# **Tokenisation is NP-Complete**

## Anonymous ACL submission

#### Abstract

In this work, we prove the NP-completeness of two variants of tokenisation, defined as the problem of compressing a dataset to at most  $\delta$ symbols by either finding a vocabulary directly (*direct* tokenisation), or selecting a sequence of merge operations (*bottom-up* tokenisation).

#### 1 Introduction

011

017

019

024

027

Tokenisation is at the heart of natural language processing (NLP) being the first step required to use a language model (LM). Given a string of characters c, a tokeniser converts it into a string of **subwords** s. Language models are then trained to estimate distributions over subword strings—never seeing the original character strings. Despite its prominent role, however, much remains unknown about tokenisation. We still do not know, for instance, what makes a good tokeniser (Gowda and May, 2020; Cognetta et al., 2024): which characteristics should its produced subwords s have to be a good starting point for language modelling? If we knew this, then we could define an **objective function** with which we could evaluate tokenisers.

Another open question is how to—given such an objective function—efficiently find a tokeniser which maximises it. Byte pair encoding (BPE; Gage, 1994; Sennrich et al., 2016), for instance, is a greedy solution to find a tokeniser which maximises a text's compression. UnigramLM (Kudo, 2018) is a heuristic method to find a tokeniser that maximises its tokenised text's unigram logprobability. Both these methods, however, are approximate: they do not necessarily find an optimal tokeniser according to their objective function. This raises the question of whether finding such optimal tokenisers efficiently is even possible.

In this paper, we answer this question (at least partially), proving the NP-completeness of several variants of this tokenisation problem. Specifically, we focus on finding tokenisers that maximise the **compression** of a text.<sup>1</sup> Given this objective, we then define the **tokenisation problem** as the task of finding a tokeniser which compresses a dataset to at most  $\delta$  symbols. Notably, prior work imposes different constraints on how tokenisers are defined; here we consider two variants. In **direct tokenisation**, the desired compression must be reached by choosing a vocabulary (i.e., a set of subwords) which is directly used to represent the text. In **bottom-up tokenisation**, the desired compression must be reached by finding a sequence of merge operations instead, which we apply to the input text.

041

042

043

044

045

047

049

052

054

056

058

060

061

062

063

064

065

066

067

068

069

070

071

072

We prove the NP-hardness of both of these tokenisation problems (as well as of some variants thereof) by reducing from the **max 2-satisfiability** problem.<sup>2</sup> Practically speaking, our results imply that we are unlikely to discover an efficient algorithm for the problem of finding optimal tokenisers, and that we should focus on approximate algorithms (such as BPE or UnigramLM) instead.

## 2 How to Choose a Tokeniser?

Given a subword-level language model, we can extract word-level (Pimentel and Meister, 2024; Oh and Schuler, 2024) or character-level (Phan et al., 2024; Giulianelli et al., 2024) distributions from it. Further, regardless of which tokeniser is used, a sufficiently expressive language model should be able to represent the exact distributions over characters or words that we are interested in. In theory, thus, a researcher's choice of tokeniser should not influence their language model's quality.

In practice, however, a bad choice of tokeniser

<sup>&</sup>lt;sup>1</sup>The compression achieved by a tokeniser correlates with downstream language modelling performance (Gallé, 2019; Zouhar et al., 2023a) and computational efficiency.

<sup>&</sup>lt;sup>2</sup>We note two related concurrent works. Kozma and Voderholzer (2024) also prove the NP-completeness of bottom-up tokenisation; in fact, they prove something stronger: its APXhardness. Lim et al. (2025) prove the NP-completeness of a restricted variant of direct tokenisation, in which a set of candidate tokens is previously specified.

073can have undesirable effects on downstream appli-<br/>cations. For instance, performing standard arith-<br/>metic tasks (e.g., 317 + 421) can be difficult even<br/>for large models (Nogueira et al., 2021; Muffo<br/>et al., 2022) due to the arbitrary splitting of num-<br/>bers into subwords. Indeed, simple changes in how<br/>numbers are tokenised can improve performance<br/>in such tasks (Singh and Strouse, 2024). Similar<br/>issues arise when prompting LMs to count letters in<br/>words, where even advanced models such as GPT-4<br/>infamously cannot correctly count the number of<br/>occurrences of the letter r in the word strawberry.

This raises the question of how to select a good tokeniser. Ideally, we would choose the tokeniser which maximises downstream language modelling performance. Unfortunately, we do not know how to measure such performance without fully training a model, making its direct maximisation computationally infeasible. Rather, we thus optimise proxy objectives—assumed to correlate with downstream performance. Among these are unigram logprobability (Kudo, 2018), Rényi efficiency (Zouhar et al., 2023a), and compression (Gallé, 2019).

We focus on compression in this paper. Denoting our tokenisation's **objective function** as  $\mathfrak{G}$ , we write this objective as:  $\mathfrak{G}(\mathbf{s}) = -|\mathbf{s}|$ . Improved compression leads to: (i) more efficient training and inference, due to shortened inputs;<sup>3</sup> (ii) improved downstream performance, at least to a certain extent (Gallé, 2019; Rust et al., 2021; Zouhar et al., 2023a; Goldman et al., 2024);<sup>4</sup> and (iii) fairer multilingual treatment—assuming similar compression across languages—given models' limited context lengths and the per-token costs to use proprietary models (Petrov et al., 2023; Ahia et al., 2023).

#### Our Notation's Colour-coding

100

101

103

104

105

107

108

109

110

111

- Green for raw data (i.e., characters  $\mathbf{c} \in \Sigma^*$ );
- Purple for tokeniser-specific data (i.e., subwords s ∈ S\* and merges m ∈ M\*);
- Blue for functions (e.g., tok).

#### **3** Defining a Tokeniser

A tokeniser can be defined as a 3-tuple  $\langle S, tok, detok \rangle$ , composed of a vocabulary, a to-

kenisation and a detokenisation function. Before defining these terms, however, we require some notation. Let  $\mathbf{c} = c_1 c_2 \cdots c_{|\mathbf{c}|} \in \Sigma^*$  be a **character-string**, i.e., a sequence of characters cfrom alphabet  $\Sigma$ . Further, let  $\mathcal{D} = \{\mathbf{c}_n\}_{n=1}^N$  be a dataset of character-strings.<sup>5</sup> A subword  $s \in S$ represents a non-empty character-string  $\mathbf{c}$  (where sequence  $\mathbf{c}$  can have length one). Finally, let  $\mathbf{s} = \langle s_1, s_2, \cdots, s_{|\mathbf{s}|} \rangle \in S^*$  be a **subword-string**. Just like a single subword, a subword-string  $\mathbf{s} \in S^*$ represents a character-string via the concatenation of its subwords' characters: 112

113

114

115

116

117

118

119

120

121

122

123

124

125

126

127

128

129

130

131

132

133

134

135

136

137

138

140

141

142

143

144

145

146

147

148

149

150

151

$$\operatorname{concat}(\mathbf{s}) = s_1 \circ s_2 \circ \dots \circ s_{|\mathbf{s}|} \tag{1}$$

and we say that a pair of character and subword strings are equivalent if:

$$\mathbf{c} \stackrel{\circ}{=} \mathbf{s} \iff \mathbf{c} = \texttt{concat}(\mathbf{s})$$
 (2)

Given the notation above, we can now define the items in tuple  $\langle S, tok, detok \rangle$ . A tokeniser's **vocabulary** is a set of subwords  $S \subset \Sigma^+$  such that  $\Sigma \subseteq S$ ;<sup>6</sup> we say its size is  $|S| = |\Sigma| + K$ . Further, a **detokenisation function** is defined as detok :  $S^* \to \Sigma^*$  and given a subword-string it outputs the character-string it represents. This function thus is simply defined as detok(s)  $\stackrel{\text{def}}{=} \text{concat}(s)$ .

Finally, we are left with defining a **tokenisation function** tok :  $\Sigma^* \to S^*$ , which maps from character- to subword-strings. Notably, these functions always ensure the equivalence  $\mathbf{c} \stackrel{\circ}{=} \mathbf{s}$  for  $\mathbf{s} = \operatorname{tok}(\mathbf{c})$ . Several tokenisation functions, however, are compatible with this constraint, as given a vocabulary, many subword-strings may be equivalent to the same character-string. For instance, given  $S = \{a, c, t, at\}$ , the string  $\mathbf{c} = \langle c, a, t \rangle$  could be tokenised as  $\mathbf{s} = \langle c, a, t \rangle$  or as  $\mathbf{s} = \langle c, at \rangle$ . Most researchers define tokenisation functions in one of two ways, which we term direct and bottom-up tokenisation functions here; we define these next.

#### 3.1 Direct Tokenisation Functions

In direct tokenisation, a character-string is directly replaced by an optimal subword-string. To

<sup>&</sup>lt;sup>3</sup>Recent work tries to improve the computational efficiency of byte-level models (Yu et al., 2023; Pagnoni et al., 2024).

<sup>&</sup>lt;sup>4</sup>Although, see also Ali et al. (2024), who argue that compression might be a necessary but not sufficient condition for good downstream performance, and Schmidt et al. (2024), who argue that compression and downstream performance have a more complex relationship than prior work suggests.

<sup>&</sup>lt;sup>5</sup>We note that we use set notation here, but our datasets are actually multisets—datasets can include the same string c multiple times. We show that tokenisation is still NP-complete for datasets with no repetitions in §6.3. Further, we impose no constraint on the kind of string present in these datasets: each  $c_n$  can be either a raw or pre-tokenised character-string (i.e., either a full document or a whitespace-separated word).

 $<sup>{}^{6}\</sup>Sigma \subseteq S$  is typically enforced to guarantee that every character-string can be represented by at least one token-string.

implement this, one must thus first define what *optimal* means; this is done through an objective function  $\mathfrak{G}$  which, given a subword-string, returns a score. Given a previously chosen vocabulary  $\mathcal{S}$  (we discuss how to find  $\mathcal{S}$  in §5), a direct tokenisation function then encodes string c as:

$$tok_{\phi}[S](\mathbf{c}) = \operatorname*{arg\,max}_{\mathbf{s}\in S^{*}} \mathfrak{G}(\mathbf{s})$$
(3)  
s.t.  $\mathbf{s} \stackrel{\circ}{=} \mathbf{c}$ 

In words, given a vocabulary S, function  $tok_{\diamond}$  returns the optimal subword-string  $s \in S^*$  which is equivalent to the input character-string c. We then set  $tok(c) \stackrel{\text{def}}{=} tok_{\diamond}[S](c)$ . Different choices of  $\mathfrak{G}$  recover methods such as UnigramLM (Kudo, 2018) or PathPiece (Schmidt et al., 2024). Notably, in general, this function is not efficiently computable.<sup>7</sup>

In this paper, we are concerned with tokenisers that use compression as their objective: that is, for which  $\mathfrak{G}(\mathbf{s}) = -|\mathbf{s}|$ . In this case, we can rewrite the direct tokenisation function as:

$$tok_{\diamond}[\mathcal{S}](\mathbf{c}) = \underset{\mathbf{s}\in\mathcal{S}^{*}}{\arg\min|\mathbf{s}|}$$
(4)  
s.t.  $\mathbf{s} \stackrel{\circ}{=} \mathbf{c}$ 

Importantly, in the case of compression, this equation can be computed efficiently (as shown in §5.1).

## **3.2 Bottom-up Tokenisation Functions**

In bottom-up tokenisation, one starts with a set of character-strings, and merges their symbols bottom-up, one pair at a time.<sup>8</sup> Formally, let  $m \in \mathcal{M}$  be a **merge**, defined as a pair of subwords:  $m = \langle s_1, s_2 \rangle$ . Further, let  $\mathcal{M} \stackrel{\text{def}}{=} \Sigma^+ \times \Sigma^+$ . Now, let **merge** be a functional; given merge  $m = \langle s_1, s_2 \rangle$ , it returns a function merge[m] :  $\mathcal{S}^* \to (\mathcal{S} \cup \{s_1 \circ s_2\})^*$  which operates on string **s** left-to-right, replacing every occurrence of  $s_1$  followed by  $s_2$  in it with subword  $s' = s_1 \circ s_2$ . E.g., given  $\mathbf{s} = \langle wo, r, ld \rangle$  and  $m = \langle wo, r \rangle$ , the output of  $\text{merge}[m](\mathbf{s})$  is  $\langle wor, ld \rangle$ .

Consider now  $\mathbf{m} \in \mathcal{M}^*$ , a sequence of merges. Given a character-string  $\mathbf{c} \in \Sigma^*$ , a bottom-up tokenisation function compresses it as:

$$\mathsf{tok}_{\uparrow}[\mathbf{m}](\mathbf{c}) = \left( \bigcup_{z=1}^{|\mathbf{m}|} \mathsf{merge}[m_z] \right)(\mathbf{c}) \quad (5)$$

where  $\bigcirc$  represents function composition, e.g.,  $\bigcirc_{z=1}^{2} \operatorname{merge}[m_{z}] = \operatorname{merge}[m_{2}] \odot \operatorname{merge}[m_{1}].$ Bottom-up tokenisers then set  $\operatorname{tok} \stackrel{\text{def}}{=} \operatorname{tok}_{\uparrow}[\mathbf{m}].$ Further, a merge sequence  $\mathbf{m}$  is also used to set a bottom-up tokeniser's vocabulary as:

$$S = \Sigma \cup \{s_1 \circ s_2 \mid \langle s_1, s_2 \rangle \in \mathbf{m}\}$$
(6)

where  $|\mathbf{m}| = K$  implies this vocabulary has size  $|\mathcal{S}| = |\Sigma| + K$ , as before.

# 4 Maximum 2-Satisfiability

Our paper's goal is to prove the NP-completeness of tokenisation. To show this, we must reduce an NP-hard problem to tokenisation in polynomial time. We will rely on the **maximum 2-satisfiability** problem (max-2-SAT) for this, whose definition we provide here. The NP-hardness of max-2-SAT was proven by Garey et al. (1974).

**Definition 1.** Let  $\mathcal{X} = \{X_j\}_{j=1}^J$  be a set of variables; each of these variables are assigned values  $x_j \in \{F, T\}$ , and we write  $\chi = \{x_j\}_{j=1}^J \in \{F, T\}^J$ . Let  $\mathcal{L} = \{(L_i^1 \lor L_i^2)\}_{i=1}^I$  be a set of clauses,<sup>9</sup> where each literal L represents either a variable  $X_j$  or its negation  $\neg X_j$ . The max-2-SAT decision problem requires deciding whether there exists an assignment for which at least  $\psi$  clauses are satisfied:

$$\psi \le \max_{\chi \in \{\mathsf{F},\mathsf{T}\}^J} \sum_{i=1}^{I} \mathbb{1}\{L_i^1 \lor L_i^2\}$$
(7)

where  $\mathbb{1}$  is an indicator function which evaluates the clause and returns one if the clause is satisfied by  $\chi$  and zero otherwise.

For mathematical convenience, we will write  $M2S(\mathcal{X}, \mathcal{L}, \psi)$  for a function which returns T if its input is satisfiable under a max-2-SAT decision problem, and F otherwise. As a concrete example, consider the set of variables  $\mathcal{X} = \{X_1, X_2\}$  and the set of clauses  $\mathcal{L} = \{X_1 \lor X_2, \neg X_1 \lor X_2, X_1 \lor \neg X_2, \neg X_1 \lor \neg X_2\}$ . The assignment  $x_1 = T, x_2 = T$  leads to 3 clauses being satisfied, which is the optimum. For this example, we thus have that  $M2S(\mathcal{X}, \mathcal{L}, 3) = T$ , but that  $M2S(\mathcal{X}, \mathcal{L}, 4) = F$ .

# 5 Finding an Optimal Direct Tokeniser

We are now left with the task of finding an optimal tokeniser. We do this by selecting either: its vocabulary in direct tokenisation, since  $tok = tok_{\Rightarrow}[S]$ ;

<sup>&</sup>lt;sup>7</sup>In fact, Geh et al. (2024) shows that it is NP-complete for  $\mathfrak{G}(\mathbf{s}) = \sum_{t=1}^{|\mathbf{s}|} \log p_{\theta}(s_t | \mathbf{s}_{< t})$ , where  $p_{\theta}$  is a language model.

<sup>&</sup>lt;sup>8</sup>Currently, this is likely the most common tokenisation function, being used in popular tokenisers such as, e.g., GPT-4's (OpenAI et al., 2024), LLaMA's (Touvron et al., 2023a,b), and Pythia's (Biderman et al., 2023).

 $<sup>^{9}</sup>$ max-2-SAT also allows clauses to have a single literal  $L_i$ . In this case, we can always rewrite the clause as  $(L_i \lor L_i)$  with no change to the solution of this decision problem.

304

305

306

307

308

310

311

312

314

315

316

317

318

9

279

281

or its merge sequence in bottom-up tokenisation, since  $tok = tok_{\uparrow}[m]$  and since its vocabulary is chosen according to Eq. (6). (Note that in §3, we only showed how to apply tokenisers at inference time, but not how to find them.) In this section, we focus on direct tokenisation, defining its optimisation and decision problems; we then prove its NP-completeness. The optimisation problem is defined as follows.

**Definition 2.** Given a dataset  $\mathcal{D}$  and a vocabulary size K, the **direct tokenisation optimisation problem** is to find a vocabulary  $S^* \subset \Sigma^+$  which maximally compresses  $\mathcal{D}$ :

244

247

249

252

258

259

261

263

267

271

272

273 274

275

276

278

$$S^{\star} = \underset{S \subset \Sigma^{+}}{\operatorname{arg\,min}} \sum_{\mathbf{c} \in \mathcal{D}} |\operatorname{tok}_{\phi}[S](\mathbf{c})| \qquad (8)$$
  
s.t.  $|S| = |\Sigma| + K$ 

We can similarly define direct tokenisation's decision problem.

**Definition 3.** Given a dataset  $\mathcal{D}$  and a vocabulary size K, the **direct tokenisation decision problem** requires deciding whether there exists a vocabulary  $S \subset \Sigma^+$  which compresses  $\mathcal{D}$  to at most  $\delta$  symbols:

$$\delta \ge \min_{\boldsymbol{\mathcal{S}} \subset \Sigma^{+}} \sum_{\mathbf{c} \in \mathcal{D}} |\mathsf{tok}_{\phi}[\boldsymbol{\mathcal{S}}](\mathbf{c})| \qquad (9)$$
  
s.t.  $|\boldsymbol{\mathcal{S}}| = |\Sigma| + K$ 

We write  $\operatorname{Tok}_{\diamond}(\mathcal{D}, K, \delta)$  for a function which returns T if a direct tokenisation decision problem with those inputs is satisfiable, and F otherwise. Note that, whenever  $|\mathcal{D}| \leq K$ , the solution to the problem above is trivial, as an optimal solution simply requires including all strings  $\mathbf{c}_n$  in vocabulary S. As we show next, however, in the general case the above decision problem is NP-complete. We now state this as a theorem, which we will prove in the next two sections.

**Theorem 1.** *The direct tokenisation decision problem, as in Definition 3, is NP-complete.* 

*Proof.* A decision problem is considered to be NPcomplete if: (i) it is in NP; (ii) it is NP-hard. We prove these conditions in §5.1 and §5.2.  $\Box$ 

#### 5.1 Direct Tokenisation is in NP

A decision problem is in the nondeterministic polynomial time class (NP) if, given a **certificate** of polynomial length, one can verify that certificate in polynomial time. Specifically, a certificate usually encodes a decision problem's solution, allowing us to verify its satisfiability. In the case of direct tokenisation, this certificate would be a vocabulary S which leads a dataset D to be compressed to at most  $\delta$  symbols. Verifying this certificate simply requires computing the sum in Eq. (9), i.e.:

$$\sum_{\mathbf{c}\in\mathcal{D}}|\mathsf{tok}_{\diamond}[\mathcal{S}](\mathbf{c})| \tag{10}$$

**Lemma 1.** The direct tokenisation decision problem, as in Definition 3, is in NP.

*Proof.* As noted above, whenever  $|\mathcal{D}| \leq K$ , each  $\mathbf{c}_n \in \mathcal{D}$  can be included in the vocabulary  $\mathcal{S}$  and fully compressed to a single symbol; we can thus verify the problem's satisfiability by simply checking that  $\delta \geq |\mathcal{D}|$  as this is the best reachable compression. Assuming K to be bounded by  $|\mathcal{D}|$ —and therefore polynomial in the input—we have that the certificate S also has polynomial length. Given such a certificate  $\mathcal{S}$ , verifying it simply requires computing the sum in Eq. (10). In turn, computing this sum requires  $|\mathcal{D}|$  calls to function tok. It follows that, if function toke runs in polynomial time, then direct tokenisation is in NP. Luckily, this function can indeed be computed efficiently using Schmidt et al.'s (2024) PathPiece method, which runs in  $O(|\mathbf{c}|^2)$  time. 

## 5.2 Direct Tokenisation is NP-hard

We now use a reduction from max-2-SAT to prove the NP-hardness of direct tokenisation.

**Reduction 1.** Let us have an instance of the max-2-SAT decision problem as in Definition 1. To reduce this instance to an instance of the direct tokenisation decision problem, as in Definition 3, we first define alphabet  $\Sigma = \{ \odot \} \cup \{x_j^{\mathsf{T}}, x_j^{\mathsf{F}}\}_{j=1}^J$ . We then construct three sets of strings:

$$\mathcal{D}_1 = \{ \odot x_j^{\mathsf{T}} \odot \}_{j=1}^J \cup \{ \odot x_j^{\mathsf{F}} \odot \}_{j=1}^J \qquad (11a)$$

$$\mathcal{D}_2 = \{ \odot x_j^{\mathsf{T}} \odot x_j^{\mathsf{F}} \odot \}_{j=1}^J \tag{11b}$$
 313

$$\mathcal{D}_3 = \{ \textcircled{O}L_i^1 \textcircled{O}L_i^2 \textcircled{O}\}_{i=1}^I \tag{11c}$$

In these strings  $L_i$  is replaced by either character  $x_j^{\mathsf{T}}$  or  $x_j^{\mathsf{F}}$ , depending on whether it represents  $X_j$  or  $\neg X_j$ , respectively. We then construct our dataset  $\mathcal{D}$ , and choose K and  $\delta$  as:

$$\mathcal{D} = \left(\bigcup_{=1}^{f} \mathcal{D}_{1}\right) \cup \left(\bigcup_{=1}^{f'} \mathcal{D}_{2}\right) \cup \mathcal{D}_{3}$$
(12a) 31

$$K = J, \quad \delta = (4f + 3f') J + 5 I - 2\psi \quad (12b)$$
 320

where we set  $f' \stackrel{\text{def}}{=} 4I + 1$  and  $f \stackrel{\text{def}}{=} 4f'J + 4I + 1$ . 321

394

395

396

397

398

399

400

401

403

404

405

406

407

408

409

410

411

412

368

370

371

372

We write  $R1(\mathcal{X}, \mathcal{L}, \psi)$  to represent a function which, given an instance of max-2-SAT, returns an instance of the tokenisation problem given by our reduction (i.e.,  $\mathcal{D}, K, \delta$ ). For our reduction to be correct, we must have that:

324

327

329

331

332

334

335

337

341

346

347

348

354

364

367

$$M2S(\mathcal{X}, \mathcal{L}, \psi) \iff Tok_{\diamond}(R1(\mathcal{X}, \mathcal{L}, \psi))$$
 (13)

meaning that a max-2-SAT instance is satisfiable if and only if its reduced direct tokenisation instance is as well. We now set out to prove this. We start by proving the forward direction of this iff clause.

**Lemma 2.** If a max-2-SAT instance is satisfiable, then the direct tokenisation instance output by Reduction 1 is also satisfiable. Formally:

$$M2S(\mathcal{X}, \mathcal{L}, \psi) \implies Tok_{\phi}(R1(\mathcal{X}, \mathcal{L}, \psi))$$
 (14)

Proof sketch. See a formal proof in App. A. Our proof works by first fixing a satisfying solution to max-2-SAT with values  $x_j^*$ . Given this solution, for each variable, we add to our vocabulary S a subword  $\odot x_j^T \odot$  if  $x_j^*$  is true, or  $\odot x_j^F \odot$  if  $x_j^*$  is false. Given these subwords, strings in  $\mathcal{D}_1$  and  $\mathcal{D}_2$  occupy a total length of (4f + 3f') J. Further, since at least  $\psi$  of the max-2-SAT clauses are satisfied by  $x_j^*$ , the strings in  $\mathcal{D}_3$  will occupy a total length smaller or equal to  $5I - 2\psi$ . This solution to the tokenisation problem thus gives us a total length which is smaller or equal to  $\delta = (4f + 3f') J + 5I - 2\psi$ .

Now, we are left with proving the backward direction of the iff clause in Eq. (13). We do so in the following lemma.

**Lemma 3.** If the direct tokenisation instance output by Reduction 1 is satisfiable, the max-2-SAT instance which generated it is as well. Formally:

 $\operatorname{Tok}_{\phi}(\operatorname{R1}(\mathcal{X},\mathcal{L},\psi)) \Longrightarrow \operatorname{M2S}(\mathcal{X},\mathcal{L},\psi)$  (15)

**Proof sketch.** See a formal proof in App. B. Our proof works in three steps. First, we show that all satisfying solutions must only have subwords of the form  $@x_j^T @$  or  $@x_j^F @$ , since this is required to compress strings in  $\mathcal{D}_1$  to at most 4fJ symbols. Second, we show that all satisfying solutions must only have either subword  $@x_j^T @$  or  $@x_j^F @$  for any variable  $X_j$ ; this is required to compress strings in  $\mathcal{D}_2$  to at most 3f'J symbols. Finally, we show that if a tokeniser compresses strings in  $\mathcal{D}_3$  to  $5I - 2\psi$ , then there is an assignment  $\chi$  which satisfies at least  $\psi$  of the original max-2-SAT problem. Given both lemmas above, we can now trivially prove that direct tokenisation is NP-hard. **Lemma 4.** *The direct tokenisation decision problem, as in Definition 3, is NP-hard.* 

*Proof.* First, it is easy to see that Reduction 1 runs in polynomial time. Second, max-2-SAT is an NP-hard problem (Garey et al., 1974). This lemma then follows trivially from Lemmas 2 and 3, which together show that an instance of the tokenisation problem generated through Reduction 1 is satisfiable if and only if the max-2-SAT instance used to produce it is also satisfiable.

## 6 Finding Optimal Bottom-up Tokenisers

In this section, we shift our attention to bottomup tokenisation. We define both its optimisation and decision problems, and then prove its NPcompleteness. We start with defining the optimisation problem.

**Definition 4.** Given a dataset  $\mathcal{D}$  and a vocabulary size K, the **bottom-up tokenisation optimisation** problem is to find a merge sequence  $\mathbf{m}^* \in \mathcal{M}^*$  which maximally compresses  $\mathcal{D}$ :

$$\mathbf{m}^{\star} = \operatorname*{arg\,min}_{\mathbf{m}\in\mathcal{M}^{\star}} \sum_{\mathbf{c}\in\mathcal{D}} |\mathsf{tok}_{\uparrow}[\mathbf{m}](\mathbf{c})| \qquad (16)$$

s.t. 
$$|\mathbf{m}| = K$$
 3

As can be seen, this optimisation problem is similar to the direct tokenisation problem, albeit its target is to find a merge sequence instead of a vocabulary. We similarly define a decision problem.

**Definition 5.** Given a dataset  $\mathcal{D}$  and a vocabulary size K, the **bottom-up tokenisation decision prob***lem* requires deciding whether there exists a merge sequence  $\mathbf{m} \in \mathcal{M}^*$  which compresses  $\mathcal{D}$  to at most  $\delta$  symbols:

$$\delta \ge \min_{\mathbf{m} \in \mathcal{M}^*} \sum_{\mathbf{c} \in \mathcal{D}} |\mathsf{tok}_{\uparrow}[\mathbf{m}](\mathbf{c})|$$
(17)

s.t. 
$$|\mathbf{m}| = K$$
 40

We write  $\operatorname{Tok}_{\uparrow}(\mathcal{D}, K, \delta)$  for a function which returns T if a bottom-up tokenisation decision problem with those inputs is satisfiable, and F otherwise. We spend the rest of this section showing that bottom-up tokenisers are NP-complete.

**Theorem 2.** *The bottom-up tokenisation decision problem, as in Definition 5, is NP-complete.* 

*Proof.* We prove this in two steps below. We first prove that this problem is in NP, in §6.1. We then prove that this problem is NP-hard, in §6.2.  $\Box$ 

r

#### 6.1 Bottom-up Tokenisation is in NP

413

414

415

416

417

418

419

420

421

422

423

446

447

448

449

450

451

452

453 454

455

456 457

458

We can verify this using a solution, the merge sequence  $\mathbf{m} \in \mathcal{M}^*$ , as a certificate. By showing that this certificate has polynomial length and that it can be verified in polynomial time, we prove this problem is in NP. To verify this certificate, we simply need to compute the sum in Eq. (17), i.e.:

 $\sum_{\mathbf{c}\in\mathcal{D}} |\mathsf{tok}_{\uparrow}[\mathbf{m}](\mathbf{c})| \tag{18}$ 

which we show now can be done efficiently.

**Lemma 5.** The bottom-up tokenisation decision problem, as in Definition 5, is in NP.

*Proof.* First, if K is larger than the total number 424 of characters in  $\mathcal{D}$ , i.e.,  $\sum_{\mathbf{c}\in\mathcal{D}} |\mathbf{c}|$ , then this dataset 425 can be compressed to  $|\mathcal{D}|$  by merging each string 426 down to a single symbol; further, compressing  $\mathcal{D}$ 427 more than that is not possible independently of K. 428 Verifying the satisfiability of such an instance of the 429 tokenisation problem is thus trivial, only requiring 430 checking if  $\delta \geq |\mathcal{D}|$ . Second, if K is bounded by 431  $|\mathcal{D}|$ —and therefore polynomial in the input—the 432 433 certificate **m** has polynomial length. Given such a certificate **m**, verifying it then simply requires 434 computing the sum in Eq. (18). In turn, comput-435 ing this sum requires  $|\mathcal{D}|$  calls to function tok<sub>1</sub>. It 436 follows that, if function tok<sub>↑</sub> runs in polynomial 437 438 time, then bottom-up tokenisation is in NP. The 439 computation of  $tok_{\uparrow}$ , can be done in polynomial time following the structure described in §3.2. For 440 each  $m = \langle s_1, s_2 \rangle \in \mathbf{m}$ , scan the current **c** and 441 replace each occurrence of  $s_1, s_2$  by s'. This takes 442 time  $\mathcal{O}(|\mathbf{c}|)$  for each merge. Afterwards, the result-443 ing string can be compared against the desired size. 444 We obtain a total runtime of  $O(|\mathcal{D}||\mathbf{c}||\mathbf{m}|)$ . 445

#### 6.2 Bottom-up Tokenisation is NP-hard

As before, we use a reduction from max-2-SAT to prove bottom-up tokenisation's NP-hardness.

**Reduction 2.** Let us have an instance of the max-2-SAT decision problem as in Definition 1. To reduce this instance to an instance of the bottom-up tokenisation decision problem, as in Definition 5, we first define alphabet  $\Sigma = \{ \odot, \otimes \} \cup \{x_j^T, x_j^F\}_{j=1}^J$ . We then construct five sets of strings:

$$\mathcal{D}_{1} = \{ \odot x_{j}^{\mathsf{T}} \}_{j=1}^{J} \cup \{ x_{j}^{\mathsf{F}} \odot \}_{j=1}^{J} \cup \{ x_{j}^{\mathsf{T}} \odot \}_{j=1}^{J}$$
(19)
$$\cup \{ \odot x_{i}^{\mathsf{F}} \}_{i=1}^{J} \cup \{ x_{i}^{\mathsf{T}} \otimes \}_{i=1}^{J} \cup \{ \otimes x_{i}^{\mathsf{F}} \}_{i=1}^{J}$$

$$\mathcal{D}_{2} = \{ \odot x_{j}^{\mathsf{T}} \odot \}_{j=1}^{J} \cup \{ \odot x_{j}^{\mathsf{F}} \odot \}_{j=1}^{J} \cup \{ \odot x_{j}^{\mathsf{F}} \odot \}_{j=1}^{J} \cup \{ \odot x_{j}^{\mathsf{T}} \odot \}$$

$$\mathcal{D}_3 = \{ \odot x_j^{\mathsf{T}} \odot x_j^{\mathsf{F}} \odot \}_{j=1}^J \cup \{ \bigotimes x_j^{\mathsf{F}} \odot x_j^{\mathsf{T}} \otimes \}_{j=1}^J$$

$$459$$

$$\mathcal{D}_4 = \{ \odot x_j^{\mathsf{F}} \odot x_j^{\mathsf{T}} \otimes \}_{j=1}^J \cup \{ \otimes x_j^{\mathsf{F}} \odot x_j^{\mathsf{T}} \odot \}_{j=1}^J$$

$$460$$

$$\mathcal{D}_{5} = \begin{cases} \bigotimes_{j}^{\mathbb{T}} \odot x_{j}^{\mathbb{F}} \odot & \text{if } L_{i}^{1} = X_{j} \text{ and } L_{i}^{2} = \neg X_{j'} \\ \bigotimes_{j}^{\mathbb{T}} \odot x_{j}^{\mathbb{F}} \odot & \text{if } L_{i}^{1} = \neg X_{j} \text{ and } L_{i}^{2} = X_{j'} \\ \bigotimes_{j}^{\mathbb{F}} \odot x_{j'}^{\mathbb{F}} \odot x_{j'}^{\mathbb{F}} \odot \text{ if } L_{i}^{1} = \neg X_{j} \text{ and } L_{i}^{2} = \neg X_{j'} \\ \bigotimes_{j}^{\mathbb{T}} \odot x_{j'}^{\mathbb{T}} \odot x_{j'}^{\mathbb{T}} \otimes \text{ if } L_{i}^{1} = X_{j} \text{ and } L_{i}^{2} = X_{j'} \end{cases} \end{cases}$$

462

463

466

470

471

472

473

474

475

476

477

478

479

480

481

482

483

484

485

486

487

488

489

491

492

493

We then construct our dataset D, and choose K and  $\delta$  as:

$$\mathcal{D} = \bigcup_{=1}^{f} \mathcal{D}_1 \cup \bigcup_{=1}^{f'} \mathcal{D}_2 \cup \bigcup_{=1}^{f''} \mathcal{D}_3 \cup \bigcup_{=1}^{f'''} \mathcal{D}_4 \cup \mathcal{D}_5 \quad (20)$$

$$K = 8J, \ \delta = (6f + 6f' + 4f'' + 4f''') \ J + 3 \ I - \psi$$

where we set:

$$f''' \stackrel{\text{def}}{=} 5I, \quad f'' \stackrel{\text{def}}{=} 10 f''' J + 5I$$
 (21a) 40

$$f' \stackrel{\text{def}}{=} (10f'' + 10f''') J + 5I \tag{21b}$$

$$f \stackrel{\text{\tiny def}}{=} (12f' + 10f'' + 10f''') J + 5I \quad (21c)$$

As before, we write  $R2(\mathcal{X}, \mathcal{L}, \psi)$  for a function which, given an instance of the max-2-SAT problem, returns an instance of the bottom-up tokenisation problem. For our reduction to be correct, we must have that:

$$M2S(\mathcal{X}, \mathcal{L}, \psi) \iff Tok_{\uparrow}(R2(\mathcal{X}, \mathcal{L}, \psi))$$
 (22)

We follow the same proof strategies as before, starting by proving the forward direction of this iff statement.

**Lemma 6.** If a max-2-SAT instance is satisfiable, then the bottom-up tokenisation instance output by Reduction 2 is also satisfiable. Formally:

$$M2S(\mathcal{X}, \mathcal{L}, \psi) \implies Tok_{\uparrow}(R2(\mathcal{X}, \mathcal{L}, \psi)) \quad (23)$$

*Proof sketch.* See a formal proof in App. C. Without loss of generality, let a satisfying solution to max-2-SAT have values  $x_j^*$ . Our proof works by first defining the three following lists of merges, which must be included in any satisfying solution to this tokenisation problem:

$$\mathbf{n}_{1} = \bigcirc_{j=1}^{J} [\langle \otimes, x_{j}^{\mathsf{F}} \rangle, \langle x_{j}^{\mathsf{T}}, \otimes \rangle]$$
(24a)

$$\mathbf{m}_{3} = \bigcirc_{j=1}^{J} [\langle x_{j}^{\mathsf{F}}, \odot \rangle, \langle \odot, x_{j}^{\mathsf{T}} \rangle]$$
(24b)

$$\mathbf{m}_{5} = \bigcirc_{j=1}^{J} [\langle \odot, x_{j}^{\mathsf{F}} \rangle, \langle x_{j}^{\mathsf{T}}, \odot \rangle]$$
(24c)

We then construct two other lists of merges, which depend on the satisfying assignments to

495

496

497

498

499

501

502

504

510

511

512

513

514

516

517

518

521

522

524

525

527

530

531

533

max-2-SAT:

1

$$\mathbf{m}_{2} = \bigcirc_{j=1}^{J} \begin{bmatrix} \langle \odot, x_{j}^{\mathsf{T}} \otimes \rangle & \text{if } x_{j}^{\star} = \mathsf{T} \\ \langle \otimes x_{j}^{\mathsf{F}}, \odot \rangle & \text{else} \end{bmatrix}$$
(25a)

 $\mathbf{m}_4 = \bigcirc_{j=1}^J \left[ \begin{array}{cc} (\textcircled{o} x_j, \textcircled{o}) & \text{II } x_j = \mathbf{I} \\ \langle \textcircled{o}, x_j^{\mathsf{F}} \textcircled{o} \rangle & \text{else} \end{array} \right]$ (25b)

Finally, we create a merge sequence by concatenating these lists in order:

$$\mathbf{m} = \mathbf{m}_1 \circ \mathbf{m}_2 \circ \mathbf{m}_3 \circ \mathbf{m}_4 \circ \mathbf{m}_5 \qquad (26)$$

Note that we have exactly K = 8J merges in this list. Given this merge sequence, it is easy to verify that strings in  $\mathcal{D}_1$  to  $\mathcal{D}_4$  will use exactly (6f + 6f' + 4f'' + 4f''') J symbols after being tokenised. Further, since at least  $\psi$  of the max-2-SAT problem are satisfied by  $x_j^{\star}$ , the strings in  $\mathcal{D}_5$  will occupy a total length smaller or equal to  $3I - \psi$ . This solution to the tokenisation problem thus gives us a tokeniser which will compress  $\mathcal{D}$  to at most  $\delta = (6f + 6f' + 4f'' + 4f''') J + 3I - \psi$ .

We now prove the backward direction of the iff clause in Eq. (22).

**Lemma 7.** If the bottom-up tokenisation instance output by Reduction 2 is satisfiable, the max-2-SAT instance which generated it is as well. Formally:

$$\operatorname{Tok}_{\uparrow}(\operatorname{R2}(\mathcal{X},\mathcal{L},\psi)) \Longrightarrow \operatorname{M2S}(\mathcal{X},\mathcal{L},\psi)$$
 (27)

**Proof sketch.** See a formal proof in App. D. Our proof works in five steps. First, we show that all satisfying solutions must include merges  $\mathbf{m}_1$ ,  $\mathbf{m}_3$ , and  $\mathbf{m}_5$  from Eq. (24), since this is required to compress strings in  $\mathcal{D}_1$  to at most 6fJ symbols. Second, we show the other merges of any satisfying solution must be of the form:

$$\mathbf{m}_{j}^{\circledast} = \left\{ \begin{array}{l} \langle \odot x_{j}^{\mathsf{T}}, \odot \rangle, \langle \odot, x_{j}^{\mathsf{F}} \odot \rangle \\ \langle \odot, x_{j}^{\mathsf{T}} \odot \rangle, \langle \odot x_{j}^{\mathsf{F}}, \odot \rangle \end{array} \right\}$$
(28a)

$$\mathbf{m}_{j}^{\otimes} = \left\{ \begin{array}{l} \langle \odot, x_{j}^{\mathsf{T}} \otimes \rangle, \langle \otimes x_{j}^{\mathsf{F}}, \odot \rangle \\ \langle \odot x_{j}^{\mathsf{T}}, \otimes \rangle, \langle \otimes, x_{j}^{\mathsf{F}} \odot \rangle \end{array} \right\}$$
(28b)

this is required to compress strings in  $\mathcal{D}_2$  to at most 6f'J symbols. Third, we show that any satisfying solution will have at least one merge of each set  $\mathbf{m}_j^{\otimes}$  and one of each set  $\mathbf{m}_j^{\otimes}$ ; this is required to compress strings in  $\mathcal{D}_3$  to at most 4f''J symbols. Fourth, we show that any satisfying solution will have—for each  $j \in \{1, J\}$ —both its merges in sets  $\mathbf{m}_j^{\otimes}$  and  $\mathbf{m}_j^{\otimes}$  containing character  $x_j^{\mathrm{T}}$  or containing character  $x_j^{\mathrm{F}}$ ; this is required to compress strings in

 $\mathcal{D}_4$  to at most 4f'''J symbols. Finally, we show that if a tokeniser compresses strings in  $\mathcal{D}_5$  to  $3I - \psi$ , then there is an assignment  $\chi$  which satisfies at least  $\psi$  of the original max-2-SAT problem.  $\Box$ 

534

535

536

537

538

539

540

541

542

543

544

545

547

548

549

551

552

553

554

555

556

557

558

559

560

561

562

563

564

565

566

567

569

570

571

572

573

574

575

576

577

578

579

580

582

Finally, given both lemmas above, we can now prove that bottom-up tokenisation is NP-hard.

**Lemma 8.** The bottom-up tokenisation decision problem, as in Definition 5, is NP-hard.

*Proof.* First, it is easy to see that Reduction 2 runs in polynomial time. Second, max-2-SAT is an NP-hard problem (Garey et al., 1974). This lemma then follows trivially from Lemmas 6 and 7.

#### 6.3 Other Definitions of Tokenisation

We now expand our discussion to consider variations of the above tokenisation problems.

Deduped Datasets. Our definitions of both direct and bottom-up tokenisation allow datasets  $\mathcal{D}$  to include repeated entries. It is common, however, to deduplicate datasets in NLP-thus removing repeated entries. A small change to both our reductions is enough to adapt it to this deduplicated dataset case: simply append each string in the repeated datasets (either  $\mathcal{D}_1$  and  $\mathcal{D}_2$  in Reduction 1 or  $\mathcal{D}_1$  to  $\mathcal{D}_4$  in Reduction 2) with a unique character  $\{a_y\}_{y=1}^\infty$  and increase the target compression size  $\delta$  accordingly (by f + f' or f + f' + f'' + f''', respectively). These new characters will never be included in optimal tokenisers' solutions, and thus the previous proofs hold, with the difference that each dataset will require extra symbols once compressed.

A Single Long String. In the previous sections, we considered tokenisers trained on a dataset  $\mathcal{D}$ . Work on compression, however, usually considers a single long string c as its input. It is easy to see that direct tokenisation is not an NP-complete problem if its input is a single long string; including this string in vocabulary  $\mathcal{S}$  already achieves optimal compression. Bottom-up tokenisation, however, is still NP-complete even when given a single string as input. As before, this can be shown with a similar strategy to Reduction 2, but where we first append each string in dataset  $\mathcal{D}$  with a unique character  $\{a_y\}_{y=1}^{\infty}$  and then concatenate all these strings. As in the deduped case above, characters  $a_{y}$  will never be merged by any optimal tokeniser; they will thus serve as virtual string delimiters and will not affect our proofs beyond an increase to the target compression size  $\delta$ .

638

639

640

641

642

643

644

645

646

647

648

649

650

651

652

653

654

655

656

657

658

659

660

661

662

663

664

665

666

667

668

669

670

671

672

673

674

627

628

629

A Hybrid Approach. Finally, the last variant we consider is a hybrid between direct and bottom-up tokenisation, where we find a merge sequence **m** which—when we extract a vocabulary from it as  $S = \Sigma \cup \{s_1 \circ s_2 \mid \langle s_1, s_2 \rangle \in \mathbf{m}\}$ —optimally compresses a dataset  $\mathcal{D}$  using the direct tokenisation function in Eq. (4). We can easily prove the NP-hardness of this tokenisation variant by relying on Reduction 2; as our proof in Lemma 8 did not make use of the order of merges in **m**, only of the subwords composed by it, this lemma's proof strategy can be similarly applied to this hybrid variant.

583

584

585

588

589

593

594

598

602

604

610

611

612

613

614

615

616

617

618

625

# 7 Tokenisation's Connection to Compression

The variants of tokenisation that we consider here with compression as their objective function—are closely related to the field of dictionary compression. In both fields, we wish to reduce the size of an input (c or D) by exploiting repetitive elements. In fact, the most popular tokenisation algorithm to date, BPE, was originally proposed as a compression algorithm (Gage, 1994) and has only somewhat recently been ported into NLP to find tokenisers (by Sennrich et al., 2016).

Not surprisingly, prior work has also considered, from a theoretical perspective, the compression tokenisers achieve. Zouhar et al. (2023b), for instance, analyse bottom-up tokenisation and prove an approximation bound on the compression achieved by the tokenisers found using BPE. More recently, Kozma and Voderholzer (2024) also analyses bottom-up tokenisation, proving a tighter bound on this compression achieved by BPE.

A popular dictionary compression method, the **straight-line program** (SLP; Kieffer and Yang, 2000; Charikar et al., 2005), can be used to illustrate the similarities and differences between tokenisers and compressors.<sup>10</sup> Given a string **c**, an SLP describes a context-free grammar from which **c** can be uniquely derived. Formally, an SLP in Chomsky normal form (CNF) is a set of rules of form  $X \rightarrow a$  or  $X \rightarrow AB$ , where X, A, B are called nonterminals and *a* is a terminal.<sup>11</sup> Starting from a special nonterminal *S*, applying these

rules exhaustively—until only terminals are left produces exactly the desired string c. Notably, given a string c, it is NP-complete to find the smallest SLP which generates it (Charikar et al., 2005).

On the one hand, SLPs in CNF are closely linked to bottom-up tokenisation; each of its rules expands to two nonterminals, and thus corresponds to a merge. However, while SLPs must find the minimum number of merges (or rules) to fully compress a string into a single symbol, bottom-up tokenisers must maximally compress the string given a fixed number of merges. On the other hand, SLPs which are not in CNF (that is, for which other contextfree production rules are allowed, as long as the decoding stays unique) are closely linked to direct tokenisation. In this case, a direct tokeniser could be converted into an SLP with depth two; this grammar has a start rule  $S \rightarrow s$ , and a rule from each subword to its characters  $s \rightarrow c$ . Again, while SLPs must find a minimal grammar representing the string, direct tokenisers must minimise the size of rule  $S \rightarrow \mathbf{s}$  given a fixed number of rules  $s \rightarrow \mathbf{c}$ .

The paragraphs above highlights two important differences between tokenisers and compressors. First, tokenisers aim to reduce only the size of the resulting tokenised text (i.e.,  $|\mathbf{s}|$ ), whereas compressors also consider the size of the compression information (e.g., considering the size required to store  $\mathcal{S}$ , which would be  $\sum_{s \in \mathcal{S}} |\det(s)|$ . This is because tokenisers must create shorter inputs for NLP algorithms, while compressors must make information compact. Second, tokenisers and compressors have different optimisation parameters. Compression algorithms always compress a string to the best extent possible (e.g., for SLPs, until a single nonterminal is reached), whereas tokenisation algorithms are given a maximum vocabulary size (i.e., K) and find tokenisers which only compress their input as much as possible until this limit is reached.

# 8 Conclusion

In this work, we proved the NP-completeness of two variants of tokenisation. These results underline that finding optimal tokenisers most likely will remain a difficult quest and that research should focus on approximate algorithms instead. Regarding those, there is potential both in improving the analysis of currently used algorithms, such as BPE, as well as in designing new ones.

<sup>&</sup>lt;sup>10</sup>See Lohrey (2012) for an overview of straight-line programs, and Kempa and Prezza (2018); Kociumaka et al. (2023) for a more detailed overview of compression in general.

<sup>&</sup>lt;sup>11</sup>Although not originally defined that way, SLP's grammars are typically assumed to be in CNF, for simplicity. This does not make a big difference for compression, but will be important for our purposes.

# 675 Limitations

While we prove the NP-completeness of multi-676 ple variants of the tokenisation problem-which is 677 an important part of modern language modelling 678 pipelines-we must note a few limitations in our 679 work here. First, we only prove NP-completeness of tokenisation with compression as its objective. This is an important and popular objective function, frequently used to judge the quality of tokenisers; however, it is not perfectly correlated with downstream language modelling performance, as discussed in §2. Second, our proofs do not assume a fixed alphabet size, so for fixed alphabets tokenisation might not be NP-complete. Tokenisers are frequently run at the byte level, for which efficient algorithms might exist. Finally, while we investigated the complexity of two variants of the 691 tokenisation problem, similar results for other variants (e.g., with other tokenisation functions) remain open; this would be exciting future work.

## References

704

705

712

714

715

716

717

718

721

722

724

725

727

728

- Orevaoghene Ahia, Sachin Kumar, Hila Gonen, Jungo Kasai, David Mortensen, Noah Smith, and Yulia Tsvetkov. 2023. Do all languages cost the same? Tokenization in the era of commercial language models. In *Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing*, pages 9904–9923, Singapore. Association for Computational Linguistics.
  - Mehdi Ali, Michael Fromm, Klaudia Thellmann, Richard Rutmann, Max Lübbering, Johannes Leveling, Katrin Klug, Jan Ebert, Niclas Doll, Jasper Buschhoff, Charvi Jain, Alexander Weber, Lena Jurkschat, Hammam Abdelwahab, Chelsea John, Pedro Ortiz Suarez, Malte Ostendorff, Samuel Weinbach, Rafet Sifa, Stefan Kesselheim, and Nicolas Flores-Herr. 2024. Tokenizer choice for LLM training: Negligible or crucial? In *Findings of the Association* for Computational Linguistics: NAACL 2024, pages 3907–3924, Mexico City, Mexico. Association for Computational Linguistics.
- 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.
- M. Charikar, E. Lehman, Ding Liu, R. Panigrahy, M. Prabhakaran, A. Sahai, and A. Shelat. 2005. The smallest grammar problem. *IEEE Transactions on Information Theory*, 51(7):2554–2576.

Marco Cognetta, Vilém Zouhar, Sangwhan Moon, and Naoaki Okazaki. 2024. Two counterexamples to tokenization and the noiseless channel. In *Proceedings* of the 2024 Joint International Conference on Computational Linguistics, Language Resources and Evaluation (LREC-COLING 2024), pages 16897–16906, Torino, Italia. ELRA and ICCL. 729

730

731

732

733

736

737

738

739

740

741

742

743

744

745

746

747

749

750

751

752

753

754

755

756

757

758

759

760

761

762

763

764

765

766

767

768

769

770

771

772

773

774

776

777

778

779

780

781

782

783

784

785

- Philip Gage. 1994. A new algorithm for data compression. *C Users Journal*, 12(2):23–38.
- 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.
- M. R. Garey, D. S. Johnson, and L. Stockmeyer. 1974. Some simplified NP-complete problems. In *Proceedings of the Sixth Annual ACM Symposium on Theory of Computing*, STOC '74, page 47–63, New York, NY, USA. Association for Computing Machinery.
- Renato Geh, Honghua Zhang, Kareem Ahmed, Benjie Wang, and Guy Van Den Broeck. 2024. Where is the signal in tokenization space? In *Proceedings* of the 2024 Conference on Empirical Methods in Natural Language Processing, pages 3966–3979, Miami, Florida, USA. Association for Computational Linguistics.
- 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.
- Omer Goldman, Avi Caciularu, Matan Eyal, Kris Cao, Idan Szpektor, and Reut Tsarfaty. 2024. Unpacking tokenization: Evaluating text compression and its correlation with model performance. In *Findings of the Association for Computational Linguistics: ACL* 2024, pages 2274–2286, Bangkok, Thailand. Association for Computational Linguistics.
- Thamme Gowda and Jonathan May. 2020. Finding the optimal vocabulary size for neural machine translation. In *Findings of the Association for Computational Linguistics: EMNLP 2020*, pages 3955–3964, Online. Association for Computational Linguistics.
- Dominik Kempa and Nicola Prezza. 2018. At the roots of dictionary compression: string attractors. In *Proceedings of the 50th Annual ACM SIGACT Symposium on Theory of Computing*, STOC 2018, page 827–840, New York, NY, USA. Association for Computing Machinery.
- J.C. Kieffer and En-Hui Yang. 2000. Grammar-based codes: a new class of universal lossless source codes. *IEEE Transactions on Information Theory*, 46(3):737–754.

- 787 792 795 796 797 800 801

- 806
- 810 811 812
- 813
- 814 815
- 816 817
- 818 819 820 821

833

834 835

836

837

839

841

842

843

- Tomasz Kociumaka, Gonzalo Navarro, and Nicola Prezza. 2023. Toward a definitive compressibility measure for repetitive sequences. IEEE Transactions on Information Theory, 69(4):2074–2092.
  - László Kozma and Johannes Voderholzer. 2024. Theoretical analysis of byte-pair encoding. Preprint, arXiv:2411.08671.
  - Taku Kudo. 2018. Subword regularization: Improving neural network translation models with multiple subword candidates. In Proceedings of the 56th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers), pages 66–75, Melbourne, Australia. Association for Computational Linguistics.
  - Jia Peng Lim, Davin Choo, and Hady W. Lauw. 2025. A partition cover approach to tokenization. Preprint, arXiv:2501.06246.
  - Markus Lohrey. 2012. Algorithmics on SLPcompressed strings: A survey. Groups - Complexity -Cryptology, 4(2):241-299.
  - Matteo Muffo, Aldo Cocco, and Enrico Bertino. 2022. Evaluating transformer language models on arithmetic operations using number decomposition. In Proceedings of the Thirteenth Language Resources and Evaluation Conference, pages 291-297, Marseille, France. European Language Resources Association.
  - Rodrigo Nogueira, Zhiying Jiang, and Jimmy Lin. 2021. Investigating the limitations of transformers with simple arithmetic tasks. Preprint, arXiv:2102.13019.
  - Byung-Doh Oh and William Schuler. 2024. Leading whitespaces of language models' subword vocabulary pose a confound for calculating word probabilities. In Proceedings of the 2024 Conference on Empirical Methods in Natural Language Processing, pages 3464-3472, Miami, Florida, USA. Association for Computational Linguistics.
- OpenAI, Josh Achiam, Steven Adler, Sandhini Agarwal, Lama Ahmad, Ilge Akkaya, Florencia Leoni Aleman, Diogo Almeida, Janko Altenschmidt, Sam Altman, Shyamal Anadkat, Red Avila, Igor Babuschkin, Suchir Balaji, Valerie Balcom, Paul Baltescu, Haiming Bao, Mohammad Bavarian, Jeff Belgum, Irwan Bello, Jake Berdine, Gabriel Bernadett-Shapiro, Christopher Berner, Lenny Bogdonoff, Oleg Boiko, Madelaine Boyd, Anna-Luisa Brakman, Greg Brockman, Tim Brooks, Miles Brundage, Kevin Button, Trevor Cai, Rosie Campbell, Andrew Cann, Brittany Carey, Chelsea Carlson, Rory Carmichael, Brooke Chan, Che Chang, Fotis Chantzis, Derek Chen, Sully Chen, Ruby Chen, Jason Chen, Mark Chen, Ben Chess, Chester Cho, Casey Chu, Hyung Won Chung, Dave Cummings, Jeremiah Currier, Yunxing Dai, Cory Decareaux, Thomas Degry, Noah Deutsch, Damien Deville, Arka Dhar, David Dohan, Steve Dowling, Sheila Dunning, Adrien Ecoffet, Atty Eleti, Tyna Eloundou, David Farhi, Liam Fedus, Niko Felix,

Simón Posada Fishman, Juston Forte, Isabella Ful-844 ford, Leo Gao, Elie Georges, Christian Gibson, Vik 845 Goel, Tarun Gogineni, Gabriel Goh, Rapha Gontijo-846 Lopes, Jonathan Gordon, Morgan Grafstein, Scott 847 Gray, Ryan Greene, Joshua Gross, Shixiang Shane 848 Gu, Yufei Guo, Chris Hallacy, Jesse Han, Jeff Harris, 849 Yuchen He, Mike Heaton, Johannes Heidecke, Chris 850 Hesse, Alan Hickey, Wade Hickey, Peter Hoeschele, 851 Brandon Houghton, Kenny Hsu, Shengli Hu, Xin 852 Hu, Joost Huizinga, Shantanu Jain, Shawn Jain, 853 Joanne Jang, Angela Jiang, Roger Jiang, Haozhun 854 Jin, Denny Jin, Shino Jomoto, Billie Jonn, Hee-855 woo Jun, Tomer Kaftan, Łukasz Kaiser, Ali Ka-856 mali, Ingmar Kanitscheider, Nitish Shirish Keskar, 857 Tabarak Khan, Logan Kilpatrick, Jong Wook Kim, 858 Christina Kim, Yongjik Kim, Jan Hendrik Kirch-859 ner, Jamie Kiros, Matt Knight, Daniel Kokotajlo, 860 Łukasz Kondraciuk, Andrew Kondrich, Aris Kon-861 stantinidis, Kyle Kosic, Gretchen Krueger, Vishal 862 Kuo, Michael Lampe, Ikai Lan, Teddy Lee, Jan 863 Leike, Jade Leung, Daniel Levy, Chak Ming Li, 864 Rachel Lim, Molly Lin, Stephanie Lin, Mateusz 865 Litwin, Theresa Lopez, Ryan Lowe, Patricia Lue, Anna Makanju, Kim Malfacini, Sam Manning, Todor Markov, Yaniv Markovski, Bianca Martin, Katie 868 Mayer, Andrew Mayne, Bob McGrew, Scott Mayer 869 McKinney, Christine McLeavey, Paul McMillan, 870 Jake McNeil, David Medina, Aalok Mehta, Jacob 871 Menick, Luke Metz, Andrey Mishchenko, Pamela 872 Mishkin, Vinnie Monaco, Evan Morikawa, Daniel 873 Mossing, Tong Mu, Mira Murati, Oleg Murk, David 874 Mély, Ashvin Nair, Reiichiro Nakano, Rajeev Nayak, 875 Arvind Neelakantan, Richard Ngo, Hyeonwoo Noh, 876 Long Ouyang, Cullen O'Keefe, Jakub Pachocki, Alex 877 Paino, Joe Palermo, Ashley Pantuliano, Giambat-878 tista Parascandolo, Joel Parish, Emy Parparita, Alex 879 Passos, Mikhail Pavlov, Andrew Peng, Adam Perel-880 man, Filipe de Avila Belbute Peres, Michael Petrov, 881 Henrique Ponde de Oliveira Pinto, Michael, Poko-882 rny, Michelle Pokrass, Vitchyr H. Pong, Tolly Pow-883 ell, Alethea Power, Boris Power, Elizabeth Proehl, 884 Raul Puri, Alec Radford, Jack Rae, Aditya Ramesh, 885 Cameron Raymond, Francis Real, Kendra Rimbach, 886 Carl Ross, Bob Rotsted, Henri Roussez, Nick Ryder, Mario Saltarelli, Ted Sanders, Shibani Santurkar, 888 Girish Sastry, Heather Schmidt, David Schnurr, John 889 Schulman, Daniel Selsam, Kyla Sheppard, Toki 890 Sherbakov, Jessica Shieh, Sarah Shoker, Pranav 891 Shyam, Szymon Sidor, Eric Sigler, Maddie Simens, 892 Jordan Sitkin, Katarina Slama, Ian Sohl, Benjamin 893 Sokolowsky, Yang Song, Natalie Staudacher, Fe-894 lipe Petroski Such, Natalie Summers, Ilya Sutskever, 895 Jie Tang, Nikolas Tezak, Madeleine B. Thompson, 896 Phil Tillet, Amin Tootoonchian, Elizabeth Tseng, 897 Preston Tuggle, Nick Turley, Jerry Tworek, Juan Fe-898 lipe Cerón Uribe, Andrea Vallone, Arun Vijayvergiya, 899 Chelsea Voss, Carroll Wainwright, Justin Jay Wang, 900 Alvin Wang, Ben Wang, Jonathan Ward, Jason Wei, 901 CJ Weinmann, Akila Welihinda, Peter Welinder, Ji-902 ayi Weng, Lilian Weng, Matt Wiethoff, Dave Willner, 903 Clemens Winter, Samuel Wolrich, Hannah Wong, 904 Lauren Workman, Sherwin Wu, Jeff Wu, Michael 905 Wu, Kai Xiao, Tao Xu, Sarah Yoo, Kevin Yu, 906 Qiming Yuan, Wojciech Zaremba, Rowan Zellers, 907

Chong Zhang, Marvin Zhang, Shengjia Zhao, Tianhao Zheng, Juntang Zhuang, William Zhuk, and Barret Zoph. 2024. GPT-4 technical report. *Preprint*, arXiv:2303.08774.

908

909

910

911

912

913

914

915

917

918

919

920

921

924

925

928

929

930

931

933

937

941

942

943

952

953

955

956

957

960

961

962

963 964

- Artidoro Pagnoni, Ram Pasunuru, Pedro Rodriguez, John Nguyen, Benjamin Muller, Margaret Li, Chunting Zhou, Lili Yu, Jason Weston, Luke Zettlemoyer, Gargi Ghosh, Mike Lewis, Ari Holtzman, and Srinivasan Iyer. 2024. Byte latent transformer: Patches scale better than tokens. *Preprint*, arXiv:2412.09871.
- Aleksandar Petrov, Emanuele La Malfa, Philip Torr, and Adel Bibi. 2023. Language model tokenizers introduce unfairness between languages. In *Thirtyseventh Conference on Neural Information Processing Systems.*
- Buu Phan, Marton Havasi, Matthew Muckley, and Karen Ullrich. 2024. Understanding and mitigating tokenization bias in language models. *Preprint*, arXiv:2406.16829.
- Tiago Pimentel and Clara Meister. 2024. How to compute the probability of a word. In *Proceedings of the* 2024 Conference on Empirical Methods in Natural Language Processing, Miami, Florida, USA. Association for Computational Linguistics.
- Phillip Rust, Jonas Pfeiffer, Ivan Vulić, Sebastian Ruder, and Iryna Gurevych. 2021. How good is your tokenizer? On the monolingual performance of multilingual language models. In Proceedings of the 59th Annual Meeting of the Association for Computational Linguistics and the 11th International Joint Conference on Natural Language Processing (Volume 1: Long Papers), pages 3118–3135, Online. Association for Computational Linguistics.
- Craig W. Schmidt, Varshini Reddy, Haoran Zhang, Alec Alameddine, Omri Uzan, Yuval Pinter, and Chris Tanner. 2024. Tokenization is more than compression. In *Proceedings of the 2024 Conference on Empirical Methods in Natural Language Processing*, pages 678–702, Miami, Florida, USA. Association for Computational Linguistics.
- 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.
- Aaditya K. Singh and DJ Strouse. 2024. Tokenization counts: the impact of tokenization on arithmetic in frontier LLMs. *Preprint*, arXiv:2402.14903.
- 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. 2023a. LLaMA: Open and efficient foundation language models. *Preprint*, arXiv:2302.13971.

Hugo Touvron, Louis Martin, Kevin Stone, Peter Albert, Amjad Almahairi, Yasmine Babaei, Nikolay Bashlykov, Soumya Batra, Prajjwal Bhargava, Shruti Bhosale, Dan Bikel, Lukas Blecher, Cristian Canton Ferrer, Moya Chen, Guillem Cucurull, David Esiobu, Jude Fernandes, Jeremy Fu, Wenyin Fu, Brian Fuller, Cynthia Gao, Vedanuj Goswami, Naman Goyal, Anthony Hartshorn, Saghar Hosseini, Rui Hou, Hakan Inan, Marcin Kardas, Viktor Kerkez, Madian Khabsa, Isabel Kloumann, Artem Korenev, Punit Singh Koura, Marie-Anne Lachaux, Thibaut Lavril, Jenya Lee, Diana Liskovich, Yinghai Lu, Yuning Mao, Xavier Martinet, Todor Mihaylov, Pushkar Mishra, Igor Molybog, Yixin Nie, Andrew Poulton, Jeremy Reizenstein, Rashi Rungta, Kalyan Saladi, Alan Schelten, Ruan Silva, Eric Michael Smith, Ranjan Subramanian, Xiaoqing Ellen Tan, Binh Tang, Ross Taylor, Adina Williams, Jian Xiang Kuan, Puxin Xu, Zheng Yan, Iliyan Zarov, Yuchen Zhang, Angela Fan, Melanie Kambadur, Sharan Narang, Aurelien Rodriguez, Robert Stojnic, Sergey Edunov, and Thomas Scialom. 2023b. Llama 2: Open foundation and fine-tuned chat models. Preprint, arXiv:2307.09288. 965

966

967

968

969

970

971

972

973

974

975

976

977

978

979

980

981

982

983

984

985

986

987

988

989

990

991

992

993

994

995

996

997

998

999

1000

1001

1002

- Lili Yu, Daniel Simig, Colin Flaherty, Armen Aghajanyan, Luke Zettlemoyer, and Mike Lewis. 2023. MEGABYTE: Predicting million-byte sequences with multiscale transformers. In *Thirty-seventh Conference on Neural Information Processing Systems*.
- Vilém Zouhar, Clara Meister, Juan Gastaldi, Li Du, Mrinmaya Sachan, and Ryan Cotterell. 2023a. 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.
- Vilém Zouhar, Clara Meister, Juan Gastaldi, Li Du, Tim Vieira, Mrinmaya Sachan, and Ryan Cotterell. 2023b. A formal perspective on byte-pair encoding. In *Findings of the Association for Computational Linguistics: ACL 2023*, pages 598–614, Toronto, Canada. Association for Computational Linguistics.

## A Proof of Lemma 2

1006

1008

1009

1010

1011

1022

1025

1028

1029

1030

**Lemma 2.** If a max-2-SAT instance is satisfiable, then the direct tokenisation instance output by Reduction 1 is also satisfiable. Formally:

$$M2S(\mathcal{X}, \mathcal{L}, \psi) \implies Tok_{\phi}(R1(\mathcal{X}, \mathcal{L}, \psi))$$
(14)

*Proof.* First, note that if  $M2S(\mathcal{X}, \mathcal{L}, \psi)$ , then we have that Eq. (7) holds:

$$\psi \le \max_{\chi \in \{\mathbf{F}, \mathbf{T}\}^J} \sum_{i=1}^{I} \mathbb{1}\{L_i^1 \lor L_i^2\}$$
(29)

1012 Now, without loss of generality, let a satisfying solution have values  $x_j^*$ . In this case, for each variable  $X_j$ , 1013 we construct token  $\odot x_j^{\mathrm{T}} \odot$  if  $x_j^*$  is true, or  $\odot x_j^{\mathrm{F}} \odot$  if  $x_j^*$  is false. This gives us a total of J new tokens, so 1014 satisfies the  $|\mathcal{S}| = |\Sigma| + K$  condition. Now we just need to count the symbols output by this solution 1015 to see if Eq. (9) is satisfied, since any given tokenisation  $tok(\cdot, \mathcal{S})$  will provide an upper bound on the 1016 optimal tokenisation in terms of compression:

1017  

$$\sum_{\mathbf{c}\in\mathcal{D}} |\mathsf{tok}_{\diamond}[\mathcal{S}](\mathbf{c})| \ge \min_{\mathcal{S}'\subset\Sigma^+} \sum_{\mathbf{c}\in\mathcal{D}} |\mathsf{tok}_{\diamond}[\mathcal{S}'](\mathbf{c})| \qquad (30)$$
1018  
s.t.  $|\mathcal{S}'| = |\Sigma| + K$ 

For each pair of strings  $\odot x_j^{\mathsf{T}} \odot$  and  $\odot x_j^{\mathsf{F}} \odot$  in  $\mathcal{D}_1$ , one is compressed into a single subword while the other is kept as originally—using 3 symbols. We thus have that the strings in  $\mathcal{D}_1$  will occupy a total of (1+3)Jcharacters, and:

$$\sum_{\mathbf{c}\in(\bigcup_{i=1}^{f}\mathcal{D}_{1})}|\mathsf{tok}_{\diamond}[\boldsymbol{\mathcal{S}}](\mathbf{c})| = 4fJ$$
(31)

1023 Similarly, for each string in  $\mathcal{D}_2$  of form  $\odot x_j^{\mathsf{T}} \odot x_j^{\mathsf{F}} \odot$ , we have that either token  $\odot x_j^{\mathsf{T}} \odot$  or  $\odot x_j^{\mathsf{F}} \odot$  exists. So 1024 each of these strings is compressed from 5 into 3 symbols. We thus have:

$$\sum_{\mathbf{c} \in (\bigcup_{j=1}^{f'} \mathcal{D}_2)} |\mathsf{tok}_{\diamond}[\mathcal{S}](\mathbf{c})| = 3f'J$$
(32)

Finally, we have strings in  $\mathcal{D}_3$  of form  $\odot L_i^1 \odot L_i^2 \odot$ . These strings will be compressed into 3 symbols if  $\odot L_i^1 \odot$  or  $\odot L_i^2 \odot$  (or both) exist, and kept with 5 symbols otherwise. We thus have:

$$\sum_{\mathbf{c}\in\mathcal{D}_{3}}|\mathsf{tok}_{\diamond}[\mathcal{S}](\mathbf{c})| = \sum_{\mathbf{c}\in\mathcal{D}_{3}}\left(5-2\,\mathbb{1}\left\{\begin{array}{c} \odot L_{i}^{1}\odot\in\mathcal{S}\\ \text{or}\\ \odot L_{i}^{2}\odot\in\mathcal{S} \end{array}\right\}\right) \tag{33a}$$

$$=5I-2\sum_{\mathbf{c}\in\mathcal{D}_{3}}\mathbb{1}\left\{\begin{array}{c}\otimes L_{i}^{*}\otimes\in\mathcal{S}\\ \text{or}\\\otimes L_{i}^{2}\otimes\in\mathcal{S}\end{array}\right\}$$
(33b)

$$=5I - 2\sum_{i=1}^{I} \mathbb{1}\{L_i^1 \vee L_i^2\}$$
(33c)

1031 
$$< 5I - 2\psi$$

where, by construction, we have that a subword  $\odot L_i \odot \in S$  if and only if its associated variable  $(x_j \text{ or } \neg x_j)$  is true. Summing together the lengths in Eqs. (31) to (33), we get that

034 
$$\sum_{\mathbf{c}\in\mathcal{D}}|\mathsf{tok}_{\mathfrak{F}}[\boldsymbol{\mathcal{S}}](\mathbf{c})| \leq \delta = (4f+3f')J+5I-2\psi$$
(34)

1035 which concludes the proof.

(33d)

## B Proof of Lemma 3

**Lemma 3.** If the direct tokenisation instance output by Reduction 1 is satisfiable, the max-2-SAT instance 1037 which generated it is as well. Formally:

$$\operatorname{Tok}_{\diamond}(\operatorname{R1}(\mathcal{X},\mathcal{L},\psi)) \Longrightarrow \operatorname{M2S}(\mathcal{X},\mathcal{L},\psi)$$
 (15) 103

*Proof.* First, note that the dataset  $\mathcal{D}$  output by Reduction 1 has a total number of characters:

$$\sum_{\mathbf{c}\in\mathcal{D}}|\mathbf{c}| = (6f + 5f')J + 5I \tag{35}$$
 1041

Further, let:

$$\mathsf{toklen}(\mathcal{D}, \mathcal{S}) \stackrel{\text{def}}{=} \sum_{\mathbf{c} \in \mathcal{D}} |\mathsf{tok}_{\diamond}[\mathcal{S}](\mathbf{c})|, \qquad \qquad \mathcal{S}_0 = \bigcup_{j=1}^J \{ \odot x_j^{\mathsf{T}} \odot, \odot x_j^{\mathsf{F}} \odot \}$$
(36) 104

The maximum number of symbols in this dataset after compression is set to  $\delta = (4f + 3f') J + 5I - 2\psi$ . 1044 This means that, to satisfy this objective, there must exist a vocabulary whose tokeniser compresses the text by at least  $(2f + 2f') J + 2\psi$  symbols. We now prove this lemma in three steps: (1) we show that 1046 any solution which compresses the text by at least 2fJ symbols must only have subwords of the form 1047  $\odot x_j^T \odot$  or  $\odot x_j^F \odot$ ; (2) we show that any solution which compresses the text by at least (2f + 2f')J symbols 1048 must only have either subword  $\odot x_j^T \odot$  or  $\odot x_j^F \odot$  for any variable  $X_j$ ; (3) we show that any solution which compresses the text by at least  $(2f + 2f')J + 2\psi$  symbols must be produced by a max-2-SAT instance 1050 which has at least  $\psi$  clauses that are simultaneously satisfiable.

**LemmaProofStep 1.** (Step (1)). Any solution which compresses the text by at least 2fJ symbols must only have nontrivial subwords<sup>12</sup> of the form  $\odot x_i^T \odot$  or  $\odot x_i^F \odot$ , i.e.,: 1053

$$\left(\operatorname{toklen}(\mathcal{D}, \mathcal{S}) \leq \underbrace{(4f + 5f')J + 5I}_{\sum_{\mathbf{c} \in \mathcal{D}} |\mathbf{c}| - 2fJ}\right) \implies \mathcal{S} \subset \mathcal{S}_{0}$$
(37) 1054

*Proof.* First, given a solution with  $S \subset S_0$ , each subword  $s \in S$  will replace at least f strings in  $\mathcal{D}_1$ —i.e., with form  $\otimes x_j^{\mathsf{T}} \otimes \operatorname{or} \otimes x_j^{\mathsf{F}} \otimes$ —for a single subword, thus saving 2f characters. Since we have |S| = K = Jtokens, we save exactly 2fJ symbols:

$$\mathcal{S} \subset \mathcal{S}_0 \implies \left( \operatorname{toklen}(\mathcal{D}_1, \mathcal{S}') = \underbrace{4fJ}_{\sum_{c \in \mathcal{D}_1} |c| - 2fJ} \right)$$
(38) 105

Note now that any solution S' for which  $S' \not\subset S_0$  has at least one subword which is not of the form  $\odot x_j^{\mathrm{T}} \odot$ or  $\odot x_j^{\mathrm{F}} \odot$ ; this subword  $s \notin S_0$  will thus not compress strings in  $\mathcal{D}_1$  by 2f symbols, but by at most f:

$$\mathcal{S}' \not\subset \mathcal{S}_0 \implies \left( \operatorname{toklen}(\mathcal{D}_1, \mathcal{S}') \ge \underbrace{4f(J-1) + 5f}_{\sum_{\mathbf{c} \in \mathcal{D}_1} |\mathbf{c}| - 2fJ + f} \right)$$
(39)

Even if this new subword were able to fully compress strings in  $D_2$  and  $D_3$  to a single symbol each, it would reach a compression of at most 4f'J + 4I. Since by design f = 4f'J + 4I + 1, we get that:

$$\mathcal{S}' \not\subset \mathcal{S}_0 \implies \left( \operatorname{toklen}(\mathcal{D}, \mathcal{S}') \ge 4fJ + f + f'J + I > (4f + 5f')J + 5I \right)$$
(40) 100

which concludes this step of the proof.

1065

 $\square$ 

1059

1036

1040

<sup>&</sup>lt;sup>12</sup>We define nontrivial subwords as subwords with more than one character. Remember that by definition  $\Sigma \subseteq S$ , so all characters are always included in tokenisers' vocabularies. Also note that  $|S| = |\Sigma| + K$ , so those trivial subwords are not counted towards vocabulary size K.

**LemmaProofStep 2.** (Step (2)). Any solution which compresses the text by at least (2f + 2f')J symbols must only have either subword  $\odot x_j^{T} \odot$  or  $\odot x_j^{F} \odot$  for any variable  $X_j$ , i.e.,:

1068 
$$\left( \mathsf{toklen}(\mathcal{D}, \mathcal{S}) \leq \underbrace{(4f + 3f')J + 5I}_{\sum_{c \in \mathcal{D}} |c| - (2f + 2f')J} \right) \implies \forall_{j \in \{1, \dots, J\}} |\mathcal{S} \cap \{ \odot x_j^{\mathsf{T}} \odot, \odot x_j^{\mathsf{F}} \odot \}| = 1$$
(41)

1069 Proof. In this step of the proof, we show that satisfying solutions must have one and only one of subwords 1070  $\otimes x_j^T \otimes$  and  $\otimes x_j^F \otimes$  for any variable  $X_j$ . As before, it's easy to see that a solution of the form described 1071 achieves at least (2f + 2f')J symbol compression. Each subword of form  $\otimes x_j^T \otimes$  or  $\otimes x_j^F \otimes$  saves exactly 1072 2f characters in the strings in  $\mathcal{D}_1$ . Further, because we always have either subword  $\otimes x_j^T \otimes$  or  $\otimes x_j^F \otimes$  for 1073 each value of j, we also get 2f' compression in the strings in  $\mathcal{D}_2$ :

1074 
$$\forall_{j \in \{1,...,J\}} | \mathcal{S} \cap \{ \odot x_{j}^{\mathsf{T}} \odot, \odot x_{j}^{\mathsf{F}} \odot \} | = 1$$
1075 
$$\Longrightarrow \left( \operatorname{toklen}(\mathcal{D}_{1}, \mathcal{S}) = \underbrace{4fJ}_{\sum_{c \in \mathcal{D}_{1}} |c| - 2fJ} \right) \operatorname{and} \left( \operatorname{toklen}(\mathcal{D}_{2}, \mathcal{S}) = \underbrace{3f'J}_{\sum_{c \in \mathcal{D}_{2}} |c| - 2f'J} \right)$$
(42)

1076 Now note that this is not true if both  $\odot x_j^T \odot$  and  $\odot x_j^F \odot$  exist for a single j; in this case, only one of 1077 the tokens can be applied to  $\odot x_j^T \odot x_j^F \odot$ , and thus both tokens together lead to a benefit of 2 instead 1078 of 4. If both  $\odot x_j^T \odot$  and  $\odot x_j^F \odot$  exist for any token  $X_j$ , this implies that neither of  $\odot x_{j'}^T \odot$  and  $\odot x_{j'}^F \odot$ 1079 exists for some other  $X_{j'}$ , resulting in an uncompressed string. Then, we get at most a compression of 1080 2fJ + 2f'(J - 1) + 4I:

$$\exists_{j \in \{1,\dots,J\}} |\mathcal{S}' \cap \{ \odot x_j^{\mathsf{T}} \odot, \odot x_j^{\mathsf{F}} \odot \}| \neq 1 \implies \left( \mathsf{toklen}(\mathcal{D}, \mathcal{S}') \ge \underbrace{(4f + 3f')J + f' + I}_{\sum_{\mathbf{c} \in \mathcal{D}} |\mathbf{c}| - (2f + 2f)'J + f' - 4I} \right)$$
(43)

By construction f' = 4I + 1, which leads to:

$$\exists_{j \in \{1,\dots,J\}} |\mathcal{S}' \cap \{ \odot x_j^{\mathsf{T}} \odot, \odot x_j^{\mathsf{F}} \odot \}| \neq 1 \implies \left( \mathsf{toklen}(\mathcal{D}, \mathcal{S}') > (4f + 3f')J + 5I \right)$$
(44)

1084 This concludes the proof.

1083

1092

**LemmaProofStep 3.** (Step 3). Any instance of the tokenisation problem with a solution which compresses the text by at least  $(2f + 2f')J + 2\psi$  symbols must be produced by a max-2-SAT problem with at least  $\psi$  satisfied clauses, i.e.,:

$$\left( \operatorname{toklen}(\mathcal{D}, \mathcal{S}) \leq \underbrace{(4f + 3f')J + 5I - 2\psi}_{\sum_{c \in \mathcal{D}} |c| - (2f + 2f')J + 2\psi} \right) \implies \operatorname{M2S}(\mathcal{X}, \mathcal{L}, \psi)$$
(45)

1089 *Proof.* Finally, we now know any solution with this compression must have—for any variable  $X_j$ — 1090 either subword  $\odot x_j^T \odot$  or  $\odot x_j^F \odot$ . We can thus create a bijection  $\operatorname{Conv}_{S \to \chi}$  between the set of possible 1091 vocabularies respecting this condition, and the set of T/F assignments to SAT variables  $\chi$ :

$$\operatorname{Conv}_{\mathcal{S}\to\mathcal{X}}(\mathcal{S}) = \left\{ \begin{array}{ll} \mathsf{T} & \operatorname{if} \otimes x_j^{\mathsf{T}} \otimes \in \mathcal{S} \\ \mathsf{F} & \operatorname{if} \otimes x_j^{\mathsf{F}} \otimes \in \mathcal{S} \end{array} \right\}_{j=1}^{J}$$
(46)

Further, note that vocabularies of this form (as shown in Eq. (42)) lead to exactly (2f + 2f')J symbols being compressed in  $\mathcal{D}_1$  and  $\mathcal{D}_2$ . To achieve the target compression, a solution must thus compress  $\mathcal{D}_3$ by at least  $2\psi$  symbols. Now note that for any string  $\odot L_i^1 \odot L_i^2 \odot$  in  $\mathcal{D}_3$  we have three compression options:  $\odot L_i^1 \odot$  will be compressed, saving 2 symbols;  $\odot L_i^2 \odot$  will be compressed, also saving 2 symbols; or nothing will be compressed. More specifically,  $\odot L_i^1 \odot$  can be compressed if  $L_i^1$  represents  $X_j$ and subword  $\odot x_i^T \odot$  exists, or if  $L_i^1$  represents  $\neg X_j$  and subword  $\odot x_i^F \odot$  exists; the same is true for  $\odot L_i^2 \odot$ . They cannot both be compressed, however, as there is only one symbol  $\odot$  between the literals. 1099 We thus get a compression of 2 symbols for each of these strings if at least one of its literals has an 1100 associated subword in S. Note thus that whenever a string  $\odot L_i^1 \odot$  is compressed by 2 symbols using 1101 vocabulary  $\mathcal{S}$ , the max-2-SAT disjunction  $L_i^1 \vee L_i^2$  will also be satisfied by assignment  $\chi = \operatorname{Conv}_{\mathcal{S} \to \chi}(\mathcal{S})$ ; 1102 similarly, whenever this string suffers no compression (i.e., having a compression of zero), the max-2-SAT 1103 disjunction will not be satisfied. As our condition assumes a compression of at least  $2\psi$  symbols, we know 1104 that we have at least  $\psi$  strings for which a literal has an associated subword. We can thus write: 1105

$$2\psi \leq \max_{\mathcal{S} \subset \Sigma^*} \sum_{\mathbf{c} \in \mathcal{D}_3} |\mathbf{c}| - |\mathsf{tok}_{\diamond}[\mathcal{S}](\mathbf{c})| \tag{47a}$$

s.t. 
$$|\mathcal{S}| = J$$
 and  $\forall_{j \in \{1, \dots, J\}} |\mathcal{S} \cap \{ \odot x_j^{\mathsf{T}} \odot, \odot x_j^{\mathsf{F}} \odot \}| = 1$  1107

$$= \max_{\mathcal{S} \subset \Sigma^*} \sum_{\substack{\emptyset: L_i^1 \otimes L_i^2 \otimes \in \mathcal{D}_3}} 2 \, \mathbb{I} \left\{ \begin{array}{c} \otimes L_i^1 \otimes \in \mathcal{S} \\ \text{or} \\ \otimes L_i^2 \otimes \in \mathcal{S} \end{array} \right\}$$
(47b) 110

s.t. 
$$|\mathcal{S}| = J$$
 and  $\forall_{j \in \{1, \dots, J\}} |\mathcal{S} \cap \{ \odot x_j^{\mathsf{T}} \odot, \odot x_j^{\mathsf{F}} \odot \}| = 1$  1109

$$= \max_{\chi \in \{0,1\}^J} \sum_{i=1}^{I} 2\mathbb{1}\{L_i^1 \lor L_i^2\}$$
(47c) 1110

$$\Rightarrow M2S(\mathcal{X}, \mathcal{L}, \psi) \tag{47d}$$
111

1114

1117

1131

We thus know that, if a satisfying tokenisation solution exists, then the associated max-2-SAT problem 1112 will also be satisfiable. This concludes the proof. 1113

\_

#### С Proof of Lemma 6

**Lemma 6.** If a max-2-SAT instance is satisfiable, then the bottom-up tokenisation instance output by 1115 Reduction 2 is also satisfiable. Formally: 1116

$$M2S(\mathcal{X}, \mathcal{L}, \psi) \implies Tok_{\uparrow}(R2(\mathcal{X}, \mathcal{L}, \psi))$$
(23)

*Proof.* Our proof starts by first defining the three following lists of merges, which will be included in any 1118 satisfying solution to the tokenisation problem: 1119

$$\mathbf{m}_{1} = \bigcup_{j=1}^{J} \left[ \langle \otimes, x_{j}^{\mathsf{F}} \rangle, \langle x_{j}^{\mathsf{T}}, \otimes \rangle \right], \quad \mathbf{m}_{3} = \bigcup_{j=1}^{J} \left[ \langle x_{j}^{\mathsf{F}}, \otimes \rangle, \langle \otimes, x_{j}^{\mathsf{T}} \rangle \right], \quad \mathbf{m}_{5} = \bigcup_{j=1}^{J} \left[ \langle \otimes, x_{j}^{\mathsf{F}} \rangle, \langle x_{j}^{\mathsf{T}}, \otimes \rangle \right]$$
(48)

Now, without loss of generality, let a satisfying solution to max-2-SAT have values  $x_i^{\star}$ . We then construct 1121 two other lists of merges, which depend on this max-2-SAT solution: 1122

$$\mathbf{m}_{2} = \bigcup_{j=1}^{J} \begin{bmatrix} \langle \odot, x_{j}^{\mathsf{T}} \otimes \rangle & \text{if } x_{j}^{\star} = \mathsf{T} \\ \langle \otimes x_{j}^{\mathsf{F}}, \odot \rangle & \text{else} \end{bmatrix}, \qquad \mathbf{m}_{4} = \bigcup_{j=1}^{J} \begin{bmatrix} \langle \odot x_{j}^{\mathsf{T}}, \odot \rangle & \text{if } x_{j}^{\star} = \mathsf{T} \\ \langle \odot, x_{j}^{\mathsf{F}} \odot \rangle & \text{else} \end{bmatrix}$$
(49) 1123

in words, we create merges  $\langle \odot, x_j^{\mathsf{T}} \otimes \rangle$  and  $\langle \odot x_j^{\mathsf{T}}, \odot \rangle$  if  $x_j^{\star}$  is true, or  $\langle \otimes x_j^{\mathsf{F}}, \odot \rangle$  and  $\langle \odot, x_j^{\mathsf{F}} \odot \rangle$  if  $x_j^{\star}$  is false. 1124 We then create a merge sequence by concatenating these lists in order: 1125

$$\mathbf{m} = \mathbf{m}_1 \circ \mathbf{m}_2 \circ \mathbf{m}_3 \circ \mathbf{m}_4 \circ \mathbf{m}_5 \tag{50}$$

This gives us a total of  $|\mathbf{m}| = K = 8J$  merges. Now we just need to count the symbols output by this 1127 solution to see if Eq. (17) is satisfied, since any given tokenisation  $tok_{+}[m]$  will provide an upper bound 1128 on the optimal tokenisation in terms of compression: 1129

$$\sum_{\mathbf{c}\in\mathcal{D}} |\mathsf{tok}_{\uparrow}[\mathbf{m}](\mathbf{c})| \ge \min_{\mathbf{m}'\in\mathcal{M}^*} \sum_{\mathbf{c}\in\mathcal{D}} |\mathsf{tok}_{\uparrow}[\mathbf{m}'](\mathbf{c})|$$
(51) 113  
s.t.  $|\mathbf{m}'| = K$  113

С	$\mathtt{tok}_{\uparrow}[\mathbf{m}_1](\mathbf{c})$	$\mathtt{tok}_{\uparrow}[\mathbf{m}_{1}\circ\mathbf{m}_{2}\circ\mathbf{m}_{3}](\mathbf{c})$	$\mathtt{tok}_{\uparrow}[\mathbf{m}_{1}\circ\mathbf{m}_{2}\circ\mathbf{m}_{3}\circ\mathbf{m}_{4}\circ\mathbf{m}_{5}](\mathbf{c})$	$ \texttt{tok}_{\uparrow}[\mathbf{m}](\mathbf{c}) $
$\langle \odot, x_{j}^{\mathrm{T}} \rangle$		$\langle \odot x_{j}^{\mathrm{T}} \rangle$		1
$\langle x_{j}^{\mathrm{F}}, \odot \rangle$		$\langle x_{j}^{\mathbf{F}} \circ \rangle$		1
$\langle x_{i}^{\mathrm{T}}, \odot \rangle$		•	$\langle x_{i}^{\mathrm{T}}$ $\odot \rangle$	1
$\langle \odot, x_i^{\rm F} \rangle$			$\langle \odot x_{j}^{F} \rangle$	1
$\langle x_i^{\mathrm{T}}, \overset{\mathrm{J}}{\otimes} \rangle$	$\langle x_{j}^{\mathrm{T}}\otimes \rangle$			1
$\langle \otimes, x_j^{\rm F} \rangle$	$\langle \otimes x_{i}^{\mathbf{F}} \rangle$			1

Table 1: Example of applying **m** in  $\mathcal{D}_1$  of bottom-up tokenisation problem obtained from Reduction 2. The dot symbol  $\cdot$  denotes the string not changing under the given merge.

с	$\mathtt{tok}_{\uparrow}[\mathbf{m}_1](\mathbf{c})$	$\mathtt{tok}_{\uparrow}[\mathbf{m}_{1}\circ\mathbf{m}_{2}](\mathbf{c})$		$\mathtt{tok}_{\uparrow}[\mathbf{m}_{1}\circ\mathbf{m}_{2}\circ\mathbf{m}_{3}](\mathbf{c})$	$\mathtt{tok}_{\uparrow}[\mathbf{m}_{1}\circ\mathbf{m}_{2}\circ\mathbf{m}_{3}\circ\mathbf{m}_{4}](\mathbf{c})$		$ \texttt{tok}_{\uparrow}[\mathbf{m}](\mathbf{c}) $	
		$x_j^\star = \mathbf{T}$	$x_j^\star = \mathbf{F}$		$x_j^\star = \mathtt{T}$	$x_j^\star = \mathbf{F}$	$x_j^\star = \mathtt{T}$	$x_j^\star = \mathbf{F}$
$\langle \odot, x_i^{\mathrm{T}}, \odot \rangle$				$\langle \odot x_i^{T}, \odot \rangle$	$\langle \odot x_i^{\mathrm{T}} \odot \rangle$	$\langle \odot, x_i^{T} \odot \rangle$	1	2
$\langle \odot, x_i^{\check{\mathbf{F}}}, \odot \rangle$		•		$\langle \odot, x_i^{\mathbf{F}} \odot \rangle$	$\langle \odot, x_j^{\mathbf{F}} \odot \rangle$	$\langle \odot x_{i}^{\vec{\mathbf{F}}} \odot \rangle$	2	1
$\langle \odot, x_i^{\mathrm{T}}, \otimes \rangle$	$\langle \odot, x_i^{\mathrm{T}} \otimes \rangle$	$\langle \odot x_i^{T} \otimes \rangle$	$\langle \odot, x_i^{\mathrm{T}} \otimes \rangle$				1	2
$\langle \otimes, x_j^{\mathrm{F}}, \odot  angle$	$\langle \otimes x_j^{{f F}'}, \odot  angle$	$\langle \otimes, x_j^{\mathbf{F}} \otimes \rangle$	$\langle \otimes x_j^{ec{\mathbf{F}}} \odot  angle$				2	1

Table 2: Example of applying **m** in  $D_2$  of bottom-up tokenisation problem obtained from Reduction 2. The dot symbol  $\cdot$  denotes the string not changing under the given merge.

By applying the merges  $\mathbf{m}$ , each string in  $\mathcal{D}_1$  will be compressed into a single subword; note that  $\mathbf{m}_2$ and  $\mathbf{m}_4$  will have no effect on these strings. This is easy to see by applying merges sequentially to these strings, as displayed in Tab. 1. leading to:

1135 
$$\sum_{\mathbf{c}\in(\bigcup_{j=1}^{f}\mathcal{D}_{1})}|\mathsf{tok}_{\uparrow}[\mathbf{m}](\mathbf{c})| = 6fJ$$
(52)

1136

1142

For each pair of strings  $\odot x_j^{\mathrm{T}} \odot$  and  $\odot x_j^{\mathrm{F}} \odot$  in  $\mathcal{D}_2$ , one is compressed into a single subword while the other is only compressed to two subwords—the one with  $x_j^{\mathrm{T}}$  is compressed to a single symbol if  $x_j^{\star} = \mathrm{T}$ and the one with  $x_j^{\mathrm{F}}$  otherwise. The same is true for each pair of strings  $\odot x_j^{\mathrm{T}} \otimes$  and  $\otimes x_j^{\mathrm{F}} \odot$ , also in  $\mathcal{D}_2$ . This is displayed in Tab. 2. We thus have that, for each variable  $X_j$ , the strings in  $\mathcal{D}_2$  will occupy a total of (1 + 2 + 1 + 2)J characters, and:

$$\sum_{\mathbf{c}\in(\bigcup_{j=1}^{f}\mathcal{D}_{1})}|\mathsf{tok}_{\uparrow}[\mathbf{m}](\mathbf{c})| = 6f'J$$
(53)

Similarly, each string in  $\mathcal{D}_3$  and  $\mathcal{D}_4$  will be compressed into only 2 symbols after this tokeniser is applied to it. We thus have:

1145 
$$\sum_{\mathbf{c}\in(\bigcup_{i=1}^{f''}\mathcal{D}_3)} |\mathsf{tok}_{\uparrow}[\mathbf{m}](\mathbf{c})| = 4f''J, \qquad \sum_{\mathbf{c}\in(\bigcup_{i=1}^{f'''}\mathcal{D}_4)} |\mathsf{tok}_{\uparrow}[\mathbf{m}](\mathbf{c})| = 4f'''J \tag{54}$$

Finally, we have the strings in  $\mathcal{D}_5$ . These strings are constructed such that they will be compressed into 2 symbols if either  $L_i^1$  or  $L_i^2$  evaluates to T, and kept with 3 symbols otherwise; see Tab. 4 for a detailed

$\mathcal{D}$	с	$\mathtt{tok}_{\uparrow}[\mathbf{m}_1](\mathbf{c})$	$\mathtt{tok}_{\uparrow}[\mathbf{m}_{1}\circ\mathbf{m}_{2}](\mathbf{c})$		$\mathtt{tok}_{\uparrow}[\mathbf{m}_{1}\circ\mathbf{m}_{2}\circ\mathbf{m}_{3}](\mathbf{c})$		$\mathtt{tok}_{\uparrow}[\mathbf{m}_{1}\circ\mathbf{m}_{2}\circ\mathbf{m}_{3}\circ\mathbf{m}_{4}](\mathbf{c})$		$\mathtt{tok}_{\uparrow}[\mathbf{m}_{1}\circ\mathbf{m}_{2}\circ\mathbf{m}_{3}\circ\mathbf{m}_{5}](\mathbf{c})$		$ \texttt{tok}_{\uparrow}[\mathbf{m}](\mathbf{c}) $
			$x_j^{\star} = \mathbf{T}$	$x_j^\star = \mathbf{F}$	$x_j^{\star} = \mathbf{T}$	$x_j^\star = \mathbf{F}$	$x_j^{\star} = \mathbf{T}$	$x_j^\star = \mathbf{F}$	$x_j^\star = \mathtt{T}$	$x_j^{\star} = \mathbf{F}$	
$\mathcal{D}_3$	$\langle \odot, x_i^{T}, \odot, x_i^{F}, \odot \rangle$				$\langle \odot x_i^T, \odot \rangle$	$(a, x_i^F \otimes)$	$\langle \odot x_i^{T} \odot, x_i^{F} \odot \rangle$	$\langle \odot x_i^{\mathrm{T}}, \odot x_i^{\mathrm{F}} \odot \rangle$			2
$\mathcal{D}_3$	$\langle \otimes, x_i^{\mathbf{F}}, \odot, x_i^{\mathbf{T}}, \otimes \rangle$	$\langle \otimes x_i^{\mathbf{F}}, \odot, x_i^{\mathbf{T}} \otimes \rangle$	$\langle \otimes x_{i}^{\mathrm{F}}, \odot x_{i}^{\mathrm{T}} \otimes \rangle$	$\langle \otimes x_{i}^{F} \otimes, x_{i}^{T} \otimes \rangle$		· ·					2
$\mathcal{D}_4$	$\langle \odot, x_i^{\mathbf{F}}, \odot, x_i^{\mathbf{T}}, \otimes \rangle$	$\langle \odot, \tilde{x}_{i}^{\mathrm{F}}, \odot, \tilde{x}_{i}^{\mathrm{T}} \otimes \rangle$	$\langle \odot, x_j^{F}, \odot x_j^{T} \otimes \rangle$	•		$\langle \odot, x_j^{F} \odot, x_j^{T} \otimes \rangle$		$\langle \odot x_{i}^{F} \odot, x_{i}^{T} \otimes \rangle$	$\langle \odot x_{j}^{F}, \odot x_{j}^{T} \otimes \rangle$		2
$\mathcal{D}_4$	$\langle \otimes, x_j^{\rm F}, \odot, x_j^{\rm T}, \odot \rangle$	$\langle \otimes x_j^{\mathbf{F}}, \odot, x_j^{\mathbf{T}}, \odot \rangle$	· · ·	$\langle \otimes x_{j}^{\mathrm{F}} \circledcirc, x_{j}^{\mathrm{T}}, \circledcirc \rangle$	$\langle \otimes x_j^{\rm F}, \odot x_j^{\rm T}, \odot \rangle$	· · ·	$\langle \otimes x_j^{\rm F}, \odot x_j^{\rm T} \odot \rangle$	· · ·	· · ·	$\langle \otimes x_{j}^{\mathrm{F}} \circledcirc, x_{j}^{\mathrm{T}} \circledcirc \rangle$	2

Table 3: Example of applying **m** in  $\mathcal{D}_3$  and  $\mathcal{D}_4$  of the bottom-up tokenisation problem obtained from Reduction 2. The dot symbol  $\cdot$  denotes the string not changing under the given merge.

Assignment	Condition	с	$\mathtt{tok}_{\uparrow}[\mathbf{m}_1](\mathbf{c})$	$\mathtt{tok}_{\uparrow}[\mathbf{m}_{1}\circ\mathbf{m}_{2}](\mathbf{c})$	$\mathtt{tok}_{\uparrow}[\mathbf{m}_{1}\circ\mathbf{m}_{2}\circ\mathbf{m}_{3}](\mathbf{c})$	$\mathtt{tok}_{\uparrow}[\mathbf{m}_{1}\circ\mathbf{m}_{2}\circ\mathbf{m}_{3}\circ\mathbf{m}_{4}](\mathbf{c})$	$ \texttt{tok}_{\uparrow}[\textbf{m}](c) $
	$x_j^\star = \mathtt{T} \wedge x_{j'}^\star = \mathtt{T}$	$\big\langle \odot, x_j^{T}, \odot, x_{j'}^{F}, \odot \big\rangle$			$\langle \circledcirc x_j^{\mathtt{T}}, \circledcirc, x_{j'}^{\mathtt{F}} \circledcirc \rangle$	$\langle \odot x_{j}^{T} \odot, x_{j'}^{F} \odot \rangle$	2
$L^1_i = X_i$ and $L^2_i = \neg X_{i'}$	$x_j^\star = \mathbb{F} \wedge x_{j'}^\star = \mathbb{T}$					$\langle \odot x_{j}^{\mathrm{T}}, \odot, x_{j'}^{\mathrm{F}} \odot \rangle$	3
<i>i J i i J</i>	$x_j^\star = \mathbf{T} \wedge x_{j'}^\star = \mathbf{F}$					$\langle \odot x_j^{T} \odot, x_{j'}^{F} \odot \rangle$	2
	$x_j^\star = \mathbb{F} \wedge x_{j'}^\star = \mathbb{F}$		•	•		$\langle \odot x_j^{\mathrm{T}}, \odot x_{j'}^{\mathrm{F}} \odot \rangle$	2
	$x_j^\star = \mathtt{T} \wedge x_{j'}^\star = \mathtt{T}$	$\langle \circledcirc, x_{j'}^{T}, \circledcirc, x_j^{F}, \circledcirc \rangle$	•		$\langle \odot x_{j'}^{\mathbb{T}}, \odot, x_j^{\mathbb{F}} \odot \rangle$	$\langle \odot x_{j'}^{T} \odot, x_{j}^{F} \odot \rangle$	2
$L^1 = \neg X$ and $L^2 = X$	$x_j^{\star} = \mathbf{F} \wedge x_{j'}^{\star} = \mathbf{T}$					$\langle \odot x_{j'}^{T} \odot, x_{j}^{F} \odot \rangle$	2
$E_i = i A_j \text{ and } E_i = A_j$	$x_j^{\star} = \mathbf{T} \wedge x_{j'}^{\star} = \mathbf{F}$					$\langle \odot x_{j'}^{\dagger}, \odot, x_{j}^{F} \odot \rangle$	3
	$x_j^\star = \mathbb{F} \wedge x_{j'}^\star = \mathbb{F}$					$\langle \odot x_{j'}^{T}, \odot x_{j}^{F} \odot \rangle$	2
	$x_i^\star = \mathtt{T} \wedge x_{i'}^\star = \mathtt{T}$	$\langle \otimes, x^{\rm F}_j, \odot, x^{\rm F}_{j'}, \odot \rangle$	$\langle \otimes x_j^{\rm F}, \odot, x_{j'}^{\rm F}, \odot \rangle$		$\begin{array}{ll} \langle \otimes x_{j}^{\mathrm{F}}, \odot, x_{j'}^{\mathrm{F}}, \odot \rangle \\ \langle \otimes x_{j}^{\mathrm{F}} \odot, x_{j'}^{\mathrm{F}}, \odot \rangle \end{array} \\ \begin{array}{l} \langle \otimes x_{j}^{\mathrm{F}} \odot, x_{j'}^{\mathrm{F}}, \odot \rangle \end{array}$		3
$I^1 = \neg Y$ and $I^2 = \neg Y$	$x_{i}^{\star} = F \wedge x_{i'}^{\star} = T$			$\langle \otimes x_{i}^{F} \odot, x_{i'}^{F}, \odot \rangle$			2
$E_i = \langle X_j \rangle$ and $E_i = \langle X_{j'} \rangle$	$x_{i}^{\star} = T \wedge x_{i'}^{\star} = F$				$\langle \otimes x_{i}^{\mathbf{F}}, \odot, x_{i'}^{\mathbf{F}} \odot \rangle$	$\langle \otimes x_{i}^{F}, \odot x_{i'}^{F} \odot \rangle$	2
	$x_j^\star = \mathbf{F} \wedge x_{j'}^\star = \mathbf{F}$			$\langle \otimes x_j^{F} \odot, x_{j'}^{F}, \odot \rangle$	$\langle \otimes \tilde{x}_{j}^{F} \odot, x_{j'}^{\check{F}} \odot \rangle$		2
	$x_i^\star = T \wedge x_{i'}^\star = T$	$\langle \odot, x_j^{T}, \odot, x_{j'}^{T}, \otimes \rangle$			(Og <sup>T</sup> Og <sup>T</sup> O)		2
$L^1 = X$ and $L^2 = X$	$x_{i}^{\star} = \mathbf{F} \wedge x_{i'}^{\star} = \mathbf{T}$		$( \bigcirc r^{\mathrm{T}} \bigcirc r^{\mathrm{T}} \oslash)$	$( \odot, x_j, \odot x_{j'} \otimes )$	$\begin{array}{c} (\odot, x_j, \odot, x_j, \odot) \\ (\odot, x_j, \odot, x_j, \otimes) \end{array}$	$\langle \otimes x_j, \otimes x_{j'} \otimes r \rangle$	2
$L_i = \Lambda_j$ and $L_i = \Lambda_{j'}$	$\dot{x_j^{\star}} = T \wedge \dot{x_{j'}^{\star}} = F$		$\langle \odot, x_j, \odot, x_j \rangle \otimes \rangle$			$\langle \odot x_j^{T} \odot, x_{j'}^{T} \otimes \rangle$	2
	$x_j^\star = \mathbf{F} \wedge x_{j'}^\star = \mathbf{F}$						3

Table 4: Example of applying **m** in  $\mathcal{D}_5$  of the bottom-up tokenisation problem obtained from Reduction 2. The dot symbol  $\cdot$  denotes the string not changing under the given merge.

simulation of why this is the case. We thus have:

$$\left( \begin{array}{c} \left( L_{i}^{2} = \neg X_{j'} \text{ and } \langle \otimes x_{j'}^{\mathbf{r}}, \otimes \rangle, \langle \otimes, x_{j'}^{\mathbf{r}} \otimes \rangle \in \mathbf{m} \end{array} \right) \right)$$
$$= 3I - \sum_{i=1}^{I} \mathbb{1} \{ L_{i}^{1} \lor L_{i}^{2} \}$$
(55b) 115

$$\leq 3I - \psi$$
 (55c) 1151

where, by construction, we have that a merge in our sequence (e.g.,  $\langle \odot, x_j^T \otimes \rangle$  or  $\langle \otimes x_j^F, \odot \rangle$ ), if and only if its value is in a satisfying assignment (e.g.,  $x_j^* = T$  or  $x_j^* = F$  respectively). Summing together the lengths in Eqs. (52) to (55), we get that:

$$\sum_{\mathbf{c}\in\mathcal{D}}|\mathsf{tok}_{\uparrow}[\mathbf{m}](\mathbf{c})| \le \delta = (6f + 6f' + 4f'' + 4f''') J + 3I - \psi$$
(56) 1155

which concludes the proof.

1157

1156

1158

# D Proof of Lemma 7

**Lemma 7.** *If the bottom-up tokenisation instance output by Reduction 2 is satisfiable, the* max-2-SAT 1159 *instance which generated it is as well. Formally:* 1160

$$\operatorname{Tok}_{\uparrow}(\operatorname{R2}(\mathcal{X},\mathcal{L},\psi)) \Longrightarrow \operatorname{M2S}(\mathcal{X},\mathcal{L},\psi)$$
 (27) 1161

*Proof.* First, note that:

$$\sum_{\mathbf{c}\in\mathcal{D}}|\mathbf{c}| = (12f + 12f' + 10f'' + 10f''')J + 5I$$
(57)

1164 Further, let:

 $\mathbf{m}_1$ 

$$\begin{aligned} \operatorname{toklen}(\mathcal{D}, \mathbf{m}) &\stackrel{\text{def}}{=} \sum_{\mathbf{c} \in \mathcal{D}} |\operatorname{tok}_{\uparrow}[\mathbf{m}](\mathbf{c})|, & \mathbf{m}_{0} = \mathbf{m}_{1} \circ \mathbf{m}_{2} \circ \mathbf{m}_{3} \circ \mathbf{m}_{4} \circ \mathbf{m}_{5} \end{aligned} \tag{58} \\ &= \bigcirc_{j=1}^{J} [\langle \otimes, x_{j}^{\mathsf{F}} \rangle, \langle x_{j}^{\mathsf{T}}, \otimes \rangle], & \mathbf{m}_{3} = \bigcirc_{j=1}^{J} [\langle x_{j}^{\mathsf{F}}, \otimes \rangle, \langle \otimes, x_{j}^{\mathsf{T}} \rangle], & \mathbf{m}_{5} = \bigcirc_{j=1}^{J} [\langle \otimes, x_{j}^{\mathsf{F}} \rangle, \langle x_{j}^{\mathsf{T}}, \otimes \rangle] \\ &\mathbf{m}_{2} = \bigcirc_{j=1}^{J} \begin{bmatrix} \langle \otimes, x_{j}^{\mathsf{T}} \otimes \rangle & \operatorname{if} x_{j}^{\mathsf{*}} = \mathsf{T} \\ \langle \otimes x_{j}^{\mathsf{F}}, \otimes \rangle & \operatorname{else} \end{bmatrix}, & \mathbf{m}_{4} = \bigcirc_{j=1}^{J} \begin{bmatrix} \langle \otimes x_{j}^{\mathsf{T}}, \otimes \rangle & \operatorname{if} x_{j}^{\mathsf{*}} = \mathsf{T} \\ \langle \otimes, x_{j}^{\mathsf{F}} \otimes \rangle & \operatorname{else} \end{bmatrix} \\ &\mathbf{m}_{j}^{\otimes} = \begin{cases} \langle \otimes x_{j}^{\mathsf{T}}, \otimes \rangle, \langle \otimes, x_{j}^{\mathsf{F}} \otimes \rangle \\ \langle \otimes, x_{j}^{\mathsf{T}} \otimes \rangle, \langle \otimes, x_{j}^{\mathsf{F}} \otimes \rangle \end{cases} \end{cases}, & \mathbf{m}_{j}^{\otimes} = \begin{cases} \langle \otimes, x_{j}^{\mathsf{T}} \otimes \rangle, \langle \otimes, x_{j}^{\mathsf{F}} \otimes \rangle \\ \langle \otimes x_{j}^{\mathsf{T}}, \otimes \rangle, \langle \otimes, x_{j}^{\mathsf{F}} \otimes \rangle \end{cases} \end{cases} \end{aligned}$$

1185

1191

1192

1193

1168

1162

1163

1165

1166

1167

1169 
$$\mathbf{m}_{j}^{\mathrm{T}} = \left\{ \begin{array}{l} \langle \odot x_{j}^{\mathrm{T}}, \odot \rangle, \langle \odot, x_{j}^{\mathrm{T}} \otimes \rangle, \\ \langle \odot, x_{j}^{\mathrm{T}} \odot \rangle, \langle \odot x_{j}^{\mathrm{T}}, \otimes \rangle, \end{array} \right\}, \qquad \mathbf{m}_{j}^{\mathrm{F}} = \left\{ \begin{array}{l} \langle \odot, x_{j}^{\mathrm{F}} \odot \rangle, \langle \otimes x_{j}^{\mathrm{F}}, \odot \rangle, \\ \langle \odot x_{j}^{\mathrm{F}}, \odot \rangle, \langle \otimes, x_{j}^{\mathrm{F}} \odot \rangle, \end{array} \right\}$$

The maximum number of symbols in this dataset after compression is set to  $\delta = (6f + 6f' + 4f'' + 4f''') J +$ 1170  $3 I - \psi$ . This means that to satisfy this objective, there must exist a vocabulary whose tokeniser compresses 1171 the text by at least  $(6f + 6f' + 6f'' + 6f''') J + 2I + \psi$  symbols. We now prove this lemma in five 1172 steps: (1) we show that any solution which compresses the text by at least 6fJ symbols must include 1173 all merges in  $m_1$ ,  $m_3$ , and  $m_5$ ; 2) we show that any solution which compresses the text by at least 1174 (6f + 6f')J symbols must only include either merges in  $\mathbf{m}_1, \mathbf{m}_3, \mathbf{m}_5$ , or in either  $\mathbf{m}_i^{\odot}$  or  $\mathbf{m}_i^{\otimes}$ ; ③ we 1175 show that any solution which compresses the text by at least (6f + 6f' + 6f'')J symbols must include, 1176 for each  $j \in \{1, J\}$ , exactly one merge in set  $\mathbf{m}_{i}^{\otimes}$  and one in set  $\mathbf{m}_{i}^{\otimes}$ ; ④ we show that any solution which 1177 compresses the text by at least (6f + 6f' + 6f'' + 6f''')J symbols must include, for each  $j \in \{1, J\}$ , 1178 exactly two merge in either set  $\mathbf{m}_j^{\mathsf{T}}$  or in set  $\mathbf{m}_j^{\mathsf{F}}$ ; (5) we show that any solution which compresses the text by at least  $(6f + 6f' + 4f'' + 4f''')J + 6f'' + 6f'')J + 2I + \psi$  symbols must be produced by a 1179 1180 max-2-SAT problem with at least  $\psi$  satisfied clauses. 1181

**LemmaProofStep 1.** (Step (1)). Any solution which compresses the text by at least 6fJ symbols must 1182 include all merges in  $m_1$ ,  $m_3$ , and  $m_5$ , *i.e.*,: 1183

1184 
$$\left(\operatorname{toklen}(\mathcal{D},\mathbf{m}) \leq \underbrace{6fJ + (12f' + 10f'' + 10f''')J + 5I}_{\sum -1}\right)$$
(59)

$$\implies \underbrace{\bigcirc_{j=1}^{J}[\langle \otimes, x_{j}^{\mathsf{F}} \rangle, \langle x_{j}^{\mathsf{T}}, \otimes \rangle]}_{\mathbf{m}_{1}} \subset \mathbf{m}, \quad \underbrace{\bigcirc_{j=1}^{J}[\langle x_{j}^{\mathsf{F}}, \otimes \rangle, \langle \otimes, x_{j}^{\mathsf{T}} \rangle]}_{\mathbf{m}_{3}} \subset \mathbf{m}, \quad \underbrace{\bigcirc_{j=1}^{J}[\langle \otimes, x_{j}^{\mathsf{F}} \rangle, \langle x_{j}^{\mathsf{T}}, \otimes \rangle]}_{\mathbf{m}_{5}} \subset \mathbf{m}$$

*Proof.* We prove this statement by contradiction. Assume that one of the subwords above is not present 1186 in the tokenisers' merge sequence **m**. In that case, the strings in  $\mathcal{D}_1$  which contain this character string 1187 will not be compressed, and will thus still be represented with 2 symbols. There will thus be at most 1188 6J-1 strings in  $\mathcal{D}_1$  represented with a single symbol, and at least one represented with two symbols. 1189 The minimum length achievable would thus be: 1190

$$\operatorname{toklen}(\mathcal{D}, \mathbf{m}) = \sum_{\substack{\mathbf{c} \in \bigcup_{j=1}^{f} \mathcal{D}_{0} \\ \geq (6J-1)f+2f \\} > (6J+1)f}} |\operatorname{tok}_{\uparrow}[\mathbf{m}](\mathbf{c})| + \sum_{\substack{\mathbf{c} \in \mathcal{D} \setminus (\bigcup_{j=1}^{f} \mathcal{D}_{0}) \\} > 0}} |\operatorname{tok}_{\uparrow}[\mathbf{m}](\mathbf{c})|$$
(60a)  
$$> (6J+1)f \qquad \text{By construction } f = (12f' + 10f'' + 10f''') J + 5I$$
(60b)  
(60b)

$$= (6f + 12f' + 10f'' + 10f''') J + 5I$$
(60c)

which contradicts the proofs statement. 1194

**LemmaProofStep 2.** (Step (2)). Any solution which compresses the text by at least (6f + 6f')J symbols must only include either merges in  $\mathbf{m}_1$ ,  $\mathbf{m}_3$ ,  $\mathbf{m}_5$ , or in either  $\mathbf{m}_j^{\odot}$  or  $\mathbf{m}_j^{\otimes}$ , i.e.,: 1196

$$\left(\operatorname{toklen}(\mathcal{D}, \mathcal{S}) \leq \underbrace{(6f + 6f')J + (10f'' + 10f''')J + 5I}_{\sum_{c \in \mathcal{D}} |c| - (6f + 6f')J}\right)$$
(61) 119

$$\implies \mathbf{m} \setminus (\mathbf{m}_{1} \circ \mathbf{m}_{3} \circ \mathbf{m}_{5}) \subseteq \underbrace{\left\{ \begin{array}{l} \langle \odot, x_{j}^{\mathsf{T}} \otimes \rangle, \langle \odot x_{j}^{\mathsf{F}}, \odot \rangle, \langle \odot x_{j}^{\mathsf{T}}, \odot \rangle, \langle \odot, x_{j}^{\mathsf{F}} \otimes \rangle, \langle \odot, x_{j}^{$$

*Proof.* We again prove this statement by contradiction. Assume that **m** has all merges  $\mathbf{m}_1, \mathbf{m}_3, \mathbf{m}_5$ , but one of its other merges is in neither of the sets  $\mathbf{m}_j^{\otimes}$  and  $\mathbf{m}_j^{\otimes}$ . This means that at least one of the sets  $\mathbf{m}_j^{\otimes}$  and  $\mathbf{m}_j^{\otimes}$  will have no merge in the solution; this is because there are 2*J* such sets, which—coupled together with the 6*J* already selected merges in  $\mathbf{m}_1, \mathbf{m}_3, \mathbf{m}_5$ —would amount to the maximum of 8*J* merges. In that case, the strings (e.g.,  $\otimes x_j^T \otimes, \otimes x_j^F \otimes, \otimes x_j^T \otimes$  and  $\otimes x_j^F \otimes$ ) in  $\mathcal{D}_2$  containing the characters this absent merge represents will not be fully compressed to a single symbol, being represented with 2 symbols instead. There will thus be at most 6J - 1 strings in  $\mathcal{D}_2$  represented with a single symbol, and at least one represented with two symbols. The minimum length achievable would thus be: 1200

$$\operatorname{toklen}(\mathcal{D}, \mathbf{m}) = \sum_{\substack{\mathbf{c} \in \bigcup_{i=1}^{f} \mathcal{D}_{1} \\ = 6fJ}} |\operatorname{tok}_{\uparrow}[\mathbf{m}](\mathbf{c})| + \sum_{\substack{\mathbf{c} \in \bigcup_{i=1}^{f'} \mathcal{D}_{2} \\ \geq (6J-1)f'+2f'}} |\operatorname{tok}_{\uparrow}[\mathbf{m}](\mathbf{c})| + \sum_{\substack{\mathbf{c} \in \mathcal{D} \setminus (\mathcal{D}_{1} \cup \mathcal{D}_{2}) \\ > 0}} |\operatorname{tok}_{\uparrow}[\mathbf{m}](\mathbf{c})| + \sum_{\substack{\mathbf{c} \in \mathcal{D} \setminus (\mathcal{D}_{1} \cup \mathcal{D}_{2}) \\ > 0}} |\operatorname{tok}_{\uparrow}[\mathbf{m}](\mathbf{c})| + \sum_{\substack{\mathbf{c} \in \mathcal{D} \setminus (\mathcal{D}_{1} \cup \mathcal{D}_{2}) \\ > 0}} |\operatorname{tok}_{\uparrow}[\mathbf{m}](\mathbf{c})| + \sum_{\substack{\mathbf{c} \in \mathcal{D} \setminus (\mathcal{D}_{1} \cup \mathcal{D}_{2}) \\ > 0}} |\operatorname{tok}_{\uparrow}[\mathbf{m}](\mathbf{c})| + \sum_{\substack{\mathbf{c} \in \mathcal{D} \setminus (\mathcal{D}_{1} \cup \mathcal{D}_{2}) \\ > 0}} |\operatorname{tok}_{\uparrow}[\mathbf{m}](\mathbf{c})| + \sum_{\substack{\mathbf{c} \in \mathcal{D} \setminus (\mathcal{D}_{1} \cup \mathcal{D}_{2}) \\ > 0}} |\operatorname{tok}_{\uparrow}[\mathbf{m}](\mathbf{c})| + \sum_{\substack{\mathbf{c} \in \mathcal{D} \setminus (\mathcal{D}_{1} \cup \mathcal{D}_{2}) \\ > 0}} |\operatorname{tok}_{\uparrow}[\mathbf{m}](\mathbf{c})| + \sum_{\substack{\mathbf{c} \in \mathcal{D} \setminus (\mathcal{D}_{1} \cup \mathcal{D}_{2}) \\ > 0}} |\operatorname{tok}_{\uparrow}[\mathbf{m}](\mathbf{c})| + \sum_{\substack{\mathbf{c} \in \mathcal{D} \setminus (\mathcal{D}_{1} \cup \mathcal{D}_{2}) \\ > 0}} |\operatorname{tok}_{\uparrow}[\mathbf{m}](\mathbf{c})| + \sum_{\substack{\mathbf{c} \in \mathcal{D} \setminus (\mathcal{D}_{1} \cup \mathcal{D}_{2}) \\ > 0}} |\operatorname{tok}_{\uparrow}[\mathbf{m}](\mathbf{c})| + \sum_{\substack{\mathbf{c} \in \mathcal{D} \setminus (\mathcal{D}_{1} \cup \mathcal{D}_{2}) \\ > 0}} |\operatorname{tok}_{\uparrow}[\mathbf{m}](\mathbf{c})| + \sum_{\substack{\mathbf{c} \in \mathcal{D} \setminus (\mathcal{D}_{1} \cup \mathcal{D}_{2}) \\ > 0}} |\operatorname{tok}_{\uparrow}[\mathbf{m}](\mathbf{c})| + \sum_{\substack{\mathbf{c} \in \mathcal{D} \setminus (\mathcal{D}_{1} \cup \mathcal{D}_{2}) \\ > 0}} |\operatorname{tok}_{\uparrow}[\mathbf{m}](\mathbf{c})| + \sum_{\substack{\mathbf{c} \in \mathcal{D} \setminus (\mathcal{D}_{1} \cup \mathcal{D}_{2}) \\ > 0}} |\operatorname{tok}_{\uparrow}[\mathbf{m}](\mathbf{c})| + \sum_{\substack{\mathbf{c} \in \mathcal{D} \setminus (\mathcal{D}_{1} \cup \mathcal{D}_{2}) \\ > 0}} |\operatorname{tok}_{\uparrow}[\mathbf{m}](\mathbf{c})| + \sum_{\substack{\mathbf{c} \in \mathcal{D} \setminus (\mathcal{D}_{1} \cup \mathcal{D}_{2}) \\ > 0} |\operatorname{tok}_{\uparrow}[\mathbf{m}](\mathbf{c})| + \sum_{\substack{\mathbf{c} \in \mathcal{D} \setminus (\mathcal{D}_{1} \cup \mathcal{D}_{2}) \\ > 0} |\operatorname{tok}_{\uparrow}[\mathbf{m}](\mathbf{c})| + \sum_{\substack{\mathbf{c} \in \mathcal{D} \setminus (\mathcal{D}_{1} \cup \mathcal{D}_{2}) \\ > 0} |\operatorname{tok}_{\downarrow}[\mathbf{m}](\mathbf{c})| + \sum_{\substack{\mathbf{c} \in \mathcal{D} \setminus (\mathcal{D}_{1} \cup \mathcal{D}_{2}) \\ > 0} |\operatorname{tok}_{\downarrow}[\mathbf{m}](\mathbf{c})| + \sum_{\substack{\mathbf{c} \in \mathcal{D} \setminus (\mathcal{D}_{1} \cup \mathcal{D}_{2}) \\ > 0} |\operatorname{tok}_{\downarrow}[\mathbf{m}](\mathbf{c})| + \sum_{\substack{\mathbf{c} \in \mathcal{D} \setminus (\mathcal{D}_{1} \cup \mathcal{D}_{2}) |\operatorname{tok}_{\downarrow}[\mathbf{m}](\mathbf{c})| + \sum_{\substack{\mathbf{c} \in \mathcal{D} \setminus (\mathcal{D}_{1} \cup \mathcal{D}_{2}) |\operatorname{tok}_{\downarrow}[\mathbf{m}](\mathbf{c})| + \sum_{\substack{\mathbf{c} \in \mathcal{D} \setminus (\mathcal{D}_{1} \cup \mathcal{D}_{2}) |\operatorname{tok}_{\downarrow}[\mathbf{m}](\mathbf{c})| + \sum_{\substack{\mathbf{c} \in \mathcal{D} \setminus (\mathcal{D}_{1} \cup \mathcal{D}_{2}) |\operatorname{tok}_{\downarrow}[\mathbf{m}](\mathbf{c})| + \sum_{\substack{\mathbf{c} \in \mathcal{D} \setminus (\mathcal{D}_{1} \cup \mathcal{D}_{2} \cup \mathcal{D}_{2} \cup \mathcal{D}_{2} \cup \mathcal{D}_{2} \mid + \sum_{\substack{\mathbf$$

$$= (6f + 6f' + 10f'' + 10f''') J + 5I$$
(62c) 12

1210

which contradicts the proofs statement.

**LemmaProofStep 3.** (Step ③). Any solution which compresses the text by at least (6f + 6f' + 6f'')J 1211 symbols must include all merges in  $\mathbf{m}_1$ ,  $\mathbf{m}_3$ ,  $\mathbf{m}_5$ , and, for each  $j \in \{1, J\}$ , exactly one merge in set  $\mathbf{m}_j^{\odot}$  1212 and one in set  $\mathbf{m}_j^{\odot}$ , i.e.,: 1213

$$\left(\operatorname{toklen}(\mathcal{D}, \mathbf{m}) \leq \underbrace{(6f + 6f' + 4f'')J + 10f'''J + 5I}_{\sum_{c \in \mathcal{D}} |c| - (6f + 6f' + 6f'')J}\right)$$
(63)

$$\implies \forall_{j \in \{1, \dots, J\}} \left| \mathbf{m} \cap \underbrace{\left\{ \begin{array}{c} \langle \odot x_j^{\mathsf{T}}, \odot \rangle, \langle \odot, x_j^{\mathsf{F}} \odot \rangle \\ \langle \odot, x_j^{\mathsf{T}} \odot \rangle, \langle \odot x_j^{\mathsf{F}}, \odot \rangle \end{array} \right\}}_{\mathbf{m}_j^{\odot}} \right| = 1 \text{ and } \left| \mathbf{m} \cap \underbrace{\left\{ \begin{array}{c} \langle \odot, x_j^{\mathsf{T}} \otimes \rangle, \langle \odot x_j^{\mathsf{F}}, \odot \rangle \\ \langle \odot x_j^{\mathsf{T}}, \otimes \rangle, \langle \otimes, x_j^{\mathsf{F}} \odot \rangle \end{array} \right\}}_{\mathbf{m}_j^{\odot}} \right| = 1$$

Proof. We again prove this statement by contradiction. First, assume that m contains all the merges in 1216  $\mathbf{m}_1, \mathbf{m}_3, \mathbf{m}_5$ ; further, assume all its other merges are contained in sets  $\mathbf{m}_j^{\odot}$  and  $\mathbf{m}_j^{\otimes}$ . Note now that, if any 1217 merge in  $\mathbf{m}_{j}^{\otimes}$  is in the selected merges  $\mathbf{m}$ , the string  $\otimes x_{j}^{\mathsf{F}} \odot x_{j}^{\mathsf{T}} \otimes$  in  $\mathcal{D}_{3}$  will be compressed to 2 symbols 1218 (e.g.,  $\langle \otimes x_i^{\rm F}, \odot x_i^{\rm T} \otimes \rangle$ ); if none of these merges is present, however, this string will only be compressed 1219 to 3 symbols (e.g.,  $\langle \otimes x_j^{\mathsf{F}}, \odot, x_j^{\mathsf{T}} \otimes \rangle$ ). The same is true for strings  $\odot x_j^{\mathsf{T}} \odot x_j^{\mathsf{F}} \odot$  and merges in  $\mathbf{m}_j^{\odot}$ . Now, assume the contradictory case: for a value of  $j \in \{1, J\}$ , **m** does not satisfy the condition above. As, by 1220 1221 construction, our solution has K = 8J merges, and because  $|\mathbf{m}_1 \circ \mathbf{m}_3 \circ \mathbf{m}_5| = 6J$ , we know that we have 1222 2J merges in sets  $\mathbf{m}_{i}^{\otimes}$  and  $\mathbf{m}_{i}^{\otimes}$ . As there are exactly 2J such sets, if the condition above does not hold, at 1223 least one of these sets must have no merge present in m. In that case, the strings in  $\mathcal{D}_3$  which contain the 1224 character string represented by these absent merges will be compressed to three symbols, while others will be compressed to two symbols. There will thus be at most 2J - 1 strings in  $\mathcal{D}_3$  represented with two 1226 symbols, and at least one represented with three symbols. The minimum length achievable would thus be:

1229 
$$\operatorname{toklen}(\mathcal{D}, \mathbf{m}) = \sum_{\substack{\mathbf{c} \in \bigcup_{j=1}^{f} \mathcal{D}_{1} \cup \bigcup_{j=1}^{f'} \mathcal{D}_{2} \\ =(6f+6f')J}} |\operatorname{tok}_{\uparrow}[\mathbf{m}](\mathbf{c})| + \sum_{\substack{\mathbf{c} \in \bigcup_{j=1}^{f''} \mathcal{D}_{3} \\ \ge (2J-1)2f''+3f''}} |\operatorname{tok}_{\uparrow}[\mathbf{m}](\mathbf{c})| + \sum_{\substack{\mathbf{c} \in \bigcup_{j=1}^{f'''} \mathcal{D}_{4} \cup \mathcal{D}_{5} \\ \ge (2J-1)2f''+3f''}} |\operatorname{tok}_{\uparrow}[\mathbf{m}](\mathbf{c})|$$
(64a)  
1230 
$$> (6f+6f')J + (4J+1)f'' \qquad \text{By construction } f'' = 10f''' J + 5I$$
(64b)

1231

1227 1228

$$= (6f + 6f' + 4f'' + 10f''') J + 5I$$
(64c)

(66)

 $\mathbf{m}_{i}^{\mathbf{F}}$ 

which contradicts the proof's statement.

**LemmaProofStep 4.** (Step (4)). Any solution which compresses the text by at least (6f + 6f' + 6f'' + 6f'')6f''') J symbols must include all merges in  $\mathbf{m}_1$ ,  $\mathbf{m}_3$ ,  $\mathbf{m}_5$ , and, for each  $j \in \{1, J\}$ , exactly one merge in set  $\mathbf{m}_{i}^{\odot}$  and one in set  $\mathbf{m}_{i}^{\otimes}$ , such that either both these merges are in  $\mathbf{m}_{i}^{\mathrm{T}}$  or both are in  $\mathbf{m}_{i}^{\mathrm{F}}$ , i.e.,: 1235

$$1236 \qquad \left( \operatorname{toklen}(\mathcal{D}, \mathbf{m}) \leq \underbrace{(6f + 6f' + 4f'' + 4f''')J + 5I}_{\sum_{\mathbf{c} \in \mathcal{D}} |\mathbf{c}| - (6f + 6f' + 6f'' + 6f''')J} \right)$$
(65)  
$$1237 \qquad \Longrightarrow \forall_{j \in \{1, \dots, J\}} |\mathbf{m} \cap \left\{ \begin{array}{l} \langle @x_j^{\mathsf{T}}, @ \rangle, \langle @, x_j^{\mathsf{T}} \otimes \rangle, \\ \langle @, x_j^{\mathsf{T}} @ \rangle, \langle @, x_j^{\mathsf{T}} \otimes \rangle, \\ \langle @, x_j^{\mathsf{T}} @ \rangle, \langle @, x_j^{\mathsf{T}} \otimes \rangle, \\ \langle @, x_j^{\mathsf{T}} @ \rangle, \langle @, x_j^{\mathsf{T}} \otimes \rangle, \\ \langle @, x_j^{\mathsf{T}} @ \rangle, \langle @, x_j^{\mathsf{T}} \otimes \rangle, \\ \langle @, x_j^{\mathsf{T}} @ \rangle, \langle @, x_j^{\mathsf{T}} \otimes \rangle, \\ \langle @, x_j^{\mathsf{T}} \otimes \rangle, \langle @, x_j^{\mathsf{T}} \otimes \rangle, \\ \langle @, x_j^{\mathsf{T}} \otimes \rangle, \langle @, x_j^{\mathsf{T}} \otimes \rangle, \\ \end{array} \right\} | = 2 \text{ or } |\mathbf{m} \cap \left\{ \begin{array}{l} \langle @, x_j^{\mathsf{T}} \otimes \rangle, \langle @, x_j^{\mathsf{T}} \otimes \rangle, \\ \langle @, x_j^{\mathsf{T}} \otimes \rangle, \langle @, x_j^{\mathsf{T}} \otimes \rangle, \\ \langle @, x_j^{\mathsf{T}} \otimes \rangle, \langle @, x_j^{\mathsf{T}} \otimes \rangle, \\ \end{array} \right\} | = 2 \text{ or } |\mathbf{m} \cap \left\{ \begin{array}{l} \langle @, x_j^{\mathsf{T}} \otimes \rangle, \langle @, x_j^{\mathsf{T}} \otimes \rangle, \\ \langle @, x_j^{\mathsf{T}} \otimes \rangle, \langle @, x_j^{\mathsf{T}} \otimes \rangle, \\ \end{array} \right\} | = 2 \text{ or } |\mathbf{m} \cap \left\{ \begin{array}{l} \langle @, x_j^{\mathsf{T}} \otimes \rangle, \\ \end{array} \right\} | = 2 \text{ or } |\mathbf{m} \cap \left\{ \begin{array}{l} \langle @, x_j^{\mathsf{T}} \otimes \rangle, \\ \end{array} \right\} | = 2 \text{ or } |\mathbf{m} \cap \left\{ \begin{array}{l} \langle @, x_j^{\mathsf{T}} \otimes \rangle, \\ \end{array} \right\} | = 2 \text{ or } |\mathbf{m} \cap \left\{ \begin{array}{l} \langle @, x_j^{\mathsf{T}} \otimes \rangle, \\ \end{array} \right\} | = 2 \text{ or } |\mathbf{m} \cap \left\{ \begin{array}{l} \langle @, x_j^{\mathsf{T}} \otimes \rangle, \\ \langle @, x_j^{\mathsf{T}} \otimes \rangle, \\ \langle @, x_j^{\mathsf{T}} \otimes \rangle, \\ \rangle | = 2 \text{ or } |\mathbf{m} \cap \left\{ \begin{array}{l} \langle @, x_j^{\mathsf{T}} \otimes \rangle, \\ \langle @, x_j^{\mathsf{T}} \otimes \rangle, \\ \langle @, x_j^{\mathsf{T}} \otimes \rangle, \\ \rangle | = 2 \text{ or } |\mathbf{m} \cap \left\{ \begin{array}{l} \langle @, x_j^{\mathsf{T}} \otimes \rangle, \\ \langle @, x_j^{\mathsf{T}} \otimes \rangle, \\ \rangle | = 2 \text{ or } |\mathbf{m} \cap \left\{ \begin{array}{l} \langle @, x_j^{\mathsf{T}} \otimes \rangle, \\ \langle @, x_j^{\mathsf{T}} \otimes \rangle, \\ \langle @, x_j^{\mathsf{T}} \otimes \rangle, \\ \rangle | = 2 \text{ or } |\mathbf{m} \cap \left\{ \begin{array}{l} \langle @, x_j^{\mathsf{T}} \otimes \rangle, \\ \langle @, x_j^{\mathsf{T}} \otimes \rangle, \\ \rangle | = 2 \text{ or } |\mathbf{m} \cap \left\{ \begin{array}{l} \langle @, x_j^{\mathsf{T}} \otimes \rangle, \\ \langle @, x_j^{\mathsf{T}} \otimes \rangle, \\ \rangle | = 2 \text{ or } |\mathbf{m} \cap \left\{ \begin{array}{l} \langle @, x_j^{\mathsf{T}} \otimes \rangle, \\ \langle @, x_j^{\mathsf{T}} \otimes \rangle, \\ \rangle | = 2 \text$$

 $\mathbf{m}_{i}^{\mathrm{T}}$ 

1244

1253

1254 1255

*Proof.* First, note that the conditions of the step of our proof are stricter than previous ones, so we assume 1239 the conditions of steps (1) to (3) hold—i.e., m contains all merges in  $m_1, m_3, m_5$ ; further, it has one and 1240 only one merge from each set  $\mathbf{m}_{j}^{\otimes}$  and  $\mathbf{m}_{j}^{\otimes}$ . (Note that  $\mathbf{m}_{j}^{\otimes} \cup \mathbf{m}_{j}^{\otimes} = \mathbf{m}_{j}^{\mathsf{T}} \cup \mathbf{m}_{j}^{\mathsf{F}}$ , and that the just-mentioned 1241 condition implies  $|\mathbf{m} \cap (\mathbf{m}_i^{\mathsf{T}} \cup \mathbf{m}_i^{\mathsf{F}})| = 2$ .) We now again prove this statement by contradiction. Consider now the case: 1243

$$\left| \mathbf{m} \cap \underbrace{\left\{ \begin{array}{c} \langle \odot x_{j}^{\mathsf{T}}, \odot \rangle, \langle \odot, x_{j}^{\mathsf{T}} \otimes \rangle, \\ \langle \odot, x_{j}^{\mathsf{T}} \odot \rangle, \langle \odot x_{j}^{\mathsf{T}}, \otimes \rangle, \end{array} \right\}}_{\mathbf{m}_{j}^{\mathsf{T}}} \right| = 2 \text{ or } \left| \mathbf{m} \cap \underbrace{\left\{ \begin{array}{c} \langle \odot, x_{j}^{\mathsf{F}} \odot \rangle, \langle \otimes x_{j}^{\mathsf{F}}, \odot \rangle, \\ \langle \odot x_{j}^{\mathsf{F}}, \odot \rangle, \langle \otimes, x_{j}^{\mathsf{F}} \odot \rangle, \end{array} \right\}}_{\mathbf{m}_{j}^{\mathsf{F}}} \right| = 2 \tag{67}$$

If this is true, then strings  $\otimes x_i^{\mathsf{F}} \otimes x_i^{\mathsf{T}} \otimes$  and  $\otimes x_i^{\mathsf{F}} \otimes x_i^{\mathsf{T}} \otimes$  in  $\mathcal{D}_4$  will be compressed to 2 symbols each (e.g., to 1245  $\langle \odot x_j^{\mathsf{F}}, \odot x_j^{\mathsf{T}} \otimes \rangle$  and  $\langle \otimes x_j^{\mathsf{F}}, \odot x_j^{\mathsf{T}} \odot \rangle$  or  $\langle \odot x_j^{\mathsf{F}} \odot, x_j^{\mathsf{T}} \otimes \rangle$  and  $\langle \otimes x_j^{\mathsf{F}} \odot, x_j^{\mathsf{T}} \odot \rangle$  ); if this condition is false, however, one of these strings will only be compressed to 3 symbols (e.g., to  $\langle \odot x_j^{\mathsf{F}}, \odot x_j^{\mathsf{T}} \otimes \rangle$  and  $\langle \otimes x_j^{\mathsf{F}}, \odot, x_j^{\mathsf{T}} \odot \rangle$ ). 1246 1247 Now, assume the contradictory case: for a value of  $j \in \{1, J\}$ , **m** does not satisfy the condition above. In 1248 that case, the strings in  $\mathcal{D}_4$  for which the condition does not hold will be compressed to 3 + 2 symbols, 1249 while others will be compressed to 2+2 symbols. There will thus be at most 2J-1 strings in  $\mathcal{D}_4$ 1250 represented with two symbols, and at least one represented with three symbols. The minimum length 1251 achievable would thus be: 1252

$$toklen(\mathcal{D}, \mathbf{m}) = \sum_{\substack{c \in \bigcup_{i=1}^{f} \mathcal{D}_{1} \cup \bigcup_{i=1}^{f'} \mathcal{D}_{2} \cup \bigcup_{i=1}^{f''} \mathcal{D}_{3} \\ = (6f + 6f' + 4f'')J}} |tok_{\uparrow}[\mathbf{m}](\mathbf{c})| + \sum_{\substack{c \in \mathcal{D}_{5} \\ = (2J-1)2f''' + 3f'''}} |tok_{\uparrow}[\mathbf{m}](\mathbf{c})| + \sum_{\substack{c \in \mathcal{D}_{5} \\ >0}} |tok_{\uparrow}[\mathbf{m}](\mathbf{c})| |tok_{\uparrow}[\mathbf{m}](\mathbf{c})| + \sum_{\substack{c \in \mathcal{D}_{5} \\ >0}} |tok_{\uparrow}[\mathbf{m}](\mathbf{c})| |tok_{\uparrow}[\mathbf{m}](\mathbf{c})| + \sum_{\substack{c \in \mathcal{D}_{5} \\ \geq (2J-1)2f''' + 3f'''}} |tok_{\uparrow}[\mathbf{m}](\mathbf{c})| + \sum_{\substack{c \in \mathcal{D}_{5} \\ >0}} |tok_{\uparrow}[\mathbf{m}](\mathbf{c})| + \sum_{\substack{c \in \mathcal{D}_{5} \\ \geq (2J-1)2f''' + 3f'''}} |tok_{\uparrow}[\mathbf{m}](\mathbf{c})| + \sum_{\substack{c \in \mathcal{D}_{5} \\ \geq 0}} |tok_{\uparrow}[\mathbf{m}](\mathbf{c})| + \sum_{\substack{c \in \mathcal{D}_{5} \\ \geq 0}} |tok_{\uparrow}[\mathbf{m}](\mathbf{c})| + \sum_{\substack{c \in \mathcal{D}_{5} \\ \geq 0}} |tok_{\uparrow}[\mathbf{m}](\mathbf{c})| + \sum_{\substack{c \in \mathcal{D}_{5} \\ \geq (2J-1)2f''' + 3f'''}} |tok_{\uparrow}[\mathbf{m}](\mathbf{c})| + \sum_{\substack{c \in \mathcal{D}_{5} \\ \geq 0}} |tok_{\downarrow}[\mathbf{m}](\mathbf{c})| + \sum_{\substack{c \in \mathcal{D}_{5} \\ \geq 0}} |tok_{$$

which contradicts the proof's statement. 1256

**LemmaProofStep 5.** (Step (5)). Any instance of the tokenisation problem with a solution which com-1257 presses the text by at least  $(6f + 6f' + 6f'' + 6f'')J + 2I + \psi$  symbols must be produced by a max-2-SAT 1258 problem with at least  $\psi$  satisfied clauses, i.e.,: 1259

$$\left(\operatorname{toklen}(\mathcal{D}, \mathcal{S}) \leq \underbrace{(6f + 6f' + 4f'' + 4f''')J + 3I - \psi}_{\sum_{c \in \mathcal{D}} |c| - (6f + 6f' + 6f''')J - 2I - \psi}\right) \implies \operatorname{M2S}(\mathcal{X}, \mathcal{L}, \psi)$$
1260

*Proof.* Finally, we now know any solution with this compression must have—for any variable  $X_j$ —either 1261 two merges in  $\mathbf{m}_{i}^{\mathrm{T}}$  or in  $\mathbf{m}_{i}^{\mathrm{F}}$  (and never both). We can thus create a bijection  $\mathrm{Conv}_{\mathbf{m}\to\chi}$  between the set 1262 of possible merge sequences respecting this condition, and the set of T/F assignments to SAT variables  $\chi$ : 1263

$$\operatorname{Conv}_{\mathbf{m}\to\chi}(\mathbf{m}) = \left\{ \begin{array}{ll} \mathsf{T} & \operatorname{if} |\mathbf{m}\cap\mathbf{m}_{j}^{\mathsf{T}}| = 2\\ \mathsf{F} & \operatorname{if} |\mathbf{m}\cap\mathbf{m}_{j}^{\mathsf{F}}| = 2 \end{array} \right\}_{j=1}^{J}$$
(69) 1264

Further, note that merge sequences of this form (as shown in Eq. (42)) lead to exactly (6f + 6f' + 6f'' + 6f'')1265 6f''')J symbols being compressed in datasets  $\mathcal{D}_1$  to  $\mathcal{D}_4$ . To achieve the target compression, a solution 1266 must thus compress  $\mathcal{D}_5$  by at least  $2I + \psi$  symbols. Now note that for any string, e.g.,  $\odot x_j^{\mathsf{T}} \odot x_{j'}^{\mathsf{F}} \odot$ , in  $\mathcal{D}_5$  we have three compression options:  $\odot x_i^{\mathsf{T}} \odot$  and  $x_{i'}^{\mathsf{F}} \odot$  will be compressed, saving 3 symbols;  $\odot x_i^{\mathsf{T}}$  and 1268  $\odot x_{j'}^{\mathsf{F}} \odot$  will be compressed, also saving 3 symbols; or only  $\odot x_{j}^{\mathsf{T}}$  and  $x_{j'}^{\mathsf{F}} \odot$  will be compressed saving only 1269 2 symbols. More specifically,  $\odot x_j^{\mathrm{T}} \odot$  will be compressed to a single symbol if merge  $\langle \odot, x_j^{\mathrm{T}} \odot \rangle$  exists; similarly,  $\odot x_{i'}^{\mathsf{F}} \odot$  will be compressed to a single symbol if merge  $\langle \odot x_{i'}^{\mathsf{F}}, \odot \rangle$  exists. They cannot both 1271 be compressed, however, as there is only one symbol o between the literals. We thus get a reduction 1272 of 3 symbols for each of these strings if at least one of its literals has an associated merge in m. Note 1273 thus that whenever a string  $\odot x_i^{\rm T} \odot x_{i'}^{\rm F} \odot$  is compressed by 3 symbols using merges **m**, the max-2-SAT 1274 disjunction  $X_j \vee \neg X_{j'}$  will also be satisfied by assignment  $\chi = \text{Conv}_{\mathbf{m} \to \chi}(\mathbf{m})$ ; similarly, whenever this 1275 string is only compressed by two symbols, the max-2-SAT disjunction will not be satisfied. A similar logic applies to all potential strings in  $\mathcal{D}_5$ :  $\odot x_j^{\mathsf{T}} \odot x_{j'}^{\mathsf{F}} \odot$ ,  $\odot x_j^{\mathsf{T}} \odot x_j^{\mathsf{F}} \odot$ ,  $\otimes x_j^{\mathsf{F}} \odot x_{j'}^{\mathsf{F}} \odot$ , and  $\odot x_j^{\mathsf{T}} \odot x_{j'}^{\mathsf{T}} \otimes$ . As our condition assumes a compression of at least  $2I + \psi$  symbols, we know that we have at least  $\psi$  strings 1278 for which a literal has an associated merge. We can thus write: 1279

$$2I + \psi \le \max_{\mathbf{m} \in \mathcal{M}^*} \sum_{\mathbf{c} \in \mathcal{D}_5} |\mathbf{c}| - |\mathsf{tok}_{\uparrow}[\mathbf{m}](\mathbf{c})|$$
(70a) 1280

$$=2I + \max_{\mathbf{m} \in \mathcal{M}^{*}} \sum_{\mathbf{c} \in \mathcal{D}_{5}} \mathbb{1} \left\{ \begin{array}{c} \left(x_{j}^{\mathsf{T}} \in \mathbf{c}\right) \text{ and } \left(|\mathbf{m}_{j}^{\mathsf{T}} \cap \mathbf{m}| = 2\right) \\ \text{or} \\ \left(x_{j}^{\mathsf{F}} \in \mathbf{c}\right) \text{ and } \left(|\mathbf{m}_{j}^{\mathsf{F}} \cap \mathbf{m}| = 2\right) \\ \text{or} \\ \left(x_{j'}^{\mathsf{T}} \in \mathbf{c}\right) \text{ and } \left(|\mathbf{m}_{j'}^{\mathsf{T}} \cap \mathbf{m}| = 2\right) \\ \text{or} \\ \left(x_{j'}^{\mathsf{F}} \in \mathbf{c}\right) \text{ and } \left(|\mathbf{m}_{j'}^{\mathsf{F}} \cap \mathbf{m}| = 2\right) \end{array} \right\}$$
(70b)

1281

 $=2I + \max_{x \in \{0,1\}^J} \sum_{i=1}^{I} \mathbb{1}\{L_i^1 \lor L_i^2\}$ (70c)1282

$$\implies$$
 M2S $(\mathcal{X}, \mathcal{L}, \psi)$  (70d) 1283

We thus know that, if a satisfying tokenisation solution exists, then the associated max-2-SAT problem 1284 will also be satisfiable. This concludes the proof. 1285