

# A Formal Perspective on Byte-Pair Encoding

Vilém Zouhar<sup>E</sup> Clara Meister<sup>E</sup> Juan Luis Gastaldi<sup>E</sup> Li Du<sup>J</sup>  
Tim Vieira<sup>J</sup> Mrinmaya Sachan<sup>E</sup> Ryan Cotterell<sup>E</sup>

ETH Zürich<sup>E</sup> Johns Hopkins University<sup>J</sup>  
{vzouhar,meistecl,gjuan,msachan,ryan.cotterell}@ethz.ch  
{leodu,timv}@cs.jhu.edu

## Abstract

Byte-Pair Encoding (BPE) is a popular algorithm used for tokenizing data in NLP, despite being devised initially as a compression method. BPE appears to be a greedy algorithm at face value, but the underlying optimization problem that BPE seeks to solve has not yet been laid down. We formalize BPE as a combinatorial optimization problem. Via submodular functions, we prove that the iterative greedy version is a  $\frac{1}{\sigma(\mu^*)}(1 - e^{-\sigma(\mu^*)})$ -approximation of an optimal merge sequence, where  $\sigma(\mu^*)$  is the total backward curvature with respect to the optimal merge sequence  $\mu^*$ . Empirically the lower bound of the approximation is  $\approx 0.37$ .

We provide a faster implementation of BPE which improves the runtime complexity from  $\mathcal{O}(NM)$  to  $\mathcal{O}(N \log M)$ , where  $N$  is the sequence length and  $M$  is the merge count. Finally, we optimize the brute-force algorithm for optimal BPE using memoization.

## 1 Introduction

Byte-Pair Encoding (BPE) is a popular technique for building and applying an encoding scheme to natural language texts. It is one of the most common tokenization methods used for language models (Radford et al., 2019; Bostrom and Durrett, 2020; Brown et al., 2020; Scao et al., 2022) as well as for various other conditional language modeling tasks, e.g., machine translation (Ding et al., 2019) and chatbots (Zhang et al., 2020). Despite having been popularized by Sennrich et al. (2016) in NLP as a tokenization scheme, BPE has its roots in the compression literature, where Gage (1994) introduced the method as a faster alternative to Lempel–Ziv–Welch (Welch, 1984; Cover and Thomas, 2006, 13.4). However, the ubiquity of BPE notwithstanding, the formal underpinnings of the algorithm are underexplored, and there are no existing proven guarantees about BPE’s performance.

The training and applying of BPE are traditionally presented as greedy algorithms, but the exact

optimization problems they seek to solve are neither presented in the original work of Gage (1994) nor in the work of Sennrich et al. (2016). We fill this void by offering a clean formalization of BPE training as maximizing a function we call compression utility<sup>1</sup> over a specific combinatorial space, which we define in Definition 2.3. Unexpectedly, we are then able to prove a bound on BPE’s approximation error using total backward curvature  $\sigma(\mu^*)$  (Zhang et al., 2015). Specifically, we find the ratio of compression utilities between the greedy method and the optimum is bounded below by  $\frac{1}{\sigma(\mu^*)}(1 - e^{-\sigma(\mu^*)})$ , which we find empirically  $\approx 0.37$  for  $\hat{\sigma}(\mu^*) = 2.5$ . Our proof of correctness hinges on the theory of submodular functions (Krause and Golovin, 2014; Bilmes, 2022).<sup>2</sup> Indeed, we are able to prove that compression utility is a special kind of submodular function (Malekian, 2009) over a constrained space. And, despite the presence of the length constraint, which we expound upon formally in §3, we are able to prove a similar bound to  $1 - 1/e$  as in the unconstrained case (Alaei et al., 2010).

Additionally, we give a formal analysis of greedy BPE’s runtime and provide a speed-up over the original implementation (Gage, 1994; Sennrich et al., 2016). Our runtime improvement stems from the development of a nuanced data structure that allows us to share work between iterations of the greedy procedure and that lends itself to an amortized analysis. Specifically, given a string with  $N$  characters with a desired merge count of  $M$  (usually  $N \gg M$ ), our implementation runs in  $\mathcal{O}(N \log M)$ , an improvement over the  $\mathcal{O}(NM)$ -time algorithm presented by Sennrich et al. (2016) and the  $\mathcal{O}(N \log N)$  analysis presented by Kudo

<sup>1</sup>How much space the compression saves (Definition 2.5).

<sup>2</sup>The proof further relies on a specific property of problem which BPE optimizes that we term hierarchical sequence submodularity. Hierarchical sequence submodularity neither follows from nor implies sequence submodularity, but, nevertheless, bears some superficial similarity to sequence submodularity—hence, our choice of name.

and Richardson (2018). Finally, our formalism allows us to construct an exact program for computing an optimal solution to the BPE training problem. Unfortunately, the algorithm runs in exponential time, but it is still significantly faster than a naïve brute-force approach.

Our work should give NLP practitioners confidence that BPE is a wise choice for learning a subword vocabulary based on compression principles. In general, such constrained submodular maximization problems are hard (Lovász, 1983). While we do not have a proof that the BPE problem specifically is NP-hard, it does not seem likely that we could find an efficient algorithm for the problem. Regarding the runtime, our implementation of greedy BPE runs nearly linearly in the length of the string which would be hard to improve unless we plan to not consider the entire string.

## 2 Formalizing Byte-Pair Encoding

We first provide a brief intuition for the BPE training problem and the greedy algorithm that is typically employed to solve it. Then, we will develop a formalization of BPE using the tools of combinatorial optimization, rather than as a procedure.<sup>3</sup>

### 2.1 A Worked Example

merge 1	<b>p i</b> c k e d	p <b>i c k</b> l e d	p i <b>c k l e</b> s
merge 2	p i <b>c k</b> e d	p i c k <b>l e</b> d	p i c k <b>l e</b> s
merge 3	p i c k <b>e d</b>	p i c k <b>l e</b> d	p i c k <b>l e</b> s
merge 4	p i c k <b>e d</b>	p i c k <b>l e</b> d	p i c k <b>l e</b> s
merge 5	p i c k <b>e d</b>	p i c k <b>l e</b> d	p i c k <b>l e</b> s
final	p i c k e d	p i c k l e d	p i c k l e s

Example 1: Compression of the text *picked pickled pickles* using 5 greedy merges according to the greedy BPE algorithm. The most frequently occurring pair of vocabulary items is highlighted and subsequently merged. The merge sequence is  $\langle [p,i], [c,k], [pi,ck], [e,d], [pick,l] \rangle$  (notation simplified for clarity).

Consider the string in Example 1: *picked pickled pickles*. We wish to create a compact representation of this string, where compactness is quantified in terms of the number of symbols (i.e., vocabulary units) required to precisely encode the string. The free parameter is the vocabulary that we will use to construct this representation, albeit the total size

<sup>3</sup>Note that we are interested in the optimality of the algorithm for creating the subword vocabulary in terms of compression and not the optimality of the encoding in terms of coding-theoretic metrics such as *efficiency*. This aspect of BPE is explored by Zouhar et al. (2023).

of the chosen vocabulary is often a constraint.<sup>4</sup> In our example, let’s assume we are allowed a maximum number of 13 symbols in the vocabulary<sup>5</sup> with which we can encode our string. The question is: “How can we select these symbols to achieve our goal of compactness under this constraint?”

Let us first consider the simple choice of using all the characters present in the string as our vocabulary: This scheme leads to a representation with a length of 22 units, including spaces. In order to decrease this length (while retaining all information present in the original string), we would need to add an additional symbol to our vocabulary: one with which we can replace co-occurrences of two symbols. But how should we choose this entry? One strategy—the one employed by the BPE algorithm—is to use the concatenation of the adjacent units  $a b$  that occur with the highest frequency in our string; all occurrences of these adjacent units could then be replaced with a single new unit  $ab$ . We refer to this as a **merge**, which we later define and denote formally as  $[a, b]$ . In Example 1, the first merge is  $[p, i]$ , and leads to a representation of length 19 with vocabulary size of 9+1. We can iteratively repeat the same process; the application of 5 total merges results in the vocabulary units *pick*, *pickl*, *ed*, *e*, and *s*. These **subwords**<sup>6</sup> allow us to represent our original string using just 9+1 symbols. If we continued merging, the text representation would become shorter (in terms of number of symbols required to create the representation) but the merge count (and vocabulary size) would grow. Therefore, the number of merges  $M$ , or also the merge count, is a hyperparameter to the whole procedure. The procedure outlined above is exactly the greedy algorithm for BPE proposed by Gage (1994). We provide a minimal implementation in Python in Code 1.

We will define the compression gain of a merge at any given step of the algorithm, corresponding to the number of occurrences where a merge can be applied. The compression gain of a merge does not always correspond to the frequency of adjacent merge components in that string, due to possible overlaps. Consider, for instance, the string *aaa* and

<sup>4</sup>We require a unique encoding for each item, which implies that encoding size will be dependent on the total number of vocabulary items (e.g. the dimension of a one-hot encoding or the number of bits required to encode the text).

<sup>5</sup>Typically, all the symbols in  $\Sigma$  are part of the vocabulary so that all texts can be represented, even with lower efficiency.

<sup>6</sup>The term **subword** corresponds to a merge yield (Definition 2.4). We use ‘subword’ and ‘merge’ interchangeably.

the merge  $[a, a]$ . The frequency of  $aa$  is 2, but the merge can be applied only once ( $[a, a]a$ ). While [Gage \(1994\)](#) and [Sennrich et al. \(2016\)](#) admit overlapping pair counts, [Kudo and Richardson \(2018\)](#)’s popular implementation adjusts the algorithm to disregard the overlaps. We stick to the latter, which is more suitable from the optimization standpoint adopted here.

```

1 from collections import Counter
2 from typing import Union, Tuple, List
3
4 def bpe(xs: Union[str, List], V: int):
5     for _ in range(V):
6         pairs = Counter(zip(xs, xs[1:]))
7         top_pair = pairs.most_common(1)[0][0]
8         xs = merge(list(xs), top_pair)
9     return xs
10
11 def merge(xs: List, pair: Tuple):
12     ys = []
13     while xs:
14         if tuple(xs[:2]) == pair:
15             ys.append(pair)
16             xs = xs[2:]
17         else:
18             ys.append(xs.pop(0))
19     return ys

```

Code 1: A minimal implementation of [Sennrich et al.’s \(2016\)](#) greedy algorithm for BPE in Python. See Code 2 for a version with overlap-adjusted counts.

## 2.2 Merges

The fundamental building block of the BPE problem is a merge, which we define formally below. Informally, a merge is the action of creating a new symbol out of two existing ones. Out of convention, we also refer to the resulting object as a merge.

**Definition 2.1.** Let  $\Sigma$  be an alphabet, a finite, non-empty set. The set of all **merges** over  $\Sigma$  is the smallest set of pairs  $\Upsilon_\Sigma$  with the following closure property:

- $\sigma \in \Sigma \implies \sigma \in \Upsilon_\Sigma$  (called **trivial merges**);
- $\mu', \mu'' \in \Upsilon_\Sigma \implies [\mu', \mu''] \in \Upsilon_\Sigma$

where we denote the non-trivial elements of  $\Upsilon_\Sigma$  as  $\mu = [\mu', \mu'']$ . A **merge sequence** is a sequence of merges, which we denote  $\boldsymbol{\mu} = \langle \mu_1, \dots, \mu_N \rangle \in \Upsilon_\Sigma^*$ .<sup>7</sup>

It is perhaps easiest to understand the concept of a merge through an example.

**Example 2.2.** Given the alphabet  $\Sigma = \{a, b, c\}$ , the following are some of the elements of  $\Upsilon_\Sigma$ :  $[a, b]$ ,  $[a, [a, b]]$ , and  $[[a, b], [a, c]]$ . We obtain a merge sequence by arranging these merges into an ordering  $\boldsymbol{\mu} = \langle [a, b], [a, [a, b]], [[a, b], [a, c]] \rangle \in \Upsilon_\Sigma^*$ .

<sup>7</sup> $(\cdot)^*$  is the Kleene closure.

Note that the strings corresponding to the merges in a merge sequence—along with the characters that make up the set of trivial merges—determine a **vocabulary**, to be used in downstream applications.<sup>8</sup> The greedy BPE algorithm constructs a merge sequence iteratively by picking each merge as the pairing of neighbouring symbols in the current sequence of symbols that is being processed. For instance, the sequence  $\boldsymbol{\mu}$  in Example 2.2 is not valid since it does not contain the merge  $[a, c]$  before the third element  $[[a, b], [a, c]]$ .

**Definition 2.3.** We define a merge sequence  $\boldsymbol{\mu} = \langle \mu_1, \dots, \mu_N \rangle \in \Upsilon_\Sigma^*$  to be **valid** if, for every  $\mu_n$ , it holds that  $\mu_n = [\mu', \mu'']$ , where for  $\mu \in \{\mu', \mu''\}$ ,  $\mu = \mu_k$  with  $k < n$ , or  $\mu \in \Sigma$ . We denote the set of valid merge sequences  $\mathcal{M}_{\Upsilon_\Sigma}$ .

Note that  $\mathcal{M}_{\Upsilon_\Sigma}$  is closed under concatenation, i.e., for two valid merge sequences  $\boldsymbol{\mu}', \boldsymbol{\mu}'' \in \mathcal{M}_{\Upsilon_\Sigma}$ , we have that  $\boldsymbol{\mu}'\boldsymbol{\mu}'' \in \mathcal{M}_{\Upsilon_\Sigma}$ ,<sup>9</sup> where we use  $\boldsymbol{\mu}\boldsymbol{\mu}'$  to denote the sequence concatenation of  $\boldsymbol{\mu}$  and  $\boldsymbol{\mu}'$ .

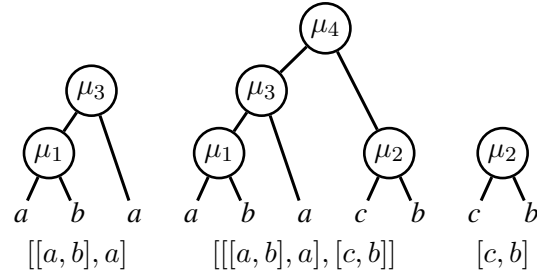


Figure 1: Application of the merge sequence  $\boldsymbol{\mu} = \langle [a, b], [c, b], [[a, b], a], [[[a, b], a], [c, b]] \rangle$  on the string  $x = abaabacbc$ . The result can be represented as an ordered forest. Each tree is associated with a subword in the text:  $aba$ ,  $abacb$ , and  $cb$ .

**Applying Merge Sequences.** Given some string  $x \in \Sigma^*$ , we can derive the representation of that string according to the merge sequence  $\boldsymbol{\mu} = \langle \mu_1, \dots, \mu_N \rangle$  by iteratively **applying** each merge  $\mu_n$ . Note that by the definition of  $\Upsilon_\Sigma^*$ , we can trivially lift a string  $x = \langle \sigma_1, \sigma_2, \dots \rangle$  to a merge sequence by treating each of its characters  $\sigma_i \in \Sigma$  as merges. Thus, we define this procedure more generally in terms of some arbitrary  $\bar{\boldsymbol{\mu}} \in \Upsilon_\Sigma^*$ . Concretely, we denote the application of a merge  $\mu_n$  to  $\bar{\boldsymbol{\mu}}$  as  $\text{APPLY}_{\mu_n}(\bar{\boldsymbol{\mu}})$ . As suggested by Code 1 (line 11), this action consists of replacing all  $\bar{\mu}_k, \bar{\mu}_{k+1}$  in  $\bar{\boldsymbol{\mu}}$  such that  $(\bar{\mu}_k, \bar{\mu}_{k+1}) = \mu_n$  by  $\mu_n$  itself, in a left-to-right fashion. We thus obtain a

<sup>8</sup>I.e., the size of the vocabulary is  $|\boldsymbol{\mu}| + |\Sigma|$ .

<sup>9</sup>The merge sequence can contain the same merges multiple times and still be valid. Only the later occurrences of the merge will not reduce the representation size.

new  $\bar{\mu} \in \Upsilon_{\Sigma}^*$ , to which a new single merge can be applied. We lift `APPLY` to a merge sequence  $\mu$  by simply repeating the application of `APPLY` on  $\bar{\mu}^{(n)}$  for the successive  $\mu_n$  in  $\mu$ ; accordingly, we denote this procedure as `APPLY` $_{\mu}(\bar{\mu})$ . As a result, we obtain  $\bar{\mu}^{(\mu)}$ , which is a non-overlapping ordered forest, i.e., a partial bracketing of the original string  $x$ . We provide an example in Fig. 1. Note that the application of the merge sequence is deterministic.

**String Yields.** We now define a conceptually reverse operation to applying merges, i.e., deriving a string from structured  $\bar{\mu}^{(n)}$ .

**Definition 2.4.** *The yield of a single  $\bar{\mu} \in \Upsilon_{\Sigma}$ , denoted as  $\text{YIELD}(\bar{\mu})$ , is defined recursively:*

$$\text{YIELD}(\bar{\mu}) = \begin{cases} \text{YIELD}(\bar{\mu}')\text{YIELD}(\bar{\mu}'') & \text{if } \bar{\mu} = [\bar{\mu}', \bar{\mu}''] \\ \mu & \text{if } \bar{\mu} \in \Sigma \end{cases} \quad (1)$$

As an example,  $\text{YIELD}([[[a, a], [[c, b], c]])$  is *aacbc*. For a given  $\bar{\mu}$ ,  $\text{YIELD}$  is applied sequentially. The resulting characters can then be concatenated to derive a single string. The yield operation can also be used to derive vocabulary units—often referred to as subwords; explicitly, the yields of individual merges in a sequence  $\mu$  can be used to form a vocabulary.

Strictly speaking, in Sennrich et al.’s (2016) implementation of BPE, the elements of the merge sequences  $\mu$  are not of the form  $\mu_n = [\mu', \mu''] \in \Upsilon_{\Sigma}$ , but rather  $\mu_n = [\text{YIELD}(\mu'), \text{YIELD}(\mu'')] \in \Sigma^* \times \Sigma^*$ , i.e., rather than consisting of prior merges as in our formalization, the merges of Sennrich et al.’s (2016) consist of the yields of those merges. This introduces an ambiguity with respect to our formalization since: for a given merge sequence in that implementation, more than one sequence  $\mu \in \Upsilon_{\Sigma}^*$  could correspond, some of which would not be valid. As an example, consider the sequence  $\langle [a, b], [ab, c], [abc, d] \rangle$  which could correspond to either  $\langle [a, b], [[a, b], c], [[[a, b], c], d] \rangle$  or  $\langle [a, b], [[a, b], c], [[a, [b, c]], d] \rangle$ , the last of which is invalid. However, it turns out that this is not an issue for us: by construction, the successive elements of the sequence are determined by the previous ones (cf. Alg. 1), which means that, in fact there is no ambiguity, and the merge sequences in Sennrich et al.’s (2016) implementation always correspond to what our formalization defines as a valid merge sequence.

---

**Algorithm 1** Iterative Greedy BPE (slow).

**Inputs:** sequence  $x$ , merge count  $M$

**Output:** merge sequence  $\mu$ , tokenized sequence  $x$   
`PAIRFREQ` are non-overlapping pair frequencies

---

```

1:  $\mu \leftarrow \langle \rangle$ 
2: for  $i$  in  $\{0, \dots, M\}$  do
3:    $\mu \leftarrow \underset{(\mu', \mu'') \in \text{set}(x)^2}{\text{argmax}} \text{PAIRFREQ}(x, (\mu', \mu''))$ 
4:    $x \leftarrow \text{APPLY}(\mu, x)$ 
5:    $\mu \leftarrow \mu \circ \langle \mu \rangle$ 
6: end for
7: return  $\mu, x$ 

```

---

### 2.3 The BPE Training Optimization Problem

We now define the BPE training task as a combinatorial optimization problem. The objective we seek to optimize is the compression utility of the chosen merge sequence (taken with respect to a string), which we define below.

**Definition 2.5.** *Let  $x \in \Sigma^*$  be a string. We define the **compression utility** of a valid merge sequence  $\mu$  applied to  $x$  as the following function:*

$$\kappa_x(\mu) = |x| - |\text{APPLY}_{\mu}(x)| \quad (2)$$

*Note that for any merge sequence  $\mu$ ,  $\kappa_x(\mu) \geq 0$  and we take  $\kappa_x(\langle \rangle) = 0$ . Then, for any merge sequence  $\mu' = \langle \mu'_1, \dots, \mu'_{|x|-1} \rangle$  of length  $|x| - 1$  where every merge produces replacements, we have  $\kappa_x(\mu') = |x| - 1$  (see proof of Theorem 4.2).*

We can further define the **compression gain** of two merge sequences with respect to each other.

**Definition 2.6.** *The **compression gain** of  $\mu'$  with respect to a sequence  $\mu$ , denoted as  $\kappa_x(\mu' | \mu)$ , is defined as*

$$\kappa_x(\mu\mu') - \kappa_x(\mu). \quad (3)$$

*Similarly, the **compression gain** of a single merge  $\mu$  with respect to a sequence  $\mu$ , denoted as  $\kappa_x(\mu | \mu)$ , is defined as  $\kappa_x(\mu\mu) - \kappa_x(\mu)$ .*

We use the compression gain to later make a sequence of observations which leads to proving the function submodularity and eventually its approximation bound of the BPE training algorithm.

Now, armed with Definition 2.5, we can formally state our optimization problem. In words, *we seek to find a valid merge sequence  $\mu$  with length of  $M$  that maximizes the compression utility  $\kappa_x(\cdot)$  for a pre-specified string  $x \in \Sigma^*$ . We write this combinatorial optimization problem more formally*



as follows:<sup>10</sup>

$$\boldsymbol{\mu}^* = \operatorname{argmax}_{\substack{\boldsymbol{\mu} \in \mathcal{M}_{\Upsilon_{\Sigma}} \\ |\boldsymbol{\mu}|=M}} \kappa_{\mathbf{x}}(\boldsymbol{\mu}) \quad (4)$$

The most common procedure found in the NLP literature for solving Eq. (4) is a greedy algorithm (Gage, 1994; Sennrich et al., 2016). The implementation of Gage’s (1994) algorithm presented by Sennrich et al. (2016) runs in  $\mathcal{O}(NM)$  time ( $N = |\mathbf{x}|, M = |\boldsymbol{\mu}^*|$ ). We describe this greedy algorithm in detail in §3 and provide a novel theoretical result: *The algorithm comes with a  $\frac{1}{\sigma(\boldsymbol{\mu}^*)}(1 - e^{-\sigma(\boldsymbol{\mu}^*)})$  bound on its approximation error of Eq. (4).* In §4, we further offer an asymptotic speed-up to Sennrich et al.’s (2016) algorithm, reducing its runtime to  $\mathcal{O}(N \log M)$ . Finally, for completeness, we offer an exact program for finding an optimal valid merge sequence in §5. While this algorithm runs in exponential time, which prevents it to be used in real applications, it is still faster than the brute-force counterpart.

### 3 A Greedy Approximation of BPE

We demonstrate that, for any string  $\mathbf{x} \in \Sigma^*$ , the following bound holds

$$\frac{\kappa_{\mathbf{x}}(\boldsymbol{\mu}^\dagger)}{\kappa_{\mathbf{x}}(\boldsymbol{\mu}^*)} \geq \frac{1}{\sigma(\boldsymbol{\mu}^*)}(1 - e^{-\sigma(\boldsymbol{\mu}^*)}) \quad (5)$$

where, as in the previous section,  $\boldsymbol{\mu}^\dagger$  is the valid merge sequence output by the greedy algorithm and  $\boldsymbol{\mu}^*$  is an optimal valid merge sequence. To prove this bound, we rely heavily on the theory of submodularity (Krause and Golovin, 2014; Bilmes, 2022).

#### 3.1 Properties of Compression Utility ( $\kappa$ )

We start by proving some useful facts about the compression utility function  $\kappa_{\mathbf{x}}$ . Specifically, we first show that  $\kappa_{\mathbf{x}}$  is a specific type of monotone non-decreasing submodular sequence function, which we make precise in the following definitions.

**Definition 3.1.** *A real-valued function  $f$  over valid merge sequences is **monotone non-decreasing** if, for all  $\boldsymbol{\mu} \in \mathcal{M}_{\Upsilon_{\Sigma}}$  and for all  $n \in \mathbb{N}$ , it holds that  $f(\boldsymbol{\mu}_{<n}) \geq f(\boldsymbol{\mu}_{<n-1})$ , where  $\boldsymbol{\mu}_{<n} \stackrel{\text{def}}{=} \langle \mu_1, \dots, \mu_{n-1} \rangle$ .*

<sup>10</sup>In practice, it does not happen that  $|\mathbf{x}| < M$  and so we use  $|\boldsymbol{\mu}^*| = M$  for convenience instead of  $|\boldsymbol{\mu}^*| \leq M$ .

**Proposition 3.2.** *Let  $\kappa_{\mathbf{x}}$  be the compression utility function. Then, for a fixed  $\mathbf{x} \in \Sigma^*$ ,  $\kappa_{\mathbf{x}}(\cdot)$  is monotone (Definition 3.1).*

*Proof.* For all  $n \in \mathbb{N}$ , we have that  $\kappa_{\mathbf{x}}(\boldsymbol{\mu}_{<n}) = \kappa_{\mathbf{x}}(\boldsymbol{\mu}_{<n-1}) + \underbrace{\kappa_{\mathbf{x}}(\mu_n | \boldsymbol{\mu}_{<n-1})}_{\geq 0}$ . It follows that  $\kappa_{\mathbf{x}}(\cdot)$  is monotone non-decreasing. ■

Next, we turn to a definition of sequence submodularity from Alaei et al. (2010). In contrast to Alaei et al.’s (2010) definition, we add the additional constraint that a merge-sequence function must take a valid merge sequence as an argument.

**Definition 3.3.** *A real-valued function  $f$  over valid merge sequences is **submodular** if, for all  $\boldsymbol{\mu}, \boldsymbol{\mu}' \in \mathcal{M}_{\Upsilon_{\Sigma}}$  such that  $\boldsymbol{\mu}' \preceq \boldsymbol{\mu}$ ,<sup>11</sup> and for all  $\nu \in \Upsilon_{\Sigma}$  such that both  $\boldsymbol{\mu}'\nu$  and  $\boldsymbol{\mu}\nu$  are valid, we have*

$$f(\nu | \boldsymbol{\mu}') \geq f(\nu | \boldsymbol{\mu}). \quad (6)$$

**Proposition 3.4.** *Let  $\kappa_{\mathbf{x}}$  be the compression utility function. Then, for a fixed  $\mathbf{x} \in \Sigma^*$ ,  $\kappa_{\mathbf{x}}(\cdot)$  is submodular (Definition 3.3) when the domain is restricted to the set of valid merges  $\mathcal{M}_{\Upsilon_{\Sigma}}$ .*

*Proof.* Let  $\boldsymbol{\mu}, \boldsymbol{\mu}' \in \mathcal{M}_{\Upsilon_{\Sigma}}$  such that  $\boldsymbol{\mu}' \preceq \boldsymbol{\mu}$ , and let  $\nu = [\nu', \nu'']$  be any merge such that  $\boldsymbol{\mu}\nu, \boldsymbol{\mu}'\nu \in \mathcal{M}_{\Upsilon_{\Sigma}}$ . First, notice that, once a merge  $\mu_n$  in a merge sequence  $\boldsymbol{\mu}$  is applied, the number of occurrences of  $\mu_n$  in  $\kappa_{\mathbf{x}}(\boldsymbol{\mu}_{\leq n})$  cannot be increased by any sequence of further applications, because all submerges of  $\mu_n$  where applied exhaustively (i.e., to all consecutive occurrences of their immediate submerges). Now, from  $\boldsymbol{\mu}'\nu \in \mathcal{M}_{\Upsilon_{\Sigma}}$ , it follows that both  $\nu'$  and  $\nu''$  are in  $\boldsymbol{\mu}'$ . Therefore, the number of occurrences  $\nu'$  and  $\nu''$ , and *a fortiori* of successive occurrences of them, cannot be greater in  $\kappa_{\mathbf{x}}(\boldsymbol{\mu})$  than in  $\kappa_{\mathbf{x}}(\boldsymbol{\mu}')$ , and hence  $\kappa_{\mathbf{x}}(\nu | \boldsymbol{\mu}) \leq \kappa_{\mathbf{x}}(\nu | \boldsymbol{\mu}')$ , which proves the submodularity of  $\kappa_{\mathbf{x}}$  over  $\mathcal{M}_{\Upsilon_{\Sigma}}$ . ■

In the context of compression, the submodularity property means, that the compression gain achieved after adding a specific merge to a merge sequence can never increase with merge sequence length. However, the requirement that the added merge does not create an invalid merge sequence is important. We highlight this importance in the following example.

<sup>11</sup>I.e., we have that  $\boldsymbol{\mu}'$  is a prefix of  $\boldsymbol{\mu}$ .

**Example 3.5.** Consider  $\Sigma = \{a, b, c, d, e\}$ , the string  $x = abcde$ , and the valid merge sequences  $\mu' = \langle [a, a] \rangle$  and  $\mu = \langle [a, a], [c, d] \rangle$ . Note that  $\mu' \preceq \mu$ . These merge sequences have compression utilities  $\kappa_x(\mu') = 6 - 5 = 1$  and  $\kappa_x(\mu) = 6 - 4 = 2$ , respectively. Next, consider the merge sequence  $\nu = \langle [b, [c, d]], [[b, [c, d]], e] \rangle$ . Now,  $\kappa_x(\nu \mid \mu') = 0$  and  $\kappa_x(\nu \mid \mu) = 2$ , which violates submodularity because  $\mu' \preceq \mu$ . What went wrong? The problem is that  $\mu\nu$  is not a valid merge sequence.

In order to formally prove our desired guarantee regarding the approximation bound of the greedy BPE algorithm, it is not enough that compression utility is sequence submodular over valid merge sequences. For this reason, we identified another property of the compression utility function that allows us to push through our result.

**Definition 3.6.** We define the following *partial order on merges*. For merges  $\mu, \mu' \in \Upsilon_\Sigma$ , we say  $\mu' \subset \mu$  iff  $\mu'$  is a submerge of  $\mu$ . The merge  $\mu'$  is a *submerge* of  $\mu = [\mu_1, \mu_2]$  iff:

- $\mu_1 = \mu'$ , or  $\mu_2 = \mu'$ , or
- $\mu' \subset \mu_1$ , or  $\mu' \subset \mu_2$ .

**Definition 3.7.** A real-valued function over valid merge sequences is *hierarchically sequence submodular* if, for every valid merge sequence of the form  $\mu'\nu'\mu\nu$  where  $\nu' \subset \nu$  according to the partial order given in Definition 3.6, we have that

$$f(\nu' \mid \mu') \geq f(\nu \mid \mu'\nu'\mu). \quad (7)$$

Note that hierarchical sequence submodularity is a different concept from function modularity, described in Definition 3.3. Indeed, in the case of functions over valid merge sequences, neither submodularity nor hierarchical sequence submodularity implies the other. To see this, note that roughly speaking, submodularity describes the difference in the value of a function when the same element is given as an argument, albeit conditioned on the presence of two different (but related) other arguments. However, if the same argument is considered in Eq. (7), we have

$$\kappa_x(\nu' \mid \mu') \geq \kappa_x(\nu' \mid \mu'\nu'\mu) = 0, \quad (8)$$

which is a trivial bound due to the non-negativity of  $\kappa_x(\cdot)$ . The naming is inspired by the fact we require the partial order over merges, which creates the hierarchy.

**Proposition 3.8.** Let  $\kappa_x$  be the compression utility function. Then, for a fixed  $x \in \Sigma^*$ ,  $\kappa_x(\cdot)$  is

hierarchically submodular (Definition 3.3) when the domain is restricted to the set of valid merges  $\mathcal{M}_{\Upsilon_\Sigma}$ .

*Proof.* Let  $x \in \Sigma^*$  be a string and  $\mu, \mu'$  be valid merge sequences. Furthermore, let  $\nu, \nu'$  be merges such that  $\nu' \subset \nu$  and  $\mu'\nu'\mu\nu$  is itself a valid merge sequence. Combinatorially,  $\kappa_x(\nu \mid \mu'\nu'\mu)$  is the number of replacements made in  $x$  by the single merge of  $\nu$ , after applying  $\mu'\nu'\mu$ . However, since  $\nu' \subset \nu$ , every new tree in  $x$  resulting from that by applying  $\nu$  must have  $\nu'$  as a descendant. Thus,  $\kappa_x(\nu' \mid \mu')$ , which is the number of new nodes in the forest created by applying  $\nu'$ , must be at least equal to  $\kappa_x(\nu \mid \mu'\nu'\mu)$ , if not greater. ■

Proposition 3.8 gives us a different notion of submodularity, which is important for the proof of the greedy BPE training guarantee. As an illustrative example of the proposition, we return to Fig. 1. In this case,  $\mu' = \langle [a, b] \rangle$ ,  $\nu' = [[a, b], a]$ ,  $\mu = \langle [c, b] \rangle$ ,  $\nu = [[[a, b], a], [c, b]]$ . Clearly,  $\nu' \subset \nu$  and  $\nu'$  appears twice, while  $\nu$  only once.

Finally, we adapt the definition of total backward curvature from (Zhang et al., 2015) to our needs. Intuitively, the total backward curvature is related to how much the utility of  $\mu$  can decrease if  $\nu$  is applied before, at the beginning.

**Definition 3.9.** The *total backward curvature* of the compression utility function  $\kappa$  with respect to an optimal merge sequence  $\mu^*$  is denoted with  $\sigma(\mu^*)$ :

$$\sigma(\mu^*) = \max_{\substack{\mu \in \Upsilon_\Sigma^* \\ |\mu| \leq M}} \left\{ 1 - \frac{\kappa(\mu\mu^*) - \kappa(\mu^*)}{\kappa(\mu)} \right\}. \quad (9)$$

## 3.2 The Greedy Algorithm for BPE

In words, the greedy algorithm proceeds as follows: For each of the  $M$  iterations, the algorithm chooses the next merge that is both valid and (locally) maximizes the objective in Eq. (4). We give pseudocode in Alg. 1. In practice, as shown in Code 1, this is done by choosing the merge that occurs most frequently (can be adjusted for pair overlaps). The main loop occurs  $M$  times. In the subsequent theorem we show the approximation bound for the greedy algorithm.

**Theorem 3.10.** The greedy algorithm for BPE training, i.e., for learning a length  $M$  merge sequence  $\mu^\dagger$ , is  $(\frac{1}{\sigma(\mu^*)}(1 - e^{-\sigma(\mu^*)}))$ -optimal: for every string  $x \in \Sigma^*$

$$\frac{\kappa_x(\mu^\dagger)}{\kappa_x(\mu^*)} \geq \frac{1}{\sigma(\mu^*)}(1 - e^{-\sigma(\mu^*)}) \quad (10)$$

Sequence	Pair frequencies
<b>Greedy</b>	
[a,b]a[a,b]baa	ab: 2, ba: 2, aa: 2, bb: 1
[[a,b],a][a,b]baa	[a,b]a: 1, [a,b]b: 1, ba: 1, aa: 1, [a,[a,b]]: 1
<b>Optimal</b>	
a[b,a]ab[b,a]a	ab: 2, ba: 2, aa: 2, bb: 1
a[[b,a],a]b[[b,a],a]	ab: 2, a[b,a]: 1, [b,a]a: 2, b[b,a]: 1

Example 2: In case of  $x = abaabbaa$  the greedy BPE yields a suboptimal compression utility (5 vs 4 subwords). Highlighted pairs show which one was chosen.

with respect to the optimal length  $M$  merge sequence  $\mu^*$ .

*Proof.* The proof is shown in App. A.  $\blacksquare$

### 3.3 Measuring Total Backward Curvature

We do not have a formal bound for  $\sigma(\mu^*)$  and estimate it by enumerating all strings of maximum length  $|x| \leq 15$  given a finite alphabet  $|\Sigma| = 5$  and maximum merge sequence size  $|\mu^*| < 5$ . The found maximum is  $\hat{\sigma}(\mu^*) = 2.5$ , from which follows an optimality bound of  $\approx 0.37$ . When we restrict our search to texts from a natural language (English), we obtain a slightly lower estimate  $\hat{\sigma}(\mu^*)_N = 2.0$  and hence optimality bound  $\approx 0.43$ . We leave the further study of the backward curvature constant to future work.

Notice that in the main proof of Theorem 3.10 in App. A, we used  $\sigma$  to bound only one particular type of sequence that becomes the prefix to  $\mu^*$ , namely  $\mu^\dagger$ . We may then check for prefixing only greedy sequences instead of taking the maximum across  $\mu \in \Upsilon_\Sigma^*$ ,  $|\mu| \leq M$  as in Definition 3.9:

$$\sigma'(\mu^*, \mu^\dagger) = \left\{ 1 - \frac{\kappa(\mu_{<M}^\dagger \mu^*) - \kappa(\mu^*)}{\kappa(\mu_{<M}^\dagger)} \right\} \quad (11)$$

This yields  $\hat{\sigma}'(\mu^*, \mu^\dagger) = 1.5$  and therefore the bound of  $\approx 0.52$ . More important than the particular bound value is that it is constant and that the BPE training algorithm can not be arbitrarily suboptimal with sequence length.

## 4 A Runtime Speed-up

We now introduce a speed-up of the greedy BPE algorithm. Assuming constant-time comparison of strings, finding the maximum pair count over the whole string is  $\mathcal{O}(N)$ , which is the same as applying one merge. Therefore, this implementation has a runtime complexity of  $\mathcal{O}(NM)$ . A large amount of time in the slow BPE implementation,

### Algorithm 2 Iterative Greedy BPE (faster).

**Inputs:** string  $x$ , merge count  $M$

**Output:** tokenized string  $x$ , merge sequence  $\mu$

```

1:  $\mu \leftarrow \langle \rangle$ 
2:  $x \leftarrow \text{LINKEDLIST}(x)$ 
3:  $h \leftarrow \text{MAXHEAP}(\text{PAIRS}(x))$ 
4: for  $i$  in  $0..M$  do
5:    $pos \leftarrow h.\text{TOP}$ 
6:   for  $(w_1, w_2)$  in  $pos$  do
7:      $h.\text{REMOVEPOSITION}(w_1.\text{prev}, w_1)$ 
8:      $h.\text{REMOVEPOSITION}(w_2, w_2.\text{next})$ 
9:      $w_1.\text{val} \leftarrow w_1.\text{val} + w_2.\text{val}$ 
10:     $w_1.\text{next} \leftarrow w_2.\text{next}$ 
11:     $w_2.\text{next}.\text{prev} \leftarrow w_1$ 
12:     $h.\text{ADDPPOSITION}(w_1.\text{prev}, w_1)$ 
13:     $h.\text{ADDPPOSITION}(w_1, w_1.\text{next})$ 
14:   end for
15:    $\mu \leftarrow \mu \circ \langle \mu \rangle$ 
16: end for
17: return  $x, \mu$ 

```

presented by Sennrich et al. (2016) and shown in Alg. 1, is spent on (1) recalculating the frequencies of pairs (Alg. 1, line 3) which are not affected by the most recent merge, and (2) scanning the whole string to apply a single merge (Alg. 1, line 4). To make this explicit, consider the following example.

**Example 4.1.** Consider  $x = abba(cddc)^n$  and merge  $[a, b]$  for  $n \geq 1$ . We can only apply the merge at the beginning of the string, which results in the forest  $[a, b]ba(cddc)^n$ . However, Alg. 2 still scans the entirety of the sequence to recalculate the pair frequencies of  $[c, d]$ ,  $[d, c]$  and  $[c, c]$ . This additional work is unnecessary.

Our idea to speed up Alg. 1 stems from the insight that we do not have to iterate over the entire sequence, an  $\mathcal{O}(N)$  operation, on each of the  $M$  iterations.<sup>12</sup> Indeed, on the  $t^{\text{th}}$  iteration, we show that one only has to do work proportional to the number of new nodes that are added to the forest (Alg. 2, line 6). To achieve this, we introduce a more efficient data structure for BPE.<sup>13</sup> Our first step is to treat the string as a linked list of subwords, initialized as a linked list of characters, that we destructively modify at each iteration. With each possible merge, we store a list of pointers where the merge operation could happen. The max heap is then sorted by the size of the sets. Lines 6 to 14 in Alg. 2 show the necessary operations needed to be performed on the linked list. Notably REMOVE-

<sup>12</sup> $N$  is the string length  $|x|$  and  $M$  the number of merges.

<sup>13</sup>Kudo and Richardson (2018) make a similar observation, however, we prove a tighter bound on the runtime.

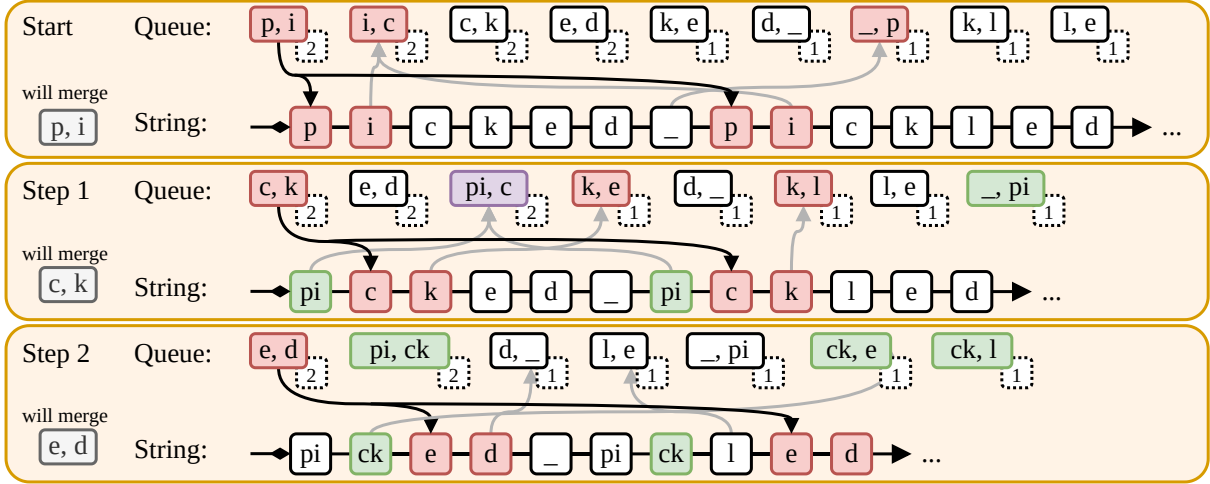


Figure 2: Visualization of linked list representation of the string and the associated priority queue (frequency values in dashed boxes) with merges. Nodes in red will be removed in the next step, nodes in green were added in contrast to the previous step and nodes in purple were just added but will be removed. Black lines from queue to the string show which nodes to merge. Grey lines show which pairs in the priority queue will have reduced frequencies.

POSITION removes the specific pair position from the set in the max heap and ADDPOSITION adds it.

See Fig. 2 for an illustration of applying a single merge in one place based on the introductory example in Example 1. The possible merge pairs are stored with a priority queue with their frequency as the sort key. During one operation, we need to remove the top merge pair and add counts for the newly created possible merge. The cost of one merge then becomes  $\mathcal{O}(R_t \log M)$  where  $R_t$  is the number of pairs in the string where the merge occurs and  $\log M$  the complexity of adding and updating frequency of a new merge pair. Note that it is not  $\log N$ , because we are keeping only top- $M$  possible pairs in the heap.

At first glance, this suggests the overall runtime of  $\mathcal{O}\left(\sum_{t=1}^M R_t \log M\right)$  with the worst case of the merge being applied along the whole string, therefore  $\mathcal{O}(MN \log M)$ .

**Theorem 4.2.** *Let  $N$  be the length of the string  $x \in \Sigma^*$  that is given as input. Then, Alg. 2 runs in  $\mathcal{O}(N \log M)$  time.*

*Proof.* Let  $R_t$  be the amount of work performed at each iteration modifying the data structure. We additionally do  $\mathcal{O}(\log M)$  work updating the priority queue on lines 6 to 14 in Alg. 2 since it has at most  $M$  elements. Thus, Alg. 2 clearly runs in  $\mathcal{O}\left(\sum_{t=1}^M R_t \log M\right)$ . We perform an amortized analysis. For this, we first make an observation about the upper bound on the number of merges and then show amortized analysis. However, for

a string  $x$  of length  $N$ , there are at most  $N - 1$  merges that can be applied to  $x$ . This implies that  $\sum_{t=1}^M R_t \leq N$ . Thus,  $\mathcal{O}\left(\sum_{t=1}^M R_t \log M\right) = \mathcal{O}(N \log M)$ , which proves the result. ■

## 5 An Exact Algorithm

In this section, we turn to developing an algorithm for exactly solving the BPE problem, i.e., Eq. (4). We change algorithmic paradigms and switch to memoization. While we are not able to devise a polynomial-time scheme, we are able to find an exact algorithm that is, in some cases, faster than the brute-force technique of enumerating all valid merge sequences. We first analyze the brute-force method of enumerating all valid merge sequences.

**Proposition 5.1.** *The set of valid merges of length  $M$  over a string  $x \in \Sigma^*$  is  $\mathcal{O}(\min(|\Sigma|^{2M}, N^M))$ .*

*Proof.* The proof can be found in App. A. ■

A simple direct enumeration of all possible merge sequences with the time complexity of one merge  $\mathcal{O}(NM)$  gives us a brute-force algorithm that runs in  $\mathcal{O}(NM \min(|\Sigma|^{2M}, N^M))$  time. The brute-force program explores all possible sequences of merges—including many that are redundant. For instance, both  $\langle [p, o], [h, a] \rangle$  and  $\langle [h, a], [p, o] \rangle$  induce the same partial bracketing when applied to another merge sequence, as in §2.2. Luckily, we are able to offer an exact characterization of when two merge sequences induce the same bracketing. To this end, we provide the following



definitions. We use the term **transposition** to refer to the swapping of items; i.e., a transposition  $(i, j)$  over a merge sequence  $\mu$  refers to the swapping of  $\mu_i$  and  $\mu_j$ .

**Definition 5.2.** A pair of merges  $\mu = [\mu_n, \mu_m]$  and  $\mu' = [\mu_{n'}, \mu_{m'}]$  **conflicts** if for a symbol  $a \in \Sigma$  and strings  $\mathbf{x}, \mathbf{x}' \in \Sigma^*$ , the yield of  $[\mu_n, \mu_m]$  is  $\mathbf{x}a$  and  $[\mu'_{n'}, \mu'_{m'}]$  is  $a\mathbf{x}'$ .

**Definition 5.3.** A transposition  $(i, j)$  is **safe** if and only if, for all  $k < j$ ,  $\mu_k$  does not conflict with  $\mu_j$  and, for all  $k > i$ ,  $\mu_k$  does not conflict with  $\mu_i$ . A permutation  $\pi = \langle \rho_1 \rho_2 \dots \rho_n \rangle$ , decomposed into transpositions, that maps one valid merge sequence  $\mu$  to another valid merge sequence  $\pi(\mu) = \mu'$  is **safe** if and only if all transpositions are safe.

Informally, Definition 5.3 says that for a permutation to produce a valid merge sequence, there should be no conflicts between the swapped merges and all merges in between. For example, given the merge sequence  $\mu = \langle [a, b], [d, d], [c, a] \rangle$ , the permutation  $\pi = \langle (1, 3) \rangle$  would not be safe.

The reason for this definition is that safe permutations characterize when two merge sequences always give the same results. Indeed, for  $\mathbf{x} = ddabcacab$ , applying the first merge sequence:  $\text{APPLY}_\mu(\mathbf{x}) = [d, d][a, b][c, a][c, a]b$ . In contrast, applying the permuted one gives an alternative outcome:  $\text{APPLY}_{\pi(\mu)}(\mathbf{x}) = [d, d][a, b][c, a]c[a, b]$ .

**Definition 5.4.** Two merge sequences  $\mu$  and  $\mu'$  are **equivalent** if and only if, for all  $\mathbf{x} \in \Sigma^*$ ,  $\text{APPLY}_\mu(\mathbf{x}) = \text{APPLY}_{\mu'}(\mathbf{x})$ . Symbolically, we write  $\mu \equiv \mu'$  if  $\mu$  and  $\mu'$  are equivalent.

**Proposition 5.5.** Two valid merge sequences  $\mu, \mu' \in \mathcal{M}_{\Upsilon_\Sigma}$  are equivalent, i.e.,  $\mu \equiv \mu'$ , if and only if there exists a safe permutation  $\pi$  such that  $\pi(\mu) = \mu'$ .

*Proof.* The proof can be found in App. A. ■

Following the previous example, it is easy to verify that  $\langle [a, b], [d, d], [c, a] \rangle \equiv \langle [a, b], [c, a], [d, d] \rangle$ . In contrast to synthetic examples with a constrained alphabet of, e.g.,  $\{a, b, c\}$ , far fewer merge conflicts arise in natural language. We can leverage this to develop a faster algorithm that only explores paths that are not equivalent to each other. We first define the concept of partial ordering between merges.

**Definition 5.6.** The **merge partial ordering**  $\mu' > \mu''$  is defined as  $\neg \text{conflicts}(\mu', \mu'') \wedge \neg(|\text{YIELD}(\mu')| < |\text{YIELD}(\mu'')|) \wedge \neg(\text{YIELD}(\mu') <_L \text{YIELD}(\mu''))$  where  $>_L$  is lexicographical ordering.

All valid merge sequences are equivalent to some merge sequence which is partially ordered using  $>$  so that no neighbouring elements violate this partial ordering. The brute-force algorithm works as depth-first search through an acyclic graph: each state corresponds to a unique sequence of merges and each transition corresponds to appending a merge to the end of the current state's merges. For the improved version, we make sure that only sequences which are ordered using  $>$  are searched and the rest are pruned. The pseudocode for the program is shown in Alg. 3. Even though the runtime is still prohibitively slow for application, Fig. 3 demonstrates how much speed is gained over the brute-force version which explores all states.

**Algorithm 3** Exact BPE with memoization guard. Removing segments marked with   would result in the brute-force version.

**Inputs:** string  $\mathbf{x}$ , merge count  $M$

**Output:** tokenized string  $\mathbf{x}$ , merge sequence  $\mu$

---

```

1:  $q \leftarrow \text{STACK}()$ 
2:  $q.\text{PUSH}(\langle \rangle, \mathbf{x})$ 
3:  $\mu^*, \mathbf{x}^* \leftarrow \langle \rangle, \mathbf{x}$ 
4: while  $|q| \neq 0$  do
5:    $\mu, \mathbf{x} \leftarrow q.\text{POP}()$ 
6:   if  $|\mu| = M$  then continue
7:   for  $\mu \in \text{PAIRS}(\mathbf{x})$  do
8:       if  $|\mu| = 0 \vee \mu > \mu_{-1}$ 
9:        $\mathbf{x}' \leftarrow \text{SINGLEAPPLY}(\mathbf{x}, \mu)$ 
10:       $\mu' \leftarrow \mu \circ \mu$ 
11:      if  $|\mathbf{x}'| < |\mathbf{x}^*|$ 
12:         $\mu^*, \mathbf{x}^* \leftarrow \mu', \mathbf{x}'$ 
13:      end if
14:       $q.\text{PUSH}(\mu', \mathbf{x}')$ 
15:      end if
16:   end for
17: end while
18: return  $\mu^*, \mathbf{x}^*$ 

```

---

## 6 Conclusion

In this paper, we developed the formalisms surrounding the training task of BPE, a very popular tokenization algorithm in NLP. This allowed us to prove a lower bound on the compression utility by greedy BPE as  $1 - e^{-\sigma(\mu^*)}$ . We further analyzed the runtime of the naive and faster greedy BPE algorithms and provided a speedup for finding an optimal BPE merge sequence. Future works should focus on providing either formal guarantees for  $\sigma(\mu^*)$  or studying  $\sigma(\mu^*)'$  across natural languages.

## 7 Limitations

Our work has focused strongly on the formal aspects of BPE. NLP practitioners should not be dissuaded from using BPE for subword tokenization, despite our presentation of examples where greedy BPE fails. Indeed, in contrast to synthetic examples on toy alphabet, on real data we made an observation that greedy BPE may be close to optimal.

## Acknowledgements

We would like to thank Andreas Krause and Giorgio Satta for discussing the proof of Theorem 3.10. Clara Meister was supported by the Google PhD Fellowship. Juan Luis Gastaldi has received funding from the European Union’s Horizon 2020 research and innovation programme under grant agreement No 839730.

## References

- Saeed Alaei, Ali Makhdoumi, and Azarakhsh Malekian. 2010. [Maximizing sequence-submodular functions and its application to online advertising](#). *arXiv preprint arXiv:1009.4153*.
- Jeff Bilmes. 2022. [Submodularity in machine learning and artificial intelligence](#). *arXiv preprint arXiv:2202.00132*.
- Kaj Bostrom and Greg Durrett. 2020. [Byte pair encoding is suboptimal for language model pretraining](#). In *Findings of the Association for Computational Linguistics: EMNLP 2020*, pages 4617–4624.
- Tom Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared Kaplan, Prafulla Dhariwal, Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, Sandhini Agarwal, Ariel Herbert-Voss, Gretchen Krueger, Tom Henighan, Rewon Child, Aditya Ramesh, Daniel M. Ziegler, Jeffrey Wu, Clemens Winter, Chris Hesse, Mark Chen, Eric Sigler, Mateusz Litwin, Scott Gray, Benjamin Chess, Jack Clark, Christopher Berner, Sam McCandlish, Alec Radford, Ilya Sutskever, and Dario Amodei. 2020. [Language models are few-shot learners](#). In *Advances in Neural Information Processing Systems*, volume 33, pages 1877–1901.
- Thomas M. Cover and Joy A. Thomas. 2006. *Elements of Information Theory*, 2 edition. Wiley-Interscience.
- Shuoyang Ding, Adithya Renduchintala, and Kevin Duh. 2019. [A call for prudent choice of subword merge operations in neural machine translation](#). In *Proceedings of Machine Translation Summit XVII: Research Track*, pages 204–213.
- Ahmed El-Kishky, Vishrav Chaudhary, Francisco Guzmán, and Philipp Koehn. 2020. [CCAligned: A massive collection of cross-lingual web-document pairs](#). In *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP)*, pages 5960–5969.
- Philip Gage. 1994. [A new algorithm for data compression](#). *The C Users Journal*, 12(2):23–38.
- Andreas Krause and Daniel Golovin. 2014. [Submodular function maximization](#). In *Tractability: Practical Approaches to Hard Problems*. Cambridge University Press.
- 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.
- László Lovász. 1983. [Submodular functions and convexity](#). In *Mathematical programming the state of the art*, pages 235–257. Springer.
- Azarakhsh Malekian. 2009. *Combinatorial Problems in Online Advertising*. Ph.D. thesis, University of Maryland, College Park.
- Alec Radford, Jeffrey Wu, Rewon Child, David Luan, Dario Amodei, and Ilya Sutskever. 2019. [Language models are unsupervised multitask learners](#).
- Teven Le Scao, Angela Fan, Christopher Akiki, Ellie Pavlick, Suzana Ilić, Daniel Hesslow, Roman Castagné, Alexandra Sasha Luccioni, François Yvon, Matthias Gallé, et al. 2022. [BLOOM: A 176B-parameter open-access multilingual language model](#). *arXiv preprint arXiv:2211.05100*.
- 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.
- Terry A. Welch. 1984. [A technique for high-performance data compression](#). *Computer*, 17(06):8–19.
- Yizhe Zhang, Siqi Sun, Michel Galley, Yen-Chun Chen, Chris Brockett, Xiang Gao, Jianfeng Gao, Jingjing Liu, and William B Dolan. 2020. [DIALOGPT: Large-scale generative pre-training for conversational response generation](#). In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics: System Demonstrations*, pages 270–278.
- Zhenliang Zhang, Edwin KP Chong, Ali Pezeshki, and William Moran. 2015. [String submodular functions with curvature constraints](#). *IEEE Transactions on Automatic Control*, 61(3):601–616.
- 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 Proofs

Our proof of approximate optimality is based on the proof of greedily sequence maximizing submodular functions by Alaei et al. (2010); Zhang et al. (2015). However, we leverage a problem-specific property, which we dub hierarchical submodularity. We restate the definition here for ease.

**Definition 3.7.** A real-valued function over valid merge sequences is **hierarchically sequence submodular** if, for every valid merge sequence of the form  $\mu' \nu' \mu \nu$  where  $\nu' \subset \nu$  according to the partial order given in Definition 3.6, we have that

$$f(\nu' \mid \mu') \geq f(\nu \mid \mu' \nu' \mu). \quad (7)$$

**Lemma A.1.** Let  $\mu', \mu \in \mathcal{M}_{\Upsilon_\Sigma}$  be valid merge sequences. Then, there exists a merge  $\nu$  in  $\mu$  such that  $\mu' \nu$  is a valid merge sequence and  $\kappa_{\mathbf{x}}(\nu \mid \mu') \geq \frac{\kappa_{\mathbf{x}}(\mu \mid \mu')}{|\mu|}$ . In words, the compression gain of some element in  $\mu$  with respect to  $\mu'$  is greater or equal to the average compression gain per element of  $\mu$  with respect to  $\mu'$

*Proof.* Let us choose one of the possible maximums,  $t = \operatorname{argmax}_{1 \leq t' \leq |\mu|} \kappa_{\mathbf{x}}(\mu_{t'} \mid \mu' \mu_{<t'})$ . Because we are taking the maximum, which is always equal to or greater than the average,<sup>14</sup> then  $\kappa_{\mathbf{x}}(\mu_t \mid \mu' \mu_{<t}) \geq \frac{1}{|\mu|} \sum_{t'=1}^{|\mu|} \kappa_{\mathbf{x}}(\mu_{t'} \mid \mu' \mu_{<t'})$ . Then, we have that either:

- $\mu \mu_t \in \mathcal{M}_{\Upsilon_\Sigma}$ , in which case the result follows by submodularity, or
- $\mu \mu_t \notin \mathcal{M}_{\Upsilon_\Sigma}$ , in which case there exists a  $\mu_{t'}$  such that:
  - $\mu_{t'} \subset \mu_t$
  - $\mu' \mu_{t'} \in \mathcal{M}_{\Upsilon_\Sigma}$
  - $\mu_{t'} \text{ in } \mu$
  - $\kappa_{\mathbf{x}}(\mu_t \mid \mu' \mu_{<t}) \leq \kappa_{\mathbf{x}}(\mu_{t'} \mid \mu' \mu_{<t'}) \leq \kappa_{\mathbf{x}}(\mu_{t'} \mid \mu')$

In particular, all trivial submerges of  $\mu_t$  (i.e., all submerges of  $\mu_t$  whose constituents are in  $\Sigma$ ) fulfill all four conditions: the first one by definition, the second by definition of  $\mathcal{M}_{\Upsilon_\Sigma}$ , the third because  $\mu \in \mathcal{M}_{\Upsilon_\Sigma}$ , and the fourth by hierarchical submodularity (first inequality) and by submodularity (second inequality). ■

We now proceed with the proof of approximate optimality of the greedy BPE merge sequence.

**Theorem 3.10.** The greedy algorithm for BPE training, i.e., for learning a length  $M$  merge sequence  $\mu^\dagger$ , is  $(\frac{1}{\sigma(\mu^*)}(1 - e^{-\sigma(\mu^*)}))$ -optimal: for every string  $\mathbf{x} \in \Sigma^*$

$$\frac{\kappa_{\mathbf{x}}(\mu^\dagger)}{\kappa_{\mathbf{x}}(\mu^*)} \geq \frac{1}{\sigma(\mu^*)}(1 - e^{-\sigma(\mu^*)}) \quad (10)$$

with respect to the optimal length  $M$  merge sequence  $\mu^*$ .

*Proof.* We make use of the sequence  $\mu_{<M}^\dagger$  (rather than  $\mu^\dagger$ ) for reasons that will subsequently become clear. From Lemma A.1, we know that we can find  $\mu_j^*$  such that  $\mu_{<M}^\dagger \mu_j^*$  is a valid merge sequence and

$$\kappa(\mu_j^* \mid \mu_{<M}^\dagger) \geq \frac{1}{M} \kappa(\mu^* \mid \mu_{<M}^\dagger) \quad (12)$$

From the greedy property of  $\mu^\dagger$ , we know:

$$\kappa(\mu_M^\dagger \mid \mu_{<M}^\dagger) \geq \kappa(\mu_j^* \mid \mu_{<M}^\dagger) \quad (13)$$

$$\kappa(\mu_M^\dagger \mid \mu_{<M}^\dagger) \geq \frac{1}{M} \kappa(\mu^* \mid \mu_{<M}^\dagger) \quad (\text{from Eq. 12}) \quad (14)$$

$$\kappa(\mu_{<M}^\dagger \mu_M^\dagger) - \kappa(\mu_{<M}^\dagger) \geq \frac{1}{M} (\kappa(\mu_{<M}^\dagger \mu^*) - \kappa(\mu_{<M}^\dagger)) \quad (\text{definition expansion}) \quad (15)$$

<sup>14</sup>Proof of this algebraic statement is omitted for brevity.



Now from backward curvature (Definition 3.9) and by substituting  $\mu_{<M}^\dagger$  for the prefix sequence:

$$\sigma(\mu^*) \geq 1 - \frac{\kappa(\mu_{<M}^\dagger \mu^*) - \kappa(\mu^*)}{\kappa(\mu_{<M}^\dagger)} \quad (16)$$

$$\sigma(\mu^*) \kappa(\mu_{<M}^\dagger) \geq \kappa(\mu_{<M}^\dagger) - \kappa(\mu_{<M}^\dagger \mu^*) + \kappa(\mu^*) \quad (17)$$

$$\kappa(\mu_{<M}^\dagger \mu^*) - \kappa(\mu_{<M}^\dagger) \geq \kappa(\mu^*) - \sigma(\mu^*) \kappa(\mu_{<M}^\dagger) \quad (18)$$

Applying this result to the right-hand side of Eq. (15), we obtain the following:

$$\kappa(\mu_{<M}^\dagger \mu_M^\dagger) - \kappa(\mu_{<M}^\dagger) \geq \frac{1}{M} (\kappa(\mu^*) - \sigma(\mu^*) \kappa(\mu_{<M}^\dagger)) \quad (\text{total backward curvature}) \quad (19)$$

$$\kappa(\mu^\dagger) - \kappa(\mu_{<M}^\dagger) \geq \frac{1}{M} (\kappa(\mu^*) - \sigma(\mu^*) \kappa(\mu_{<M}^\dagger)) \quad (\text{definition}) \quad (20)$$

$$\kappa(\mu^\dagger) \geq \frac{1}{M} (\kappa(\mu^*) - \sigma(\mu^*) \kappa(\mu_{<M}^\dagger)) + \kappa(\mu_{<M}^\dagger) \quad (\text{total backward curvature}) \quad (21)$$

$$\geq \frac{1}{M} \kappa(\mu^*) + \left(1 - \frac{\sigma(\mu^*)}{M}\right) \kappa(\mu_{<M}^\dagger) \quad (\text{algebraic manipulation}) \quad (22)$$

$$\geq \frac{1}{M} \kappa(\mu^*) \sum_{i=0}^{M-1} \left(1 - \frac{\sigma(\mu^*)}{M}\right)^i \quad (\text{recursive substitution of } \kappa(\mu_{<i}^\dagger)) \quad (23)$$

$$= \frac{1}{\sigma(\mu^*)} \left(1 - \left(1 - \frac{\sigma(\mu^*)}{M}\right)^M\right) \kappa(\mu^*) \quad (\text{geometric sum}) \quad (24)$$

$$= \frac{1}{\sigma(\mu^*)} \left(1 - \left(1 - \frac{\sigma(\mu^*)}{M}\right)^{\frac{M}{\sigma(\mu^*)}}\right)^{\sigma(\mu^*)} \kappa(\mu^*) \quad (\text{preparation}) \quad (25)$$

We substitute  $x = \frac{M}{\sigma(\mu^*)}$  in the inequality. From  $x > 0 \Rightarrow \left(1 - \frac{1}{x}\right)^x \leq \frac{1}{e}$ , we obtain and arrive at

$$\kappa(\mu^\dagger) \geq \frac{1}{\sigma(\mu^*)} \left(1 - e^{-\sigma(\mu^*)}\right) \quad (26)$$

■

**Proposition 5.5.** *Two valid merge sequences  $\mu, \mu' \in \mathcal{M}_{\Upsilon_\Sigma}$  are equivalent, i.e.,  $\mu \equiv \mu'$ , if and only if there exists a safe permutation  $\pi$  such that  $\pi(\mu) = \mu'$ .*

*Proof. ( $\Rightarrow$ ):* We prove the first implication through contrapositive, i.e., we show that if there does *not* exist such a safe permutation  $\pi$ , then the merge sequences are *not* equivalent. By supposition, all non-safe permutations mapping  $\mu$  to  $\mu'$  either have a conflict or do not preserve validity. We handle each case separately.

- **Case 1:** Suppose that the permutation  $\pi$  re-orders two conflicting merges  $\mu$  and  $\mu'$ . By the definition of a conflict,  $\mu$  has yield  $xa$  and  $\mu'$  has yield  $ax'$  for  $a \in \Sigma$  and  $x, x' \in \Sigma^*$ . Now, note the bracketing string  $xaax'$  will be different under the original and permuted merge sequence.
- **Case 2:** Suppose that the permutation  $\pi$  does not preserve validity. Then, there exists a merge  $\mu = (\mu', \mu'')$  such that either  $\mu'$  or  $\mu''$  occurs *after*  $\mu$  in the merge sequence. This also results in a different bracketing.

*( $\Leftarrow$ ):* Next, we want to show the converse, i.e., for any safe permutation  $\pi$ , we have  $\mu \equiv \pi(\mu)$ . Let  $\mu = \langle \mu_1, \dots, \mu_N \rangle$  be a merge sequence of length  $N$ , and let  $\pi$  be a safe permutation. We proceed by induction on the  $n$ .

- **Base Case:** Since  $\pi$  is safe, then for  $[a, b] = \pi(\boldsymbol{\mu})_1$ ,  $a$  and  $b$  are necessarily characters in  $\Sigma$ .
- **Inductive Step:** Suppose for  $k = n - 1$ ,  $\pi(\boldsymbol{\mu})_{\leq k}$  applies merges which are applied by  $\boldsymbol{\mu}$ . We then show  $\pi(\boldsymbol{\mu})_n$  also applies the same merges as  $\boldsymbol{\mu}$ . Consider  $\pi(\boldsymbol{\mu})_n = (\mu_m, \mu_{m'})$ ; since  $\pi$  is safe, both  $\mu_m$  and  $\mu_{m'}$  already exist in  $\text{APPLY}_{\boldsymbol{\mu}_{\leq n}}(\boldsymbol{x})$ . Moreover, since there are no conflicts, applying  $\pi(\boldsymbol{\mu})_n$  results in the same encoded sequence. ■

**Proposition 5.1.** *The set of valid merges of length  $M$  over a string  $\boldsymbol{x} \in \Sigma^*$  is  $\mathcal{O}(\min(|\Sigma|^{2M}, N^M))$ .*

*Proof.* On one hand, we note that we have an upper bound of  $N - 1$  possible merges that can occupy the first element of the sequence, assuming every symbol in  $\boldsymbol{x}$  is distinct. Next, we have  $N - 2$  possible merges that can occupy the second element of the sequence, again, assuming every symbol in  $\boldsymbol{x}$  is distinct. Continuing this pattern, we arrive at a simple upper bound on the number of merges  $\prod_{m=0}^{M-1} (N - 1 - m)$ . This quantity is recognizable as a falling factorial, which gives us the closed form  $\frac{(N-1)!}{(N-M-2)!}$ . This can be trivially bounded by  $N^M$ . However, on the other hand, we know a valid merge sequence can produce merges with a yield up to length  $M$ , and there are  $\binom{\Sigma^{\leq M}}{M}$  unique sequences. We can upper-bound the number of valid merge sequences by the total number of all possible merge sequences, of which there are  $M!$ . The size of  $\Sigma^{\leq M}$  is the sum  $|\Sigma|^1 + |\Sigma|^2 + \dots + |\Sigma|^M$  which is less than  $M|\Sigma|^M$ . Again, with  $M!$ , this leads to the falling factorial  $\frac{(M|\Sigma|^M)!}{(M|\Sigma|^M - M)!}$  which we can upper bound by  $(M|\Sigma|^M)^M$  which is in  $\mathcal{O}(|\Sigma|^{2M})$ . Taking the min of these two upper bounds gives us the overall upper bound. ■

## B BPE Modifications

In this section, we describe multiple modifications to the greedy BPE algorithm which speed up the runtime. We do not address popular heuristic modifications such as lowercasing the text or adding 20% of the most frequent words to the subword dictionary.

### B.1 (Not) Merging Space

Currently, spaces are treated as any other characters and are allowed to be part of merges. Therefore in the string “not\_that\_they\_watch\_the\_watch” the first merge is  $[_ , t]$  and the string looks as “not[\_ , t]hat[\_ , t]hey watch[\_ , t]he watch”. The next merge may be across tokens:  $[t, [_ , t]]$ . This is not desirable if we want only want to split tokens into subwords (i.e. use merges that do not contain spaces).

Furthermore, in §3 we are duplicating work by computing pair frequencies and merges multiple times across the same tokens that occur multiple times (see previous string example). In practice (Tab. 1), only 1.5% of all tokens are unique. We may then speed up our computation by considering only unique tokens. Therefore, the new runtime complexity is  $\mathcal{O}(V \cdot |\boldsymbol{x}_u|)$  where  $\boldsymbol{x}_u = \{t \mid \text{token } t \in \boldsymbol{x}\}$  which is  $\frac{|\boldsymbol{x}|}{|\boldsymbol{x}_u|} \times$  faster.

### B.2 Non-iterative BPE

A popular implementation of BPE-like algorithm in Python<sup>15</sup> uses a different speed-up mechanism to avoid  $\mathcal{O}(NV)$  runtime. This is done by:

- (1) collecting all possible merges observed in the data up until some maximum yield size which determines the maximum subword size, such as 5 and
- (2) taking top- $M$  frequent pairs as part of the subword dictionary.

Note that because of hierarchical submodularity (Definition 3.7), this will produce valid merges. This is because if  $\mu = [\mu', \mu'']$  is chosen, so must  $\mu'$  and  $\mu''$  because they have at least the same frequency as  $\mu$ . For example, for  $abcabcd$ , and maximum yield width 3, the merges would be  $[a, b]$ ,  $[[a, b], c]$ ,  $[b, c]$ ,  $[a, [b, c]]$ ,  $\dots$ . The runtime of this is  $\mathcal{O}(|\boldsymbol{x}| \log M)$  because we are scanning the whole string and at each point are modifying maximum heap.

<sup>15</sup>[pypi.org/project/bpe](https://pypi.org/project/bpe)

However, it is easy to see that this approximation algorithm is not bounded. For a constant maximum yield width of  $w$ , consider  $x = a^{wn}$  and  $V = w + k$ . The shortest possible output of this algorithm will be  $\mu^n$ . However, an optimal merge sequence can perform additional merge sequences, therefore producing  $\nu^{\frac{n}{2^k}}$ . The compressions are  $wn - n$  and  $wn - \frac{n}{2^k}$  and the ratio  $\frac{wn-n}{wn-\frac{n}{2^k}}$  with lower bound of 0 as supremum. This means that we can construct adversarial example for which the compression given by this algorithm is arbitrarily suboptimal.

### C Additional Experimental Details and Results

Sentence count (train)	13M+13M
Sentence count (dev & test)	1M+1M
Total words	324M
Unique words	5M
Average sentence length (words)	12

Table 1: Overview of the used portion of the English-German CommonCrawl dataset (El-Kishky et al., 2020).

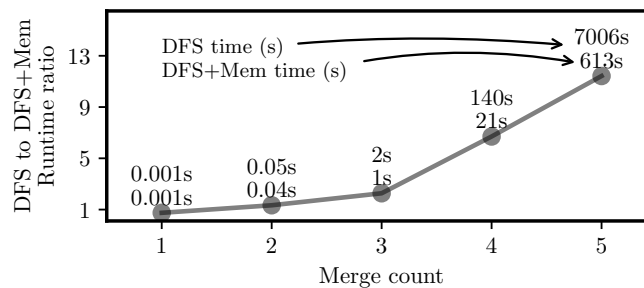


Figure 3: Comparison of runtimes for brute-force DFS and DFS with memoization. Values above 1 correspond to DFS+memoization being  $\times$  faster than DFS. Points show average<sup>16</sup> of runs on 5 different input strings (each 2 randomly sampled English sentences of 64 characters).

<sup>16</sup>Time measured on desktop AMD Ryzen 9 5900X.

```

1 from collections import Counter, defaultdict
2 from typing import Union, Tuple, List
3
4 def fixed_pair_freqs(xs: Union[str, List]):
5     pairs = defaultdict(int)
6     prev_pair = None
7     for (x, y) in zip(xs, xs[1:]):
8         # increment only if the prev suffix does not match prefix
9         # otherwise wrong estimate on `aaa`
10        if (x,y) != prev_pair:
11            pairs[x, y] += 1
12            prev_pair = (x, y)
13        else:
14            # make sure to clear it so that `aaaa` is counted twice
15            prev_pair = None
16
17        pairs = list(pairs.items())
18        pairs.sort(key=lambda x: x[1], reverse=True)
19        return pairs
20
21 def bpe(xs: Union[str, List], V: int):
22     for _ in range(V):
23         top_pair = fixed_pair_freqs(xs)[0]
24         xs = merge(list(xs), top_pair)
25     return xs
26
27 def merge(xs: List, pair: Tuple):
28     ys = []
29     while xs:
30         if tuple(xs[:2]) == pair:
31             ys.append(pair)
32             xs = xs[2:]
33         else:
34             ys.append(xs.pop(0))
35     return ys

```

Code 2: An implementation of [Sennrich et al.'s \(2016\)](#) greedy algorithm for BPE in Python with overlap-adjusted pair counts.