}_{<t}} \circ \left\{ \mathbb{S}_{\mathcal{W}^* \rightarrow \mathcal{S}^*}(\mathbf{w}') \mid \mathbf{w}' \in \mathcal{W}^* \right\} \quad \text{definition of } \mathbf{s}^{\mathbf{w}_{<t}} \quad (36e)$$

$$= \mathbf{s}^{\mathbf{w}_{<t}} \circ \left( \{\text{eos}\} \cup \left( \mathcal{S}_{\text{bow}} \circ \mathcal{S}_{\text{mid}}^* \circ \left\{ \mathbb{S}_{\mathcal{W}^* \rightarrow \mathcal{S}^*}(\mathbf{w}') \mid \mathbf{w}' \in \mathcal{W}^* \right\} \right) \right) \quad (36f)$$

$$\subseteq \mathbf{s}^{\mathbf{w}_{<t}} \circ (\{\text{eos}\} \cup (\mathcal{S}_{\text{bow}} \circ \mathcal{S}_{\text{mid}}^* \circ \mathcal{S}^*)) \quad (36g)$$

$$\subseteq \mathbf{s}^{\mathbf{w}_{<t}} \circ \overline{\mathcal{S}_{\text{bow}} \circ \mathcal{S}^*} \quad (36h)$$



We now define the set  $\Psi_S^{\mathbf{w}^{<t} \circ \mathcal{W}^*} \stackrel{\text{def}}{=} \{ \mathbb{S}_{\mathcal{W}^* \rightarrow \mathcal{S}^*}(\mathbf{w}') \mid \mathbf{w}' \in (\mathbf{w}^{<t} \circ \mathcal{W}^*) \}$ , and note that  $\mathbf{w}^{<t} \circ \mathcal{W}^* \triangleq \Psi_S^{\mathbf{w}^{<t} \circ \mathcal{W}^*}$ . We can thus split the probability we are computing into two parts:

$$\mathbb{P}_S(\mathbf{s}^{\mathbf{w}^{<t}} \circ \overline{\mathcal{S}_{\text{bow}} \circ \mathcal{S}^*}) = \mathbb{P}_S(\Psi_S^{\mathbf{w}^{<t} \circ \mathcal{W}^*}) + \mathbb{P}_S((\mathbf{s}^{\mathbf{w}^{<t}} \circ \overline{\mathcal{S}_{\text{bow}} \circ \mathcal{S}^*}) \setminus \Psi_S^{\mathbf{w}^{<t} \circ \mathcal{W}^*}) \quad (37)$$

By the same logic as in Lemma 3, if we prove that  $\mathbb{P}_S((\mathbf{s}^{\mathbf{w}^{<t}} \circ \overline{\mathcal{S}_{\text{bow}} \circ \mathcal{S}^*}) \setminus \Psi_S^{\mathbf{w}^{<t} \circ \mathcal{W}^*}) = 0$ , then we have that  $\mathbb{P}_{\mathcal{W}}(\mathbf{w}^{<t} \circ \mathcal{W}^*) = \mathbb{P}_S(\mathbf{s}^{\mathbf{w}^{<t}} \circ \overline{\mathcal{S}_{\text{bow}} \circ \mathcal{S}^*})$ . To this end, we show that  $\mathbf{s}' \in (\mathbf{s}^{\mathbf{w}^{<t}} \circ \overline{\mathcal{S}_{\text{bow}} \circ \mathcal{S}^*}) \implies \mathbb{S}_{\mathcal{S}^* \rightarrow \mathcal{W}^*}(\mathbf{s}') \in (\mathbf{w}^{<t} \circ \mathcal{W}^*)$ . As with Lemma 3, this result implies that the tokenised version of no other strings  $\mathbf{w}' \notin (\mathbf{w}^{<t} \circ \mathcal{W}^*)$  are present in the set  $(\mathbf{s}^{\mathbf{w}^{<t}} \circ \overline{\mathcal{S}_{\text{bow}} \circ \mathcal{S}^*})$ , which itself implies that  $(\mathbf{s}^{\mathbf{w}^{<t}} \circ \overline{\mathcal{S}_{\text{bow}} \circ \mathcal{S}^*}) \setminus \Psi_S^{\mathbf{w}^{<t} \circ \mathcal{W}^*} \subseteq \mathcal{S}_x$ .

$$\mathbf{s}^{\mathbf{w}^{<t}} \circ \overline{\mathcal{S}_{\text{bow}} \circ \mathcal{S}^*} = \{ \mathbf{s}^{\mathbf{w}^{<t}} \circ \mathbf{s}' \mid \mathbf{s}' \in \overline{\mathcal{S}_{\text{bow}} \circ \mathcal{S}^*} \} \quad \text{definition of } \circ \quad (38a)$$

$$\implies \left\{ \mathbb{S}_{\mathcal{S}^* \rightarrow \mathcal{W}^*}(\mathbf{s}^{\mathbf{w}^{<t}} \circ \mathbf{s}') \mid \mathbf{s}' \in \overline{\mathcal{S}_{\text{bow}} \circ \mathcal{S}^*} \right\} \quad \text{definition of } \xrightarrow{\mathcal{S}^* \rightarrow \mathcal{W}^*} \quad (38b)$$

$$= \left\{ \mathbb{S}_{\mathcal{W}^* \rightarrow \mathcal{S}^*}(\mathbf{s}^{\mathbf{w}^{<t}}) \circ \mathbb{S}_{\mathcal{W}^* \rightarrow \mathcal{S}^*}(\mathbf{s}') \mid \mathbf{s}' \in \overline{\mathcal{S}_{\text{bow}} \circ \mathcal{S}^*} \right\} \quad (38c)$$

$\mathbf{s}'$  is either empty, or starts in  $\mathcal{S}_{\text{bow}}$ ,  $\mathbb{S}_{\mathcal{W}^* \rightarrow \mathcal{S}^*}$  thus decomposes (38d)

$$= \mathbb{S}_{\mathcal{W}^* \rightarrow \mathcal{S}^*}(\mathbf{s}^{\mathbf{w}^{<t}}) \circ \left\{ \mathbb{S}_{\mathcal{W}^* \rightarrow \mathcal{S}^*}(\mathbf{s}') \mid \mathbf{s}' \in \overline{\mathcal{S}_{\text{bow}} \circ \mathcal{S}^*} \right\} \quad \text{definition of } \circ \text{ over sets} \quad (38e)$$

$$= \mathbf{w}^{<t} \circ \left\{ \mathbb{S}_{\mathcal{W}^* \rightarrow \mathcal{S}^*}(\mathbf{s}') \mid \mathbf{s}' \in \overline{\mathcal{S}_{\text{bow}} \circ \mathcal{S}^*} \right\} \quad \text{definition of } \mathbf{s}^{\mathbf{w}^{<t}} \quad (38f)$$

$$= \mathbf{w}^{<t} \circ \mathcal{W}^* \quad \text{co-domain of } \mathbb{S}_{\mathcal{W}^* \rightarrow \mathcal{S}^*} \quad (38g)$$

Since  $\mathbf{s} \in ((\mathbf{s}^{\mathbf{w}^{<t}} \circ \overline{\mathcal{S}_{\text{bow}} \circ \mathcal{S}^*}) \setminus \Psi_S^{\mathbf{w}^{<t} \circ \mathcal{W}^*}) \implies \mathbf{s} \in \mathcal{S}_x$ , we have that  $\mathbb{P}_S((\mathbf{s}^{\mathbf{w}^{<t}} \circ \overline{\mathcal{S}_{\text{bow}} \circ \mathcal{S}^*}) \setminus \Psi_S^{\mathbf{w}^{<t} \circ \mathcal{W}^*}) = 0$ , which completes the proof.  $\square$

**Theorem 2.** Let  $\mathbb{S}_{\mathcal{W}^* \rightarrow \mathcal{S}^*}$  be a bow-marking tokeniser. Further, let  $\bar{\cdot}$  represent the union of a set with eos, e.g.,  $\overline{\mathcal{S}_{\text{bow}}} = \mathcal{S}_{\text{bow}} \cup \{\text{eos}\}$ . We can show the following equivalence:

$$\mathbb{P}_{\mathcal{W}}(\mathbf{w}^{<t} \circ \mathcal{W}^*) = \mathbb{P}_S(\mathbf{s}^{\mathbf{w}^{<t}} \circ \overline{\mathcal{S}_{\text{bow}} \circ \mathcal{S}^*}) \quad (19)$$

$$\mathbb{P}_{\mathcal{W}}(\mathbf{w}^{<t} \circ w \circ \mathcal{W}^*) = \mathbb{P}_S(\mathbf{s}^{\mathbf{w}^{<t}} \circ \mathbf{s}^w \circ \overline{\mathcal{S}_{\text{bow}} \circ \mathcal{S}^*})$$

Further, we can compute a word's probability as:

$$p(w \mid \mathbf{w}^{<t}) = \underbrace{\prod_{t'=1}^{|\mathbf{s}^w|} p(\mathbf{s}_{t'}^w \mid \mathbf{s}^{\mathbf{w}^{<t}} \circ \mathbf{s}_{<t'}^w)}_{p(\mathbf{s}^w \mid \mathbf{s}^{\mathbf{w}^{<t}})} \underbrace{\frac{\sum_{\{s \in \overline{\mathcal{S}_{\text{bow}}}\}} p(s \mid \mathbf{s}^{\mathbf{w}^{<t}} \circ \mathbf{s}^w)}{\sum_{\{s \in \overline{\mathcal{S}_{\text{bow}}}\}} p(s \mid \mathbf{s}^{\mathbf{w}^{<t}})}}_{\text{Bug Fix } \textcircled{1}} \quad (20)$$

*Proof.* The first part of this theorem simply re-writes Lemma 4. We now derive the probabilities in eq. (20)

as:

$$p(w \mid \mathbf{w}_{<t}) = \frac{\mathbb{P}(\mathbf{w}_{<t} \circ w \circ \mathcal{W}^*)}{\mathbb{P}(\mathbf{w}_{<t} \circ \mathcal{W}^*)} \quad (39a)$$

$$= \frac{\sum_{\{s \in \overline{\mathcal{S}}_{\text{bow}}\}} \mathbb{P}(\mathbf{s}^{\mathbf{w}_{<t} \circ w} \circ s \circ \mathcal{S}^*)}{\sum_{\{s \in \overline{\mathcal{S}}_{\text{bow}}\}} \mathbb{P}(\mathbf{s}^{\mathbf{w}_{<t}} \circ s \circ \mathcal{S}^*)} \quad (39b)$$

$$= \frac{\sum_{\{s \in \overline{\mathcal{S}}_{\text{bow}}\}} p(s \mid \mathbf{s}^{\mathbf{w}_{<t} \circ w}) \prod_{t'=1}^{|\mathbf{s}^{\mathbf{w}_{<t} \circ w}|} p(\mathbf{s}_{t'}^{\mathbf{w}_{<t} \circ w} \mid \mathbf{s}_{<t'}^{\mathbf{w}_{<t} \circ w})}{\sum_{\{s \in \overline{\mathcal{S}}_{\text{bow}}\}} p(s \mid \mathbf{s}^{\mathbf{w}_{<t}}) \prod_{t'=1}^{|\mathbf{s}^{\mathbf{w}_{<t}}|} p(\mathbf{s}_{t'}^{\mathbf{w}_{<t}} \mid \mathbf{s}_{<t'}^{\mathbf{w}_{<t}})} \quad (39c)$$

$$= \frac{\prod_{t'=1}^{|\mathbf{s}^{\mathbf{w}_{<t} \circ w}|} p(\mathbf{s}_{t'}^{\mathbf{w}_{<t} \circ w} \mid \mathbf{s}_{<t'}^{\mathbf{w}_{<t} \circ w}) \sum_{\{s \in \overline{\mathcal{S}}_{\text{bow}}\}} p(s \mid \mathbf{s}^{\mathbf{w}_{<t} \circ w})}{\prod_{t'=1}^{|\mathbf{s}^{\mathbf{w}_{<t}}|} p(\mathbf{s}_{t'}^{\mathbf{w}_{<t}} \mid \mathbf{s}_{<t'}^{\mathbf{w}_{<t}}) \sum_{\{s \in \overline{\mathcal{S}}_{\text{bow}}\}} p(s \mid \mathbf{s}^{\mathbf{w}_{<t}})} \quad (39d)$$

$$= \frac{\prod_{t'=|\mathbf{s}^{\mathbf{w}_{<t} \circ w}|}^{|\mathbf{s}^{\mathbf{w}_{<t} \circ w}|} p(\mathbf{s}_{t'}^{\mathbf{w}_{<t} \circ w} \mid \mathbf{s}_{<t'}^{\mathbf{w}_{<t} \circ w}) \sum_{\{s \in \overline{\mathcal{S}}_{\text{bow}}\}} p(s \mid \mathbf{s}^{\mathbf{w}_{<t} \circ w})}{\sum_{\{s \in \overline{\mathcal{S}}_{\text{bow}}\}} p(s \mid \mathbf{s}^{\mathbf{w}_{<t}})} \quad (39e)$$

$$= \prod_{t'=1}^{|\mathbf{s}^w|} p(\mathbf{s}_{t'}^w \mid \mathbf{s}^{\mathbf{w}_{<t}} \circ \mathbf{s}_{<t'}^w) \frac{\sum_{\{s \in \overline{\mathcal{S}}_{\text{bow}}\}} p(s \mid \mathbf{s}^{\mathbf{w}_{<t} \circ w})}{\sum_{\{s \in \overline{\mathcal{S}}_{\text{bow}}\}} p(s \mid \mathbf{s}^{\mathbf{w}_{<t}})} \quad (39f)$$

This completes the proof.  $\square$

### D.3 Theorem of Non-eow-marking Final-word Tokeniser's

**Theorem 3.** Let  $\mathcal{W}^* \xrightarrow{\mathcal{S}} \mathcal{S}^*$  be a eow-marking tokeniser with unmarked final word. We can show the following equivalence:

$$\begin{aligned} \mathbb{P}_{\mathcal{W}}(\mathbf{w}_{<t} \circ \mathcal{W}^*) &= \mathbb{P}_{\mathcal{S}}(\mathbf{s}^{\mathbf{w}_{<t}} \circ \mathcal{S}^*) \\ \mathbb{P}_{\mathcal{W}}(\mathbf{w}_{<t} \circ w \circ \mathcal{W}^*) &= \mathbb{P}_{\mathcal{S}}((\mathbf{s}^{\mathbf{w}_{<t}} \circ \mathbf{s}^w \circ \mathcal{S}^*) \cup \{\mathbf{s}^{\mathbf{w}_{<t}} \circ \mathbf{s}_{\text{mid}}^w\}) \end{aligned} \quad (40)$$

Further, we can compute a word's probability as:

$$p(w \mid \mathbf{w}_{<t}) = \underbrace{\left( p(\mathbf{s}_{\text{mid}}^w \mid \mathbf{s}^{\mathbf{w}_{<t}}) \sum_{s \in \mathcal{S}_{!?}} p(s \mid \mathbf{s}^{\mathbf{w}_{<t}} \circ \mathbf{s}_{\text{mid}}^w) \right)}_{\text{Bug Fix } \textcircled{2}} + p(\mathbf{s}^w \mid \mathbf{s}^{\mathbf{w}_{<t}}) \quad (41)$$

### D.4 Theorem of Non-bow-marking First-word Tokeniser's

**Theorem 4.** Let  $\mathcal{W}^* \xrightarrow{\mathcal{S}} \mathcal{S}^*$  be a bow-marking tokeniser with unmarked first words. We can show the following equivalence:

$$\begin{aligned} \mathbb{P}_{\mathcal{W}}(\mathcal{W}^*) &= \mathbb{P}_{\mathcal{S}}(\overline{\mathcal{S}_{\text{mid}} \circ \mathcal{S}^*}) \\ \mathbb{P}_{\mathcal{W}}(w \circ \mathcal{W}^*) &= \mathbb{P}_{\mathcal{S}}(\mathbf{s}_{\text{mid}}^w \circ \overline{\mathcal{S}_{\text{bow}} \circ \mathcal{S}^*}) \end{aligned} \quad (42)$$

Further, we can compute a word's probability as:

$$p(w \mid \mathbf{w}_{<t}) = p(\mathbf{s}_{\text{mid}}^w \mid \text{""}) \underbrace{\frac{\sum_{\{s \in \overline{\mathcal{S}}_{\text{bow}}\}} p(s \mid \mathbf{s}^w)}{\sum_{\{s \in \overline{\mathcal{S}}_{\text{mid}}\}} p(s \mid \text{""})}}_{\text{Bug Fix } \textcircled{3}} \quad (43)$$