

# Appendix

## Contents

<b>A</b>	<b>Analysis of LZW: Proofs of Theorems 3.4 and 3.6</b>	<b>14</b>
A.1	Notation and definitions . . . . .	14
A.2	A basic result about the optimal achievable cross-entropy loss . . . . .	14
A.3	Proof of Theorem 2.1 . . . . .	15
A.4	Maximum likelihood unigram model . . . . .	15
A.5	Proof of Theorem 3.4 . . . . .	16
A.6	Heavy-hitter dictionaries and a proof of Theorem 3.6 . . . . .	17
<b>B</b>	<b>Additional Theoretical Results I: A sequential variant of BPE</b>	<b>22</b>
B.1	Analysis of Algorithm 1 . . . . .	25
B.2	Analysis for the large dictionary case: $ \text{Dict}  > d_0$ . . . . .	25
B.3	Analysis in the small dictionary case . . . . .	32
<b>C</b>	<b>Additional Theoretical Results II: Learning the likelihood model</b>	<b>35</b>
C.1	Proof of Theorem C.1 . . . . .	36
<b>D</b>	<b>Additional Theoretical Results III: The generalization ability of tokenizers</b>	<b>40</b>
<b>E</b>	<b>Additional Theoretical Results IV: Interaction between the dictionary and encoding algorithm</b>	<b>41</b>
E.1	Stochastic source and dictionary. . . . .	42
E.2	Minimal encoder achieves the optimal cross-entropy loss up to a constant. . . . .	42
E.3	Greedy-encoder achieves poor cross-entropy loss . . . . .	44
<b>F</b>	<b>Experiment details</b>	<b>46</b>
<b>G</b>	<b>NeurIPS Paper Checklist</b>	<b>48</b>

## A Analysis of LZW: Proofs of Theorems 3.4 and 3.6

### A.1 Notation and definitions

For each character  $a \in \mathcal{A}$  let  $\mathcal{T}_a^*$  denote an infinite tree, with root vertex  $\emptyset$ , and subsequent vertices labelled by strings  $t \in \mathcal{A}^*$ . The edge from a parent vertex  $t$  to any child  $ta'$  is labelled with the probability  $P(ta'|t)$  unless  $t = \emptyset$ , in which case the edge probability is  $P(a'|a)$ . An infinite trajectory sampled on the tree  $\mathcal{T}_a^*$  corresponds to an infinite string sampled from the stochastic source conditioned on the first character of the string being  $a$ . In this paper we only consider ergodic sources (Gray and Gray, 2009) for which we can define the “entropy rate”. The entropy rate fundamentally captures the compressibility of the source, and can be defined as  $H_\infty \triangleq \lim_{m \rightarrow \infty} \frac{1}{m} H(P)$  where  $s$  is a length  $m$  string drawn from the source. By Theorem 2.1,  $H_\infty$  captures  $\min_Q \mathcal{L}(Q)$ .

### A.2 A basic result about the optimal achievable cross-entropy loss

The ratio of  $H(P)$  and  $mH(\pi)$  can be made arbitrarily large for the switching Markov chains in Figure 1 as the switching probabilities  $p$  and  $q$  approach 0 or 1. See Example A.1 for more details.

*Example A.1.* Consider the switching Markov process in Figure 1 on  $\{0, 1\}$  with  $p = q = 1 - \delta$ . For this process,  $\lim_{m \rightarrow \infty} \frac{1}{m} H(P) = H_{\text{Ber}}(\delta) = \delta \log(1/\delta) + (1 - \delta) \log(1/(1 - \delta))$ , but  $\pi = \{1/2, 1/2\}$  and so  $H(\pi) = H_{\text{Ber}}(1/2) = \log(2)$ . The ratio  $\lim_{m \rightarrow \infty} \frac{mH(\pi)}{H(P)}$  goes to  $\infty$  as  $\delta \rightarrow 0$ .

### A.3 Proof of Theorem 2.1

We first characterize the minimum achievable cross-entropy loss  $\mathcal{L}_m(Q)$  without any restrictions on the likelihood model class  $\mathcal{Q}$ . Choosing  $Q(\text{enc}(s)) = Q(s) = P(s)$ , the true probability of the sequence  $s$ , we have  $\mathcal{L}_m(Q \circ \text{enc}(\cdot)) = H(s)$  where  $H(\cdot)$  is the entropy function. It is not that difficult to see that this is also the minimum cross-entropy loss that can be achieved. For any distribution  $Q$ ,

$$\begin{aligned}\mathcal{L}_m(Q) &= \mathbb{E}[\log(1/Q(s))] \\ &= \mathbb{E}[\log(P(s)/Q(s))] + \mathbb{E}[\log(1/P(s))] \\ &= H(P) + D_{\text{KL}}(P\|Q).\end{aligned}$$

On the other hand, the cross-entropy loss under any unigram model  $Q \in \mathcal{Q}_{1\text{-gram}}$  satisfies,

$$\begin{aligned}\frac{1}{m}\mathcal{L}_m(Q \circ \text{enc}(\cdot)) &\stackrel{(i)}{=} -\frac{1}{m} \sum_{i=1}^m \mathbb{E}[\log Q_{\text{tok}}(t_i)] - \frac{1}{m} \mathbb{E}[\log Q_{\#}(m)] \\ &\stackrel{(ii)}{\geq} -\sum_{a \in \mathcal{A}} \pi(a) \log Q_{\text{tok}}(a) \\ &\geq H(\pi)\end{aligned}$$

where in (i), we use the definition of the unigram model  $Q$ , and in (ii),  $\pi$  is the stationary distribution over characters induced by the stochastic source, and the ergodicity of the source is used. The last equation lower bounds  $H(X, Y) \geq H(X)$ .

### A.4 Maximum likelihood unigram model

A number of our results (Theorems 3.4 and 3.6 to name a few) are related to bounding  $\min_{Q \in \mathcal{Q}_{1\text{-gram}}} \mathcal{L}(Q \circ \text{enc}(\cdot))$  for some tokenizer  $\mathcal{T}$ . In this section we introduce the maximum likelihood unigram model which captures the optimizer over  $Q$  for any given tokenizer.

For the character level tokenizer, an examination of Theorem 2.1 shows that the optimal unigram likelihood model associates probability  $Q_{\text{tok}}(a) = \pi(a)$ , i.e. the limiting fraction of times the character  $a$  is observed in the sequence. More generally, for a non-trivial tokenizer, the corresponding optimal unigram model  $Q_{\text{tok}}^*(t)$  ends up being the limiting expected fraction of times  $t$  is observed in an encoding of a sequence. This is the maximum likelihood unigram model, which we formally define below. The unigram MLE likelihood model associates probability,

$$Q_{\text{MLE}}(t) \leftarrow \lim_{m \rightarrow \infty} \mathbb{E} \left[ \frac{n_t}{\sum_t n_t} \right] \quad (4)$$

to each token, where  $n_t$  is the random variable capturing the number of occurrences of the token  $t$  in the encoding of the length- $m$  string  $s$ . Restricting the class of likelihood models to the unigram models,  $\mathcal{Q}_{1\text{-gram}}$ ,  $Q_{\text{MLE}}$  captures the model which minimizes eq. (1).

The unigram MLE model cannot be computed without an infinite amount of data, but can be approximated well with a finite amount of data, which forms the basis for Theorem C.1. For certain encoding algorithms, we can show that the quantity  $n_t / \sum_t n_t$  asymptotically converges to its expectation (Lemma A.4). This is the reason the unigram model in eq. (4) is referred to as a “maximum likelihood” model, since  $\lim_{m \rightarrow \infty} n_t / \sum_t n_t$  is the limit as  $|s| = m \rightarrow \infty$  of the solution to the following likelihood maximization problem: given a sequence  $s$ , find the distribution over tokens,  $Q$ , which maximizes

$$\prod_{t \in \text{enc}(s)} Q(t) \equiv \prod_{t \in \text{Dict}} (Q(t))^{n_t}.$$

As discussed previously, the unigram MLE model over tokens in eq. (4) induces a joint distribution over sequences of tokens by looking at the product of the marginal probabilities of the composed tokens; in particular,

$$Q_{\text{MLE}}(t_1, \dots, t_j) = Q_{\text{MLE}}(j) \prod_{i=1}^j Q_{\text{MLE}}(t_i),$$

where  $Q_{\text{MLE}}(j)$  is a distribution on the total number of tokens generated and is instantiated as  $\text{Unif}([m])$ .

**Remark A.2.** Note that the unigram MLE model specifies a distribution over tokens which is a function of the underlying encoding algorithm,  $\text{enc}(\cdot)$ . Different encoders result in different population level distributions over tokens, and consequently different unigram MLE models.

**Definition A.3** (greedy encoder). Given a dictionary  $\text{Dict}$ , the greedy encoder  $\text{enc}_{\text{gre}}(s)$  encodes a string  $s$  into tokens by greedily matching from left to right, the largest substring that exists as a token in  $\text{Dict}$ . This substring is then removed and the process iterated on the remainder of  $s$ . The greedy decoder  $\text{dec}_{\text{gre}}(\cdot)$  is a lookup table - a sequence of tokens is decoded by replacing each occurrence of a token by the corresponding substring it maps to in  $\text{Dict}$ .

**Lemma A.4.**  $\lim_{m \rightarrow \infty} \frac{n_t}{\sum_{t'} n_{t'}} \stackrel{\text{a.s.}}{=} \lim_{m \rightarrow \infty} \mathbb{E} \left[ \frac{n_t}{\sum_{t'} n_{t'}} \right]$  for any tokenizer having a finite vocabulary and finitely long tokens, using the greedy encoder.

*Proof.* This result is essentially true because under the greedy encoder, the tokens in an encoding of a fresh string  $t$  may be generated by an  $r^{\text{th}}$ -order Markov process for some  $r$ . For such processes, the Cesàro average of the state distributions converges to a stationary distribution of the process (i.e., the Krylov–Bogolyubov argument).

Tokens are generated as follows. Suppose the previous tokens generated were  $t_1, t_2, \dots, t_i$ . The next token  $t_{i+1}$  is sampled by drawing an infinite trajectory from  $\mathcal{T}_a^*$  for  $a \sim P(\cdot | t_i)$  and returning the longest prefix  $t$  of this trajectory which is a token in  $\text{Dict}$ , conditional on satisfying the conditions,  $t_j t_{j+1} \dots t_i t \notin \text{Dict}$  for all  $j \in \{1, 2, \dots, i\}$ . This process is repeated sequentially to generate all the tokens.

Suppose the length of the longest token in the dictionary is  $\ell_{\max}$ . Then, the distribution from which a token is sampled depends on at most the previous  $\ell_{\max}$  tokens. The reason for this is that the dependency of the  $(i+1)^{\text{th}}$  token,  $t_{i+1}$ , on the previously sampled tokens emerges in the constraint  $t_j t_{j+1} \dots t_i t_{i+1} \notin \text{Dict}$ , satisfied by any candidate  $t_{i+1}$ . Since each token is of length at least one, this condition is vacuously satisfied if  $j < i - \ell_{\max}$ .

With this view, the evolution of the state, defined as  $\text{state}_r = (t_{r\ell_{\max}}, t_{r\ell_{\max}-1}, \dots, t_{(r-1)\ell_{\max}})$  evolves in a Markovian fashion. By the Krylov–Bogolyubov argument (cf. Proposition 4.2 in Chen (2018)), the time averaged visitation frequencies of a Markov chain coordinate-wise asymptotically converges to its expectation, almost surely. This expectation exists by Theorems 8.5 and 8.22 of Eisner et al. (2015) which shows that for a matrix  $A$  such that  $\sup_{t \in \mathbb{N}} \|A^t\|_{\text{op}} < \infty$  the limit  $\lim_{t \rightarrow \infty} \frac{1}{t} \sum_{i=1}^t A^i$  exists. For the finite-state Markov transition  $A$  which captures the token generation process, condition  $\sup_{t \in \mathbb{N}} \|A^t\|_{\text{op}} \leq |\text{Dict}|^{\ell_{\max}} < \infty$ . This means that the limit of the time averaged state distribution exists. Moreover, for any initial distribution  $\pi_0$  over tokens,  $\pi = \lim_{t \rightarrow \infty} \frac{1}{t} \sum_{i=1}^t \pi_0 A^i$  satisfies the condition  $\pi A = \pi$ , implying that the limiting time-averaged state distribution is a stationary distribution of  $A$ . Since the limiting time-averaged measure on the state  $\text{state}_r = (t_{r\ell_{\max}}, \dots, t_{r\ell_{\max}-1}, \dots, t_{(r-1)\ell_{\max}})$  exists, this implies that the limiting time-averaged measure of  $t_{r\ell_{\max}-r'}$  for each  $r' \in \{0, 1, \dots, \ell_{\max}\}$  exists. By taking the uniform average over  $r'$  and  $r$ , the limiting time-averaged measure of  $t_i$  over  $i \in \mathbb{N}$  exists.  $\square$

## A.5 Proof of Theorem 3.4

Consider a string  $s$  of length  $m \rightarrow \infty$  which is encoded into a sequence of tokens  $(t_i : i \in [|\text{enc}_{\text{gre}}(s)|])$ . By the Asymptotic Equipartition Property (AEP) for ergodic sources, i.e. the Shannon–McMillan–Breiman theorem,

$$\Pr \left( \lim_{m \rightarrow \infty} -\frac{1}{m} \log P(s) = H_{\infty} \right) = 1. \quad (5)$$

Here  $\lim_{m \rightarrow \infty} \frac{H(P)}{m}$  also happens to be the entropy rate of the source. We use this property to bound the length of the greedy encoding,  $|\text{enc}_{\text{gre}}(s)|$ . Indeed, the probability of  $s$  may be decomposed as,

$$P(s) = P(t_1) \prod_{i=2}^{|\text{enc}_{\text{gre}}(s)|} P(t_i | t_{i-1}) \leq \prod_{i=1}^{|\text{enc}_{\text{gre}}(s)|} \max_{a \in \mathcal{A}} P(t_i | a).$$

Noting that  $\delta \min_a P(\mathbf{t}|a) \geq \max_a P(\mathbf{t}|a)$ , up to a  $\delta$  factor we may replace the max over  $a$  by an expectation over  $a \sim \pi$  where  $\pi$  is the stationary distribution of the stochastic source. In particular,

$$P(\mathbf{s}) \leq \prod_{i=1}^{|\text{enc}_{\text{gre}}(\mathbf{s})|} P(\mathbf{t}_i)/\delta.$$

By invoking the AEP, eq. (5),

$$\lim_{m \rightarrow \infty} \frac{1}{m} \sum_{i=1}^{|\text{enc}_{\text{gre}}(\mathbf{s})|} -\log(P(\mathbf{t}_i)) - \log(1/\delta) \stackrel{\text{a.s.}}{\leq} H_\infty$$

Recall that the greedy encoder satisfies Lemma A.4 and for any  $\mathbf{t} \in \text{Dict}$ ,  $\lim_{m \rightarrow \infty} \frac{n_{\mathbf{t}}}{|\text{enc}_{\text{gre}}(\mathbf{s})|} \stackrel{\text{a.s.}}{=} Q_{\text{MLE}}(\mathbf{t})$ . Furthermore, note that for any token  $\mathbf{t} \in \text{Dict}$ ,  $P(\mathbf{t}) > \delta^{|\mathbf{t}|} > 0$ , and  $|\text{enc}_{\text{gre}}(\mathbf{s})| \leq m$  surely. By almost sure convergence,

$$\begin{aligned} \lim_{m \rightarrow \infty} \frac{|\text{enc}_{\text{gre}}(\mathbf{s})|}{m} \sum_{\mathbf{t} \in \text{Dict}} -\frac{n_{\mathbf{t}}}{|\text{enc}_{\text{gre}}(\mathbf{s})|} \left( \log(P(\mathbf{t})) - \log(1/\delta) \right) \\ \stackrel{\text{a.s.}}{=} \lim_{m \rightarrow \infty} \frac{|\text{enc}_{\text{gre}}(\mathbf{s})|}{m} \left( H(Q_{\text{MLE}}, P) - \log(1/\delta) \right) \end{aligned}$$

Furthermore, utilizing the assumption that  $\varepsilon H(Q_{\text{MLE}}, P) \geq \log(1/\delta)$  satisfied by the tokenizer,

$$\lim_{m \rightarrow \infty} \frac{(1-\varepsilon)|\text{enc}_{\text{gre}}(\mathbf{s})| \left( H(Q_{\text{MLE}}, P) \right)}{m} \stackrel{\text{a.s.}}{\leq} H_\infty. \quad (6)$$

Now we are ready to bound the expected cross-entropy loss of the tokenizer. Define the unigram model  $P_\pi(\mathbf{t}_1, \mathbf{t}_2, \dots, \mathbf{t}_j) = P_{\text{unif}}(j) \prod_{i=1}^j P(\mathbf{t}_i)$  where  $P_{\text{unif}}$  is the uniform measure over  $[m]$ . Note that we have the inequality  $\min_{Q \in \mathcal{Q}_{1\text{-gram}}} \lim_{m \rightarrow \infty} \frac{1}{m} \mathcal{L}_m(Q \circ \text{enc}_{\text{gre}}(\cdot)) \leq \lim_{m \rightarrow \infty} \frac{1}{m} \mathcal{L}_m(P_\pi \circ \text{enc}_{\text{gre}}(\cdot))$  and therefore, it suffices to upper bound the RHS. In particular,

$$\begin{aligned} \mathcal{L}_m(P_\pi \circ \text{enc}_{\text{gre}}(\cdot)) &= -\mathbb{E}[\log P_{\text{unif}}(|\text{enc}_{\text{gre}}(\mathbf{s})|)] - \mathbb{E} \left[ \sum_{\mathbf{t} \in \text{enc}_{\text{gre}}(\mathbf{s})} \log(P(\mathbf{t})) \right] \\ &\leq \log(m) - \mathbb{E} \left[ \sum_{\mathbf{t} \in \text{enc}_{\text{gre}}(\mathbf{s})} \log(P(\mathbf{t})) \right] \end{aligned} \quad (7)$$

where the last inequality uses the fact that  $P_{\text{unif}}(|\text{enc}_{\text{gre}}(\mathbf{s})|) = 1/m$ . Note that as  $m \rightarrow \infty$ , by assumption on the tokenizer, the fraction of times the token  $\mathbf{t}$  appears in the encoding of  $\mathbf{s}$  converges almost surely to  $Q_{\text{MLE}}(\mathbf{t})$ . Since  $|\text{enc}_{\text{gre}}(\mathbf{s})| \leq m$  surely and  $P(\mathbf{t}) > \delta^{|\mathbf{t}|} > 0$ , by an application of the Dominated Convergence Theorem,

$$\begin{aligned} -\lim_{m \rightarrow \infty} \frac{1}{m} \mathbb{E} \left[ \sum_{\mathbf{t} \in \text{enc}_{\text{gre}}(\mathbf{s})} \log(P(\mathbf{t})) \right] &= -\lim_{m \rightarrow \infty} \frac{1}{m} \mathbb{E} \left[ |\text{enc}_{\text{gre}}(\mathbf{s})| \cdot \sum_{\mathbf{t} \in \text{Dict}} Q_{\text{MLE}}(\mathbf{t}) \log(P(\mathbf{t})) \right] \\ &= \lim_{m \rightarrow \infty} \frac{1}{m} \mathbb{E} [|\text{enc}_{\text{gre}}(\mathbf{s})| H(Q_{\text{MLE}}, P)] \end{aligned} \quad (8)$$

Combining eq. (7) with eq. (8) and setting  $\lim_{m \rightarrow \infty} \log(m)/m = 0$ , and invoking eq. (6),

$$\begin{aligned} \min_{Q \in \mathcal{Q}_{1\text{-gram}}} \lim_{m \rightarrow \infty} \frac{1}{m} \mathcal{L}_m(Q_{\text{MLE}} \circ \text{enc}_{\text{gre}}(\cdot)) &= \lim_{m \rightarrow \infty} \frac{1}{m} \mathbb{E} [|\text{enc}_{\text{gre}}(\mathbf{s})| H(Q_{\text{MLE}}, P)] \\ &\leq \frac{H_\infty}{1-\varepsilon}. \end{aligned} \quad (9)$$

By Theorem 2.1, we have that  $\min_Q \lim_{m \rightarrow \infty} \frac{1}{m} \mathcal{L}_m(Q \circ \text{enc}_{\text{gre}}(\cdot)) = \lim_{m \rightarrow \infty} \frac{H(P)}{m} = H_\infty$ , which uses the fact that the source is ergodic. Combining with eq. (9) completes the proof.

## A.6 Heavy-hitter dictionaries and a proof of Theorem 3.6

In this section we prove Theorem 3.6 and introduce the notion of a heavy-hitting dictionary. At a high level, these dictionaries contain all the substrings which have reasonably high probability of being observed many times in a dataset of size  $n = \tilde{\Omega}_\delta(d)$ . We first show in Lemma A.6 that heavy hitting dictionaries generalize well in the sense of having  $H(Q_{\text{MLE}}, P)$  being large (in conjunction with Theorem 3.4 this implies an upper bound on the cross-entropy loss of the best unigram model). Next, we will prove that the LZW algorithm (Definition 3.5) results in a heavy hitting dictionary with high probability.

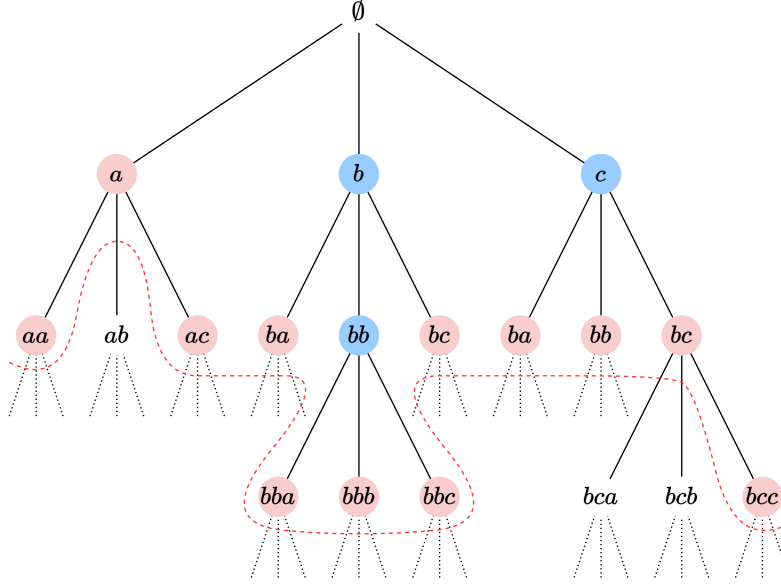


Figure 6: The circled nodes indicates substrings which are tokens in Dict. Red nodes indicate the set of “maximal tokens”, which are the set of tokens which the greedy encoder assigns, leaving out those which can only be assigned as the last token of some string. Tokens like “b” are never assigned by the greedy encoder (save as the last token of the encoding of a string) since any sufficiently long trajectory starting with  $b$  must have a longer prefix which is also a token, namely, one of  $ba$ ,  $bc$ ,  $bba$ ,  $bbb$  or  $bbc$ . The vertices of the tree which are assigned by the greedy encoder as tokens (together with all their prefixes) forms a cut of the tree, which marks the dotted red line. The heavy hitting property asserts that this cut is uniformly far away from the root node  $\emptyset$ , and that every vertex  $s$  marked red has  $P(s) \leq 1/d^\beta$ .

**Definition A.5** ( $\beta$ -heavy-hitting dictionary). A token  $t$  of a dictionary is said to be maximal if there exists an arbitrary substring containing  $t$  as a strict prefix, and in addition,  $t$  is also the largest prefix of the substring which is a token. A dictionary Dict is said to be  $\beta$ -heavy hitting if the set of maximal tokens is a subset of  $\{s' : \max_{a \in \mathcal{A}} P(s'|a) \leq 1/d^\beta\}$ .

A pictorial depiction of the heavy hitting property is illustrated in Figure 6.

**Lemma A.6.** For a  $\beta$ -heavy-hitting dictionary, with the greedy encoder,  $H(Q_{MLE}, P) \geq \beta \log(d)$ .

*Proof.* Note that the greedy encoder assigns tokens only among the set of maximal substrings (save for potentially the last token). If every maximal substring has  $\max_{a \in \mathcal{A}} P(s|a) \leq 1/d^\beta$ , by the heavy-hitting property, for any token  $t$ ,

$$P(t) \leq \max_{a \in \mathcal{A}} P(s'|a) \leq 1/d^\beta.$$

Therefore,

$$H(Q_{MLE}, P) = \mathbb{E}_{t \sim Q_{MLE}} [\log(1/P(t))] \geq \beta \log(d).$$

□

Define  $\mathcal{M}_\beta = \{t : \max_{a \in \mathcal{A}} P(t|a) \geq \delta/d^\beta\}$ . These are the set of “high-probability” substrings under the stochastic source. We will show that for  $\beta$  bounded away from 1, with high probability, every substring in  $\mathcal{M}_\beta$  is added as a token to the dictionary in a run of the LZW tokenizer (Definition 3.5). Note that if every substring in  $\mathcal{M}_\beta$  is assigned as a token by LZW, then the algorithm must be  $\beta$ -heavy hitting since there always exists a maximal token on the “boundary” of the set  $\mathcal{M}_\beta$  which is strictly contained in  $\{s' : \max_{a \in \mathcal{A}} P(s'|a) \leq 1/d^\beta\}$ .

**Lemma A.7.** Every substring in  $\mathcal{M}_\beta$  has length at most  $\ell_\star \triangleq \delta^{-1}(\beta \log(d) + \log(1/\delta))$ .

*Proof.* Note that  $\min_{a,a' \in \mathcal{A}} P(a|a') = \delta$ , which implies that the probability of any transition must be bounded away from 1, i.e.,  $\max_{a,a' \in \mathcal{A}} P(a|a') \leq 1 - \delta$ . This implies that,

$$\max_{a \in \mathcal{A}} P(\mathbf{t}|a) \leq (1 - \delta)^{|\mathbf{t}|} \leq e^{-\delta|\mathbf{t}|}. \quad (10)$$

By definition, for any substring  $\mathbf{t} \in \mathcal{M}_\beta$ ,  $\max_{a \in \mathcal{A}} P(\mathbf{t}|a) \geq \delta/d^\beta$ . In conjunction with eq. (10), this implies the statement of the lemma.  $\square$

In the remainder of this section, let  $n$  be the size of the dataset on which LZW is run. We show that the number of tokens added to the dictionary by LZW,  $d$ , is  $\tilde{\Theta}_\delta(n)$ . Rather than running the algorithm with early stopping (i.e., ceasing to add new tokens once the budget is hit), instead, we assume that the algorithm runs on a prefix of the dataset of length  $d$ . The number of tokens added this way cannot exceed  $d$ .

**Lemma A.8.** With probability  $\geq 1 - d^{-\Omega(\log(d/\delta)/\delta)}$ , in a run of the LZW algorithm, no substring  $\mathbf{t}$  added as a token to the dictionary satisfies  $|\mathbf{t}| \geq \ell_{\max} \triangleq 4 \log(d|\mathcal{A}|)/\delta$ .

*Proof.* Consider any  $s \in \mathbb{N}$  and any substring  $\mathbf{t}$  of length  $s$ . In order for  $\mathbf{t}$  to be assigned as a token, each of its prefixes must disjointly appear at least once in the string. Since there are at most  $d$  tokens, we can upper bound the probability that  $\mathbf{t}$  is assigned as a token as,

$$\begin{aligned} P(\mathbf{t} \text{ is assigned as a token}) &\leq \binom{d}{s} \prod_{i=1}^s \max_{a \in \mathcal{A}} P(\mathbf{t}_{1:i}|a) \\ &\stackrel{(i)}{\leq} \binom{d}{s} (1 - \delta)^{s(s-1)/2} \\ &\leq e^{s \log(d) - \delta s(s-1)/2}, \end{aligned}$$

where (i) uses the fact that  $\max_{a \in \mathcal{A}} P(\mathbf{t}_{1:i}) \leq \prod_{j=1}^i \max_{a \in \mathcal{A}} P(\mathbf{t}_j|a) \leq (1 - \delta)^i$ . By union bounding across the  $|\mathcal{A}|^s$  strings of length  $s$ ,

$$P(\text{any length } s \text{ string is assigned as a token}) \leq e^{s \log(|\mathcal{A}|) + s \log(d) - \delta s(s-1)/2}.$$

When  $s = 4 \log(d|\mathcal{A}|)/\delta + 1 \triangleq \ell_{\max} + 1$ , the RHS is upper bounded by  $e^{-\delta \ell_{\max}^2/4} \leq d^{-\Omega(\log(d/\delta)/\delta)}$ . With the same small probability, no substring of length  $s' > s$  can become a token, since their length- $s$  prefixes are never assigned as tokens.  $\square$

**Corollary A.9.** With probability  $\geq 1 - d^{-\Omega_s(\log(d))}$ , learns a dictionary with at least  $d^\star = d/\ell_{\max}$  tokens when run on a training sequence of length  $n$  drawn from a stochastic source satisfying Assumption 3.2.

**Lemma A.10.** For any constant  $\beta < 1$ , with probability  $\geq 1 - d^{-\Omega(\log(d/\delta)/\delta)} - \exp(-\tilde{\Omega}_\delta(d^{1-\beta}))$  over the source dataset, every substring in  $\mathcal{M}_\beta$  is added as a token to the dictionary in a run of the LZW algorithm. In other words, with the same probability, the LZW tokenizer results in a  $\beta$ -heavy hitting dictionary.

By Corollary A.9, note that with high probability the LZW tokenizer adds at least  $d^\star$  tokens to the dictionary when processing a length  $d$  training sequence in entirety. In this proof, instead of generating  $d$  samples, we sequentially sample  $d^\star$  tokens from their joint distribution, and generate a dictionary from these samples. From Corollary A.9, with high probability this results in at most  $d$  samples being generated, implying that the dictionary generated by sampling  $d^\star$  tokens is a subset of the dictionary generated by a full run of the LZW tokenizer. Here, we use the fact that the LZW tokenizer adds tokens to the dictionary in a left to right fashion, and therefore a subset of the dictionary learnt by the LZW tokenizer can be generated by processing a portion of the dataset.

Next we consider a joint view for generating the dataset from the stochastic source and the dictionary learnt by LZW simultaneously. The stochastic source is sampled as a sequence of tokens. Suppose the last character of the previous token was  $a'$ . Sample a character  $a \sim P(\cdot|a')$  and an infinite trajectory on the tree  $\mathcal{T}_{a'}^\star$ . Consider the first node visited in this trajectory which does not already

exist as a token in the dictionary. The substring corresponding to this node is added as a token in the dictionary. By repeating this process, the dictionary and the source dataset are constructed sequentially and simultaneously. As alluded to before, we truncate this token sampling process to repeat at most  $d^*$  times, which results in a subset of the dictionary output by the LZW algorithm with high probability (Corollary A.9). This is simply a variant of the ‘‘Poissonization’’ trick to avoid statistical dependencies across tokens. Denote the set of infinite trajectories generated on the forest  $\{\mathcal{T}_a^* : a \in \mathcal{A}\}$  as  $\{\text{traj}_i : i \in [d^*]\}$ .

With this view of the sampling process, observe that if the substring  $t$  sampled was a prefix of  $\text{traj}_i$  at least  $|t|$  times across different values of  $i$ , then  $t$  must be assigned as a token. In particular, in each of these  $|t|$  trajectories, each of the prefixes of  $t$  is assigned as a token. With this observation, the event that  $t$  is not assigned as a token is contained in the event that  $t$  is visited at most  $|t| - 1$  times across the  $d^*$  trajectories. Observe that,

$$P(t \text{ is not assigned as a token}) \leq \sum_{i=0}^{|t|-1} \binom{d^*}{i} \max_{a \in \mathcal{A}} (P(t|a))^i (1 - P(t|a))^{d^*-i}.$$

Since we aim to upper bound this probability across the substrings in  $t \in \mathcal{M}_\beta$ , note that (i)  $\max_{a \in \mathcal{A}} P(t|a) \geq \delta/d^\beta$ , and (ii) tokens in  $\mathcal{M}_\beta$  have length at most  $\ell_\star = \delta^{-1}(\beta \log(d) + \log(1/\delta))$  (Lemma A.7), implying there are at most  $2|\mathcal{A}|^{\ell_\star}$  substrings in this set. By union bounding,

$$P(\exists t \in \mathcal{M}_\beta \text{ not assigned as a token}) \leq 2|\mathcal{A}|^{\ell_\star} \sum_{i=0}^{\ell_\star-1} \binom{d^*}{i} \max_{x \geq \delta/d^\beta} x^i (1-x)^{d^*-i}. \quad (11)$$

**Case I.** For  $i \leq \ell_\star$  and  $x \geq 1/2$ ,

$$\begin{aligned} |\mathcal{A}|^{\ell_\star} \binom{d^*}{i} x^i (1-x)^{d^*-i} &\leq |\mathcal{A}|^{\ell_\star} \frac{(d^*)^{\ell_\star}}{2^{d^*/2}} \\ &\leq 2^{\ell_\star \log(d^*|\mathcal{A}|) - d^*/2} \\ &\leq 2^{-\Omega_{\beta,\delta}(d^*)}, \end{aligned} \quad (12)$$

where the last inequality uses the fact that  $\ell_\star = O_{\beta,\delta}(\log(d))$ .

**Case II.** For  $i \leq \ell_\star$  and  $\delta/d^\beta \leq x \leq 1/2$ ,

$$\begin{aligned} |\mathcal{A}|^{\ell_\star} \binom{d^*}{i} x^i (1-x)^{d^*-i} &\leq |\mathcal{A}|^{\ell_\star} \binom{d^*}{i} (1-x)^{d^*} \\ &\leq |\mathcal{A}|^{\ell_\star} (d^*)^{\ell_\star} e^{-d^* x} \\ &\leq e^{\ell_\star \log(|\mathcal{A}|) + \ell_\star \log(d^*) - d^* x} \\ &\leq e^{-\Omega(\delta^2 n / d^\beta / \log(d/\delta))} \\ &\leq e^{-\Omega(\delta^2 d^{1-\beta} / \log(d/\delta))}, \end{aligned} \quad (13)$$

where the last inequality uses the fact that  $\ell_\star = O(\log(d))$ ,  $x \geq \delta/d^\beta$ ,  $d^* = \Omega(d\delta/\log(d/\delta))$ . By combining eq. (12) and eq. (13) with eq. (11) completes the proof, as long as  $\beta$  is a constant bounded away from 1.

**Lemma A.11.** Fix a constant  $\gamma > 0$ . Then, with probability  $\geq 1 - d^{-\Omega_{\gamma,\delta}(\log(d))}$ , none of the substrings in the set  $\mathcal{N}_\gamma = \{s' : \max_{a \in \mathcal{A}} P(s'|a) \leq \delta/d^{1+\gamma}\}$  are assigned as tokens in a run of LZW.

*Proof.* Define the following set of substrings,

$$S_\gamma = \left\{ t : \delta/d^{1+\gamma/2} \leq \max_{a \in \mathcal{A}} P(t|a) \leq 1/d^{1+\gamma/2} \right\}$$

Since the width of this band is sufficiently large, by Assumption 3.2 every substring  $t$  such that  $\max_{a \in \mathcal{A}} P(t|a) \leq \delta/d^{1+\gamma/2}$  has at least one prefix which falls in  $S_\gamma$ , and denote the longest such

prefix  $t_\gamma$ . Define  $T_\gamma = \{t_\gamma : t \in \mathcal{N}_\gamma\}$  as the set of longest prefixes in  $S_\gamma$ . Intuitively, if we think of the strings in  $S_\gamma$  (or  $T_\gamma$ ) as being intermediate in length, the strings in  $\mathcal{N}_\gamma$  can be thought of as being particularly long: the value of  $\max_{a \in \mathcal{A}} P(t|a)$  for any  $t \in T_\gamma$  and for any  $t \in \mathcal{N}_\gamma$  are separated by a factor of at least  $1/d^{\gamma/2}$ . In particular, since the probability of any character is lower bounded by  $\delta$ , each substring in  $t \in \mathcal{N}_\gamma$  must be at least  $\Delta = \frac{\gamma \log(d)}{2 \log(1/\delta)}$  symbols longer than its corresponding longest prefix in  $T_\gamma$ ,  $t_\gamma$ . An implication of this is that for  $t$  to be assigned as a token,  $t_\gamma$  must be observed at least  $\Delta + 1$  times disjointly in  $s$ . However, note that  $t_\gamma$  already has low marginal probability to begin with ( $\ll 1/d$ ) so the odds of seeing this substring so many times disjointly is very small. Furthermore, note that  $T_\gamma$  has at most  $d^{1+\gamma/2}/\delta$  substrings, which allows the probability of this event occurring simultaneously across all substrings in  $T_\gamma$  to be controlled by union bound. Under this condition, none of the substrings in  $\mathcal{N}_\gamma$  are made into tokens.

In order to argue that the dictionary *does not* contain certain tokens, we may argue this property about any superset of the dictionary. In contrast, in Lemma A.10, we construct a subset of the dictionary by running LZW on the concatenation of  $d^*$  tokens sampled from their joint distribution. The superset we consider here is just to sample  $d$  tokens from their joint distribution and concatenate them together to result in a string of length  $\geq d$ , and running LZW on this sequence (which simply would result in these  $d$  tokens). As in Lemma A.10, let  $\{\text{traj}_i : i \in [d]\}$  denote the infinite trajectories generated from the Markov chain which are truncated to result in tokens. A sufficient condition for the event that no substring  $t \in \mathcal{N}_\gamma$  is assigned as a token by LZW is to the event that every substring  $t' \in T_\gamma$  is observed as a prefix of  $\text{traj}_i$  for  $\Delta$  or fewer choices of  $i \in [d]$ . To this end define  $\mathcal{E}(t')$  as the event that  $|i \in [d] : t' \text{ is a prefix of } \text{traj}_i| \leq \Delta$ . Then,

$$\begin{aligned} \Pr(\mathcal{E}(t')) &\leq \binom{n}{\Delta} (\max_{a \in \mathcal{A}} P(t'|a))^\Delta \\ &\stackrel{(i)}{\leq} e^{\Delta \log(n)} \left( \frac{1}{d^{1+\gamma/2}} \right)^\Delta \\ &\leq e^{-\frac{\gamma}{2} \Delta \log(d)}, \end{aligned} \tag{14}$$

where (i) uses the fact that  $\max_{a \in \mathcal{A}} P(t'|a) \leq 1/d^{1+\gamma/2}$  since the substring  $t'$  belongs to  $T_\gamma$ .

Note that the number of substrings in  $S_\gamma$  (and by extension,  $T_\gamma$ ) is at most  $O_\delta(d^{1+\gamma/2})$ . Recall that these substrings satisfy the condition  $\max_{a \in \mathcal{A}} P(t|a) \geq \delta/d^{1+\gamma/2}$ . Observe that,

$$\begin{aligned} \frac{\delta |S_\gamma|}{d^{1+\gamma/2}} &\leq \sum_{t \in S_\gamma} \max_{a \in \mathcal{A}} P(t|a) \\ &\leq \sum_{t \in S_\gamma} \sum_{a \in \mathcal{A}} P(t|a) \leq |\mathcal{A}| \leq \frac{1}{\delta}. \end{aligned}$$

This implies that there are at most  $d^{1+\gamma/2}/\delta^2$  substrings in  $S_\gamma$ . Finally, in conjunction with eq. (14),

$$P(\exists t' \in S_\gamma : \mathcal{E}(t')) \leq \frac{d^{1+\gamma/2}}{\delta^2} e^{-\frac{\gamma}{2} \Delta \log(d)} = d^{-\Omega_{\gamma, \delta}(\log(d))},$$

which implies that with high probability, no token in  $S_\gamma$  is observed as a prefix of  $s^i$  for more than  $\Delta$  choices of the index  $i \in [d]$ . Under this event, no substring in  $\mathcal{N}_\gamma$  is assigned as a token.  $\square$

### A.6.1 Proof of Theorem 3.6

Choosing  $\beta = 0.99$  in Lemma A.10, with probability  $\geq 1 - d^{-\Omega_\delta(\log(d))}$ , the LZW tokenizer results in a 0.99-heavy-hitting dictionary. As a consequence of Lemma A.6, this implies that under the same event,

$$H(Q_{\text{MLE}}, P) \geq 0.99 \log(d).$$

Finally, combining with Theorem 3.4 completes the proof.



## B Additional Theoretical Results I: A sequential variant of BPE

While the main results in the paper focused on understanding the limits of tokenization under a bound on the dictionary size, in this section we take a more practical look and try to analyze tokenizers used commonly in practice. The Byte-Pair-Encoding (BPE) algorithm (Gage, 1994; Sennrich et al., 2016), discovered in the compression literature as REPAIR (Larsson and Moffat, 2000; Navarro and Russo, 2008) was proposed as a faster alternative to LZW. It remains as one of the most commonly implemented tokenizers in natural language processing for various downstream tasks (Radford et al., 2019; Mann et al., 2020; Touvron et al., 2023). A large proportion of open source and commercial LLMs currently use BPE as the tokenization algorithm of choice, such as GPT-2/3, Llama 1/2 and Mistral-7B to name a few.

The BPE algorithm is based on constructing the dictionary iteratively by merging pairs of tokens to result in a tokens. In each iteration, the pair of tokens which appear most frequently next to each other are merged together into a single token. Subsequently, every occurrence of the pair of tokens are replaced by the newly added token, breaking ties arbitrarily. The dictionary is thus an ordered mapping of the form  $t \leftarrow (t', t'')$ . To encode a new string, the BPE encoder iterates through the dictionary and for each rule  $t \leftarrow (t', t'')$  replaces every consecutive occurrence of  $t'$  and  $t''$  by the token  $t$  breaking ties arbitrarily.

To warm up our main results, it is worth understanding the behavior of the BPE tokenizer in a bit more detail. Unlike the toy tokenizer, it is a priori unclear whether unigram models trained on sequences tokenized by BPE even asymptotically (in the dictionary size) achieve the optimal cross-entropy loss. Indeed, for  $\delta > 0$ , consider a training sequence of length  $m$  of the form,

$$s = \underbrace{\left( \underbrace{01 \cdots 01}_{2/\delta} \underbrace{10 \cdots 10}_{2/\delta} \right)}_{\times \frac{m\delta}{4}} \quad (15)$$

The probability that this sequence is generated by the order-2 switching Markov source with  $p = q = \delta$  is,

$$\approx (1 - \delta)^{\frac{m\delta}{4} \times \frac{4}{\delta} \times (1-\delta)} (\delta)^{\frac{m\delta}{4} \times 4} = e^{-H(P)},$$

which uses the fact that  $H(P) = m\delta \log(1/\delta) + m(1 - \delta) \log(1/(1 - \delta))$ . This implies that even though the string has exponentially small probability, it is one of the typical sequences for this order-2 Markov source. Let's understand what happens when the BPE tokenizer is trained on this dataset. Assuming that ties are broken arbitrarily, consider the run of the BPE algorithm detailed in Table 2. Here, we assume that  $1/\delta - 1$  is a power of 2 and denote  $r = \log_2(1/\delta - 1)$ . The algorithm first merges 0 and 1 into a single token  $t_1$ , which results in a long sequence of the form  $t_1 \cdots t_1 1 t_1 \cdots t_1 0$  repeated  $m\delta/4$  times. In subsequent rounds, the tokens  $(t_1, t_1)$  is merged into  $t_2$ , then  $(t_2, t_2)$  is merged into  $t_3$  and so on, until is no longer possible. Finally, the resulting sequence is a repeating sequence of 5 tokens where within each sequence, no pair of tokens appears more than once next to each other. Eventually these 5 tokens are merged into a single token labelled  $t_{r+4}$ , and in subsequent rounds the tokens  $(t_{r+4}, t_{r+4})$  are merged into  $t_{r+5}$ ,  $(t_{r+5}, t_{r+5})$  is merged into  $t_{r+6}$  and so on, until is no longer possible.

Observe that in the initial training dataset the substrings 0000 and 1111 never appears as a contiguous sequence. However, in a test sequence of length  $m$  sampled from the 2<sup>nd</sup>-order Markov source, with high probability these substrings disjointly occur  $\Theta(m)$  times each. The learnt dictionary associates each such disjoint occurrence of these substrings with at least 1 token, for 0000, the 3<sup>rd</sup> 0 must necessarily be tokenized as the token "0". Likewise, in 1111, the 3<sup>rd</sup> 1 must necessarily be tokenized as the token "1". Therefore, when a new test string of length  $m$  is tokenized, with high probability the tokens "0" and "1" form a constant fraction of the total collection of tokens.

Thus on freshly sampled test sequences, the BPE tokenizer appears to behave like the character-level tokenizer on a constant fraction of the input sequence. In particular, a simple calculation shows that the cross-entropy loss of any unigram model trained on this tokenizer must be far from the optimal

Initial	01.....0110.....10 ...
$t_1 \leftarrow (0, 1)$	$t_1 \dots t_1 \mathbf{1} t_1 \dots t_1 \mathbf{0}   \dots$
$t_2 \leftarrow (t_1, t_1)$	$t_2 \dots t_2 \mathbf{1} t_2 \dots t_2 \mathbf{0}   \dots$
$\vdots$	$\vdots$
$t_r \leftarrow (t_{r-1}, t_{r-1})$	$t_r t_1 \mathbf{1} t_r \mathbf{0}   \dots$
$t_{r+1} \leftarrow (t_r, t_1)$	$t_{r+1} \mathbf{1} t_r \mathbf{0}   \dots$
$t_{r+2} \leftarrow (t_r, 0)$	$t_{r+1} \mathbf{1} t_{r+2}   \dots$
$t_{r+3} \leftarrow (t_{r+1}, 1)$	$t_{r+3} t_{r+2}   \dots$
$t_{r+4} \leftarrow (t_{r+3}, t_{r+2})$	$t_{r+4}   \dots$
$t_{r+5} \leftarrow (t_{r+4}, t_{r+4})$	$t_{r+5}   \dots$
$t_{r+6} \leftarrow (t_{r+5}, t_{r+5})$	$t_{r+6}   \dots$
$\vdots$	$\vdots$

Table 2: A representation of the behavior of BPE when trained on the dataset in eq. (15). We assume that  $1/\delta - 1$  is a power of 2 and define  $r = \log_2(1/\delta - 1)$ .

bound of  $mH_{\text{BER}}(\delta)$  especially as  $\delta$  becomes smaller,

$$\begin{aligned}
& \min_{Q \in \mathcal{Q}_{1\text{-gram}}} \mathcal{L}_m(Q \circ \text{enc}(\cdot)) \\
& \geq \min_{Q \in \mathcal{Q}_{1\text{-gram}}} \mathbb{E} [n_0 \log(1/Q_{\text{tok}}(0)) + n_1 \log(1/Q_{\text{tok}}(1))] \\
& \stackrel{(i)}{\geq} \Omega(m) \cdot \min_{Q \in \mathcal{Q}_{1\text{-gram}}} (\log(1/P_{\text{tok}}(0)) + \log(1/Q_{\text{tok}}(1))) \\
& \geq \Omega(m).
\end{aligned}$$

where (i) uses the fact that  $\mathbb{E}[n_0], \mathbb{E}[n_1] \in \Omega(m)$  and the last inequality uses  $P_{\text{tok}}(0)P_{\text{tok}}(1) \leq 1/4$  (AM-GM inequality) since they sum up to at most 1. The purpose of this example is to show that there exist pathological training datasets which appear to be drawn from a stochastic source, but on which BPE fails to learn a good dictionary for the source. Thus proving a result such as Theorem 3.1 for BPE would require arguing that training datasets such as that in eq. (15) are unlikely to be seen.

The analysis of the standard variant of BPE turns out to be complicated for other reasons too. After every token is added the training dataset becomes a mix of all the previously added tokens, and arguing about the statistics of which pair of tokens appears most frequently for the next iteration becomes involved. For instance, adding 00 as a token may reduce the frequency of occurrence of the substring 01, but will not affect 11. Thus, even though 01 may a priori have been seen more frequently, it may not be chosen by BPE as the next token after 00.

To avoid dealing with these issues, we consider a sequential/sample-splitting variant of BPE. At a high level, the algorithm breaks down a dataset of size  $\Theta(d^2)$  into  $d$  chunks and learns at most 1 token from each chunk. The algorithm iterates over the chunks and finding the pair of tokens which appear most frequently next to each other in each chunk and adding it to the dictionary if it appears more than  $\log(d)$  times. Every consecutive occurrence of the pair of tokens is replaced by the newly assigned token in the dataset. Thus, in each iteration  $i$ , at most 1 token is added, depending on the statistics of the  $i^{\text{th}}$  chunk and the tokens added so far to the dictionary. Based on the final size of the dictionary a different encoder/decoder pair is used - if the algorithm adds sufficiently many tokens to the dictionary, the greedy encoder is used, and if not, a parallel implementation of BPE's encoding algorithm is used (Definition B.1). A formal description of the algorithm is in Algorithm 1.

**Definition B.1** (BPE.split encoder). The BPE.split encoder parses a new string into tokens as follows. The algorithm partitions the string into contiguous chunks of length  $d$ . Then, BPE's encoder is applied on each chunk, which iterates through DS and replaces  $t't''$  by  $t$  for every rule  $t \leftarrow (t', t'')$  in DS, breaking ties arbitrarily. The individual token sequences are finally spliced together and returned.

The main result of this section is that up to a small additive error, Algorithm 1 approaches a 2-approximation to the optimal cross-entropy loss.

---

**Algorithm 1** Sequential implementation of BPE

---

**Input:**  $\epsilon \in (0, 1)$ ; a dataset of size  $n = \Theta(d^2)$ , split into  $d$  contiguous texts  $\{\text{text}_1, \dots, \text{text}_d\}$  of length  $\Theta(d)$  each.  
**Output:** A tokenizer  $\mathcal{T}$ .  
*// Generate Dictionary*  
**for**  $i = 1, \dots, d$  **do**  
  **if**  $\exists$  a pair of tokens/characters  $(t', t'')$  appearing  $\geq \log(d)$  times consecutively in  $\text{text}_i$  **then**  
    Append the rule  $t \leftarrow (t', t'')$  to DS  
    **for**  $j = i + 1, \dots, d$  **do**  
       $\text{text}_j \leftarrow \text{APPLY}_{t \leftarrow (t', t'')}( \text{text}_j )$ ;  
  *// Can be implemented in parallel*  
  **end for**  
  **end if**  
**end for**  
  
*// Encoder and Decoder*  
**if**  $|\text{Dict}| < d_0 \triangleq \epsilon d / 2 \log(4|\mathcal{A}|)$  **then**  
   $\mathcal{T} \leftarrow (\text{Dict}, \text{DS}, \text{enc}_{\text{BPE.split}}(\cdot), \text{dec}_{\text{BPE.split}}(\cdot))$   
**else**  
   $\mathcal{T} \leftarrow (\text{Dict}, \emptyset, \text{enc}_{\text{gre}}(\cdot), \text{dec}_{\text{gre}}(\cdot))$   
**end if**  
  
**def**  $\text{APPLY}_{t \leftarrow (t_1, t_2)}(\text{text})$ :  
  Replace every consecutive occurrence of  $(t', t'')$  in  $\text{text}$  by  $t$ , breaking ties arbitrarily.

---

**Theorem B.2.** For any  $\epsilon \in (0, 1)$ , run Algorithm 1 on a dataset of  $n = \Theta(d^2)$  characters to learn a dictionary with at most  $d$  tokens. The resulting tokenizer  $\mathcal{T}$  satisfies with probability  $\geq 1 - e^{-\Omega(d\epsilon^2)}$ ,

$$\min_{Q \in \mathcal{Q}_{1\text{-gram}}} \mathcal{L}(Q \circ \text{enc}(\cdot)) \leq (2 + \epsilon) \min_Q \mathcal{L}(Q) + \epsilon$$

where  $\epsilon = O\left(\frac{\log(1/\epsilon) \log^3(1/\delta)}{\epsilon \delta^9 \log(d)}\right)$ .

While the guarantees established for the sequential BPE tokenizer are weaker than those in Theorem 3.1, the analysis turns out to be quite involved. Theorem B.2 implies that unigram models trained on the sequential BPE tokenizer asymptotically approach the optimal cross-entropy loss up to a factor of 2.

The formal proof of this result is presented in Appendix B. What is the intuition behind using a different encoder in Algorithm 1 depending on the number of tokens in the dictionary? When the number of tokens in the dictionary is smaller than  $d_0$ , we know that on a  $1 - d_0/d$  fraction of the iterations of Algorithm 1, a token is *not* added to the dictionary, i.e. every pair of tokens already appears at most  $\log(d)$  times together. This is a datapoint of “evidence” that under the dictionary in that iteration, the BPE encoder is already good at encoding new strings (of length  $\Theta(d)$ ) in a way where pairs of tokens do not appear consecutively with high frequency. Since future dictionaries only have more rules appended to them, dictionaries only get better at encoding new strings into tokens where pairs do not frequently appear consecutively. In other words, the BPE encoder satisfies a monotonicity property. It remains to show that dictionaries which encode new sequences in a way where no pair of tokens appear too frequently have large  $H(Q_{\text{MLE}}, P)$  (to invoke Theorem 3.4). This follows from ideas introduced in (Navarro and Russo, 2008).

The case where the number of tokens is large ( $\geq d_0$ ) turns out to present significant technical challenges for analyzing the BPE encoder. There is no longer much “evidence” that the dictionary in each iteration is good at encoding strings since in a large number of iterations a pair of tokens appear consecutively with high frequency. Analyzing the greedy encoder also presents its own challenges - although the algorithm has allocated a large number of tokens, it is possible that there are short tokens  $t$  which are maximal (i.e. they are not prefixes of other tokens). This is similar to the problem encountered by BPE when trained on the dataset in eq. (15) - although the algorithm has allocated a large number of tokens, the token 1 is maximal since every other token begins with the character 0.

However, it turns out that such tokens, although present in the dictionary, are not observed frequently while encoding a fresh string drawn from the source.

### B.1 Analysis of Algorithm 1

In this section we prove a rephrased version of Theorem B.2 which implies the statement in the main paper. Define  $d_0 = \frac{\epsilon d}{2 \log(4|\mathcal{A}|)}$ .

**Theorem B.3** (Rephrased Theorem B.2). *Run Algorithm 1 on a dataset of  $n = \Theta(d^2)$  characters to learn a dictionary with at most  $d$  tokens. The resulting tokenizer  $\mathcal{T}$  satisfies one of the following 3 conditions,*

1. *Either,  $|\text{Dict}| > d_0$ , and,*

$$\min_{Q \in \mathcal{Q}_{1\text{-gram}}} \mathcal{L}(Q \circ \text{enc}(\cdot)) \leq \frac{H_\infty}{1 - \varepsilon}.$$

$$\text{Here, } \varepsilon = O\left(\frac{\log^3(1/\delta) \log(1/\epsilon)}{\epsilon \delta^9 \log(d)}\right).$$

2.  $\Pr(|\text{Dict}| < d_0) = e^{-\Omega(\epsilon^2 d / \log^2(1/\delta))}$ , or,

3. *Conditional on  $|\text{Dict}| < d_0$ , with probability  $\geq 1 - e^{-\Omega(\epsilon^2 d / \log^2(1/\delta))}$ ,*

$$\min_{Q \in \mathcal{Q}_{1\text{-gram}}} \mathcal{L}(Q \circ \text{enc}(\cdot)) \leq \left(1 - \frac{2d_0}{d}\right) \left(2H_\infty + O\left(\frac{1}{\log(d)}\right)\right) + \frac{2d_0}{d} \log(4|\mathcal{A}|).$$

*With the choice of  $d_0 = \epsilon d / 2 \log(4|\mathcal{A}|)$  we get the statement of Theorem B.2.*

### B.2 Analysis for the large dictionary case: $|\text{Dict}| > d_0$

In the large dictionary case, Algorithm 1 uses the greedy encoder/decoder pair in conjunction with the dictionary. The proof of Theorem B.2 relies on establishing that the cross-entropy  $H(Q_{\text{MLE}}, P)$  of the tokenizer is large. Namely, we prove that,

**Lemma B.4.** *In Algorithm 1, assuming at least  $d_0$  tokens are allocated,*

$$H(Q_{\text{MLE}}, P) = \Omega\left(\frac{\epsilon \delta^9 \log(d)}{\log(1/\epsilon) \log^3(1/\delta)}\right).$$

To show this, it suffices to argue that conditioned on any previous set of tokens, with nontrivial probability over the underlying string generated from the stochastic source, the next token is long (i.e. having conditional probability at most  $O(1/\sqrt{d})$ ).

**Lemma B.5.** *Suppose that in a run of Algorithm 1, at least  $d_0$  tokens are allocated. Suppose a set of tokens  $t_1, \dots, t_k$  have been sampled so far by the greedy encoder. Let  $T_{i+1}$  be the random variable which denotes the next token returned by the greedy encoder, where the randomness comes from the underlying string being tokenized. Then,*

$$\Pr\left(P(T_{i+1}|t_i) \leq 1/\sqrt{C\delta d} \mid t_1, \dots, t_i\right) \geq \frac{d_0 \delta^6 (1 - \delta)^2}{8Cd\Delta|\mathcal{A}| \log(2|\mathcal{A}|)n_D} = \Omega\left(\frac{\epsilon \delta^9}{\log^3(1/\delta) \log(1/\epsilon)}\right)$$

*Proof sketch of Lemma B.5.* The proof will proceed in 2 parts. We first show in Lemma B.9 that there is a set  $D_{\text{valid}}$  of  $\Omega(d)$  tokens in the dictionary which are neither prefixes nor suffixes of any other token in  $\text{Dict}$ . The reason for considering this set of tokens is twofold,

1. Irrespective of what the previous set of tokens were, it is legal for a token  $D_{\text{valid}}$  to be sampled in the current step by the greedy encoder, since for any candidate  $t \in D_{\text{valid}}$ , by definition,  $t_j \cdots t_i t \notin \text{Dict}$  for every  $j \leq i$ .



and concatenate the corresponding substrings to get an  $m = \sum_{i=1}^{i^*} |t_i|$  length character string. Letting  $i^* \rightarrow \infty$ , we must have  $m \rightarrow \infty$  surely since  $m \geq i^*$ . In this view, eq. (16) can be rewritten as,

$$Q_{\text{MLE}}(\mathbf{t}) = \lim_{i^* \rightarrow \infty} \frac{n_{\mathbf{t}}}{i^*} = \lim_{i^* \rightarrow \infty} \frac{1}{i^*} \sum_{i=1}^{i^*} \mathbb{I}(t_i = \mathbf{t}) \stackrel{\text{a.s.}}{=} \lim_{i^* \rightarrow \infty} \frac{1}{i^*} \sum_{i=1}^{i^*} \mathbb{E}[\mathbb{I}(t_i = \mathbf{t}) | t_1, \dots, t_{i-1}] \quad (17)$$

where the last inequality follows by the sequential nature of the token sampling process and a martingale argument. Consider the set of tokens  $T$  such that  $\mathbf{t} \in T$  satisfies  $\max_{a \in \mathcal{A}} P(\mathbf{t}|a) \leq \sqrt{1/C\delta^3 d}$ . From eq. (17), summing across  $\mathbf{t} \in T$ , we have that,

$$\sum_{\mathbf{t} \in T} Q_{\text{MLE}}(\mathbf{t}) \stackrel{\text{a.s.}}{=} \lim_{i^* \rightarrow \infty} \frac{1}{i^*} \sum_{i=1}^{i^*} \Pr(t_i \in T | t_1, \dots, t_{i-1}) = \Omega\left(\frac{\epsilon \delta^9}{\log^3(1/\delta) \log(1/\epsilon)}\right) \quad (18)$$

where in the last inequality, we use Lemma B.5 and the fact that  $\delta \max_{a \in \mathcal{A}} P(\mathbf{t}|a) \geq \min_{a \in \mathcal{A}} P(\mathbf{t}|a)$ . Therefore,

$$H(Q_{\text{MLE}}, P) \geq \sum_{\mathbf{t} \in T} Q_{\text{MLE}}(\mathbf{t}) \log(1/P(\mathbf{t})) \geq \sum_{\mathbf{t} \in T} Q_{\text{MLE}}(\mathbf{t}) \log(\sqrt{C\delta^3 d})$$

where in (i) we use the fact that for  $\mathbf{t} \in T$ ,  $\max_{a \in \mathcal{A}} P(\mathbf{t}|a) \leq 1/\sqrt{C\delta^3 d}$ , which implies that  $P(\mathbf{t}) \leq 1/\sqrt{C\delta^3 d}$ . Finally, combining with eq. (18) completes the proof of Lemma B.4. Furthermore, since the cross-entropy  $H(Q_{\text{MLE}}, P)$  was established to be large, by invoking the reduction in Theorem 3.4, we complete the proof of Theorem B.3.1.

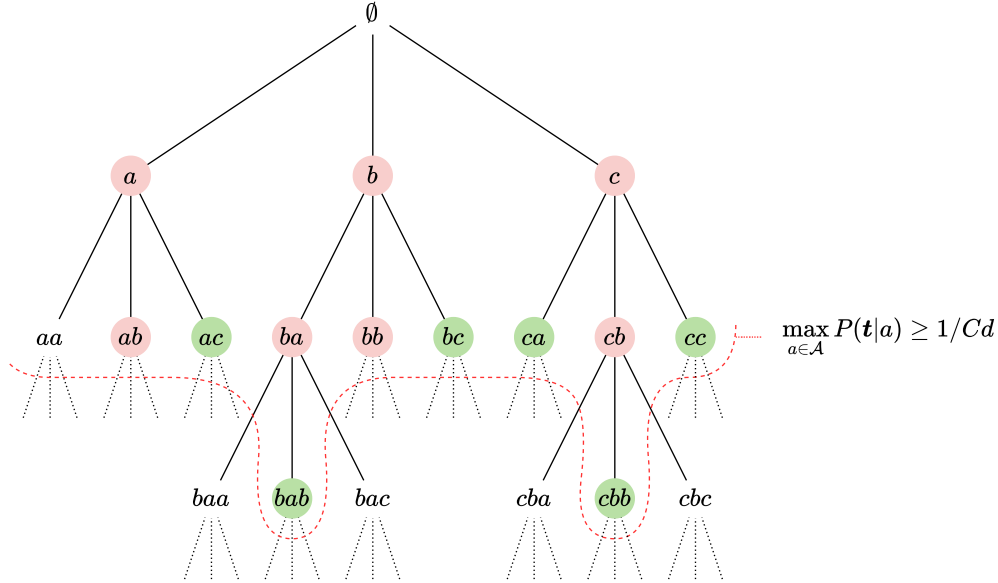


Figure 8: The circled nodes indicate substrings which are tokens in Dict. The red boundary is the set of substrings  $\mathbf{t}$  such that  $\max_{a \in \mathcal{A}} P(\mathbf{t}|a) \geq 1/Cd$ . By Lemma B.8, none of the nodes which fall outside this boundary are assigned as tokens in a run of Algorithm 1. The set of circled substrings are the set of tokens in Dict. Among them, the ones circled green are the tokens in  $D_{\text{valid}}$ , which are not prefixes or suffixes of any other tokens in Dict. Substrings such as  $cb$  or  $ba$  which are tokens in Dict do not belong to  $D_{\text{valid}}$  because they are prefixes of longer tokens (in this case,  $cbb$  and  $bab$  respectively). On the other hand, substrings like  $ab$  do not belong to  $D_{\text{valid}}$  since they are suffixes of tokens in Dict, in this case,  $bab$ . Lemma B.9 asserts that the number of tokens in  $D_{\text{valid}}$  are  $\Omega(d)$  in number, assuming that Dict has  $\Omega(d)$  tokens to begin with.

**Notation.** For each  $a \in \mathcal{A}$  and  $j \in \mathbb{N} \cup \{0\}$ , define a *level set* of substrings,

$$S_j^a = \left\{ (1 - \delta)^{j+1} < P(\mathbf{t} | t_1 = a) \leq (1 - \delta)^j \right\}$$

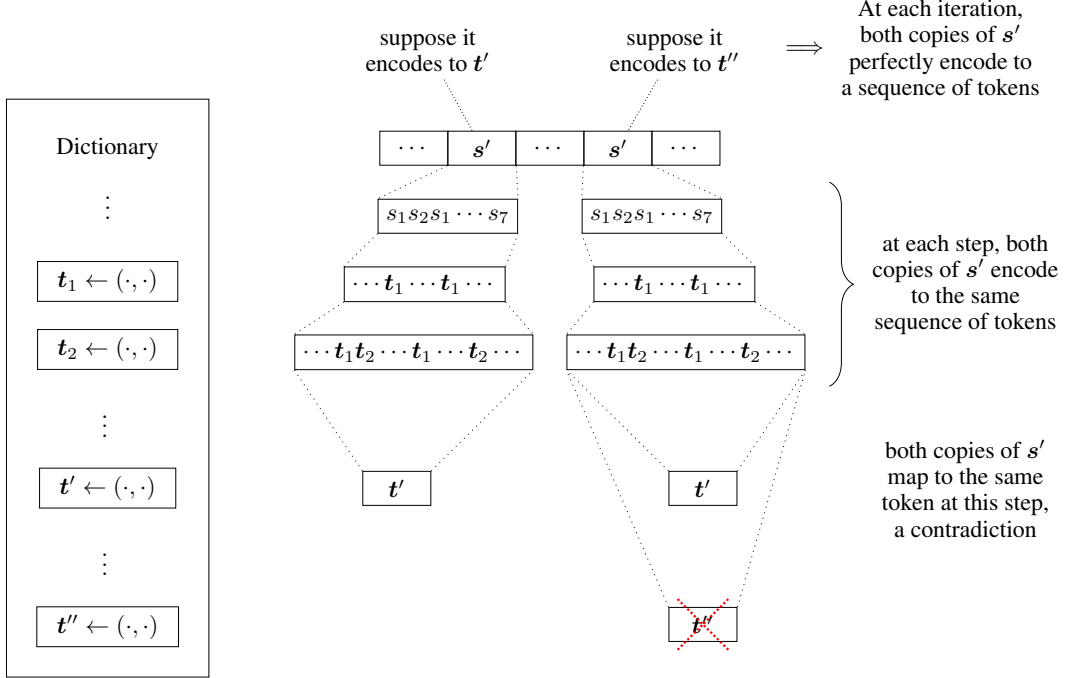


Figure 9: A pictorial representation of the proof of Lemma B.6

where  $t_1$  denotes the first character of  $t$ . And likewise, define the sets  $S_j = \cup_{a \in \mathcal{A}} S_j^a$ ,  $S_{\leq j}^a$  and  $S_{\geq j}^a$  as the union of  $S_j^a$  over  $j' \geq j$ ,  $j' \leq j$  and  $S_{\leq j}$  and  $S_{\geq j}$  as the union of  $S_{\leq j}^a$  and  $S_{\geq j}^a$  over  $a \in \mathcal{A}$ . Furthermore for a large universal constant  $C > 0$ , define parameters,

$$\Delta = \frac{\log(\delta)}{\log(1-\delta)} \asymp \Theta\left(\frac{\log(1/\delta)}{\delta}\right); \quad n_D = 1 - \frac{2 \log(4Cd/\delta d_0)}{\log(1-\delta)} \asymp \Theta\left(\frac{\log(1/\epsilon\delta)}{\delta}\right). \quad (19)$$

We first begin by stating a folklore result: every pair of tokens assigned by a merging-based dictionary generation algorithm have distinct character representations.

**Lemma B.6.** *If Algorithm 1 assigns a new token in some round, it's character representation must be distinct from that of all previously assigned tokens.*

*Proof.* A pictorial proof is in Figure 9. We will prove this result by contradiction. Suppose  $t$  and  $t'$  are tokens which decode to the same character substring,  $s'$ . Consider all occurrences of  $s'$  in the dataset which in some iteration encode into  $t'$  or  $t''$ , and denote these disjoint locations  $\mathcal{S}$ . Recall that at these locations,  $s'$  eventually is assumed to map to a singular token  $t'$  or  $t''$ . Therefore, at every step in the merging process these occurrences of  $s'$  must perfectly map to a sequence of tokens.

Now consider the merging process at the first time before any of the rules corresponding to tokens in  $t'$  or  $t''$  are implemented. Prior to this time, all the occurrences of  $s'$  corresponding to the locations in  $\mathcal{S}$  have not been tokenized yet. When the first rule corresponding to one of the tokens in  $\{t', t''\}$  is implemented, all the strings in  $\mathcal{S}$  must be modified identically. This uses the fact that we can isolate each of these occurrences of  $s'$  while carrying out the merging process, since each location must be distinct. At every step, the encodings of these copies of  $s'$  must be the same, and therefore  $t'$  and  $t''$  cannot be two distinct tokens.  $\square$

**Lemma B.7.** *The size of the level set  $S_j^a$  is bounded by  $(1-\delta)^{-(j+1)}$ .*

*Proof.* Since the probability of any transition is at most  $1-\delta$ , this implies that any infinite trajectory on the tree  $\mathcal{T}_a^*$  can intersect at most one vertex in  $S_j^a$ . Therefore,  $\sum_{t \in S_j^a} P(t|t_1 = a) \leq 1$ . By the lower bound on  $P(t|t_1 = a)$  for  $t \in S_j^a$ , this implies the statement of the lemma.  $\square$

Next we show that with high probability none of the substrings  $\mathbf{t}$  having probability mass (under  $P$ ) of at most  $\delta/Cd$  conditioned on the first character, are assigned as tokens by Algorithm 1.

**Lemma B.8.** *In a run of Algorithm 1, for a sufficiently large constant  $C > 0$ , with probability  $d^{-\Omega(1)} \text{poly}(1/\delta)$  all assigned tokens  $\mathbf{t} \in \text{Dict}$  satisfy  $\max_{a \in \mathcal{A}} P(\mathbf{t}|a) \geq 1/Cd$ . In other words, none of the substrings in  $S_{\geq j^*}$  are added as tokens to the dictionary in a run of Algorithm 1, where,*

$$j^* \triangleq \log(\delta/Cd)/\log(1-\delta)$$

*Proof.* Consider some  $j \geq j^*$  and  $a \in \mathcal{A}$  and substring  $\mathbf{t} \in S_j^a$ . In the  $i^{\text{th}}$  stage of the algorithm where  $\text{text}_i$  is being processed, for  $\mathbf{t}$  to be assigned as a token, at the very least,  $\mathbf{t}$  must appear at least  $\log(d)$  times disjointly in  $\text{text}_i$ . Therefore,

$$\begin{aligned} P(\mathbf{t} \in S_j^a \text{ is assigned as a token in } \text{text}_i) &\leq \binom{d}{\log(d)} \left( \max_{a \in \mathcal{A}} P(\mathbf{t}|a) \right)^{\log(d)} \\ &\leq d^{\log(d)} \left( \frac{1}{Cd} (1-\delta)^{j-j^*} \right)^{\log(d)} \\ &\leq d^{-\log(C)} (1-\delta)^{(j-j^*) \log(d)} \end{aligned}$$

Union bounding over  $S_j^a$  over  $j \geq j^*$  using the bound on  $|S_j^a|$  in Lemma B.7, and over  $a \in \mathcal{A}$  and  $i \in [d]$  results in the bound,

$$P(\mathbf{t} \in S_{\geq j^*} \text{ is assigned as a token in step } i \text{ for some } i \in [d]) \leq d^{-\Omega(1)} \sum_{j \geq j^*} \frac{(1-\delta)^{(j-j^*) \log(d)}}{(1-\delta)^{j+1}} \leq \frac{d^{-\Omega(1)}}{\delta(1-\delta)}$$

□

**Lemma B.9.** *Consider the set of tokens  $D_{\text{valid}}$  which are not a prefix or a suffix of any other token in  $\text{Dict}$ . That is,  $D_{\text{valid}} = \{\mathbf{t} \in \text{Dict} : \nexists s : \mathbf{t}s \in \text{Dict}\} \cap \{\mathbf{t} \in \text{Dict} : \nexists s : s\mathbf{t} \in \text{Dict}\}$ . If  $|\text{Dict}| \geq d_0$ , then,*

$$|D_{\text{valid}}| \geq \frac{d_0}{4n_D}.$$

where  $n_D$  is defined in eq. (19).

*Proof.* For any token  $\mathbf{t} \in D_{\text{valid}}$ , there may be at most  $2|\mathbf{t}|$  tokens which are suffixes or prefixes of it and belong to  $\text{Dict}$ . More importantly, every token in  $\text{Dict}$  not belonging to  $D_{\text{valid}}$  must either be a prefix or a suffix of some token in  $D_{\text{valid}}$ . Split the suffixes and prefixes of the tokens in  $D_{\text{valid}}$  into four sets,

1.  $S_{\text{suff}, \min} = \bigcup_{\mathbf{t} \in D_{\text{valid}}} \{\mathbf{t}' \in \text{Dict} : \mathbf{t}' \in \text{suff}(\mathbf{t}), |\mathbf{t}'| \leq |\mathbf{t}| - n_D\},$
2.  $S_{\text{suff}, \max} = \bigcup_{\mathbf{t} \in D_{\text{valid}}} \{\mathbf{t}' \in \text{Dict} : \mathbf{t}' \in \text{suff}(\mathbf{t}), |\mathbf{t}'| > |\mathbf{t}| - n_D\},$
3.  $S_{\text{pre}, \min} = \bigcup_{\mathbf{t} \in D_{\text{valid}}} \{\mathbf{t}' \in \text{Dict} : \mathbf{t}' \in \text{pre}(\mathbf{t}), |\mathbf{t}'| \leq |\mathbf{t}| - n_D\},$
4.  $S_{\text{pre}, \max} = \bigcup_{\mathbf{t} \in D_{\text{valid}}} \{\mathbf{t}' \in \text{Dict} : \mathbf{t}' \in \text{pre}(\mathbf{t}), |\mathbf{t}'| > |\mathbf{t}| - n_D\}.$

where  $n_D$  is defined in eq. (19). Note from Lemma B.8 that all the tokens  $\mathbf{t} \in \text{Dict}$  all satisfy  $\max_{a \in \mathcal{A}} P(\mathbf{t}|a) \geq 1/Cd$ . Therefore, the tokens in  $S_{\text{pre}, \min}$  and  $S_{\text{suff}, \min}$  all satisfy,  $\max_{a \in \mathcal{A}} P(\mathbf{t}|a) \geq d/C(1-\delta)^{n_D}$ . By summing Lemma B.7 over appropriate  $j$ , we get that  $|S_{\text{pre}, \min}| + |S_{\text{suff}, \min}| \leq 2Cd(1-\delta)^{n_D-1}/\delta$ .

On the other hand, corresponding to any  $\mathbf{t} \in D_{\text{valid}}$ , there are at most  $n_D$  tokens in  $S_{\text{pre}, \max}$  or  $S_{\text{suff}, \max}$  and therefore  $|S_{\text{pre}, \max}|, |S_{\text{suff}, \max}| \leq n_D \cdot |D_{\text{valid}}|$ . Since every token in  $\text{Dict}$  either belongs to  $D_{\text{valid}}$  or is a suffix of some token in  $D_{\text{valid}}$ ,  $S_{\text{pre}, \min} \cup S_{\text{pre}, \max} \cup S_{\text{suff}, \min} \cup S_{\text{suff}, \max} = |\text{Dict}|$  and,

$$2n_D \cdot |D_{\text{valid}}| + \frac{2C(1-\delta)^{n_D-1}d}{\delta} \geq d_0$$



Recalling the choice of  $n_D = 1 - \frac{2 \log(4Cd/\delta d_0)}{\log(1-\delta)}$ , we get that,

$$|D_{\text{valid}}| \geq \frac{d_0}{4n_D}.$$

□

**Lemma B.10.** *Suppose Algorithm 1 assigns at least  $d_0$  tokens. For any character  $a \in \mathcal{A}$ , sample an  $a' \sim P(\cdot|a)$  and an infinite trajectory on the tree  $\mathcal{T}_{a'}^*$ , denoted  $\text{traj}$ . Then,*

$$\mathbb{E}_{a' \sim P(\cdot|a)} \left[ \Pr_{\text{traj} \sim \mathcal{T}_{a'}^*} \left( \min_{t \in \text{traj} \cap D_{\text{valid}}} P(t|a) \leq \sqrt{\delta/Cd} \middle| a' \right) \right] \geq \frac{d_0 \delta^6 (1-\delta)^2}{8Cd\Delta|\mathcal{A}|n_D}.$$

where the notation  $\mathcal{T}_{a'}^*$  is used to overload the distribution over infinite trajectories on  $\mathcal{T}_{a'}^*$ . The parameters  $n_D$  and  $\Delta$  are defined in eq. (19).

*Proof.* By Lemma B.8, recall that the  $\geq d_0$  tokens assigned in a run of Algorithm 1, with high probability, are substrings in  $S_{\leq j^*}$ . For any  $a \in \mathcal{A}$ , the total number of substrings in  $S_{\leq j^*}$  can be bounded as,

$$|S_{\leq j^*}| \leq \sum_{a \in \mathcal{A}} \sum_{j=0}^{j^*} |S_j^a| \leq \sum_{a \in \mathcal{A}} \sum_{j=0}^{j^*} \frac{1}{(1-\delta)^{j+1}} \leq \frac{C|\mathcal{A}|d}{\delta(1-\delta)}. \quad (20)$$

In order to prove this result, we use a counting argument and the fact that no tokens in  $S_{> j^*}$  are assigned. Consider some character  $a$  and all the leaves in the forest  $S_{\leq j^*}$ . Since every transition has  $\geq \delta$  probability of occurring, across all leaf nodes  $t \in S_{\leq j^*}$ ,  $P(t|a')$  are within a  $\delta^2(1-\delta)$  factor of each other across different  $a' \in \mathcal{A}$ . In particular, by counting the number of paths in  $\mathcal{T}_a^*$  from  $\emptyset$  to leaf nodes in  $S_{\leq j^*}^a$ , across  $a \in \mathcal{A}$  along which a token in Dict exists in  $S_{\geq j^*/2}$ , we can also compute the probability mass across such trajectories up to a factor of  $\delta^2(1-\delta)$ .

Taking the union across  $a \in \mathcal{A}$ , consider the paths in  $\mathcal{T}_a^*$  from  $\emptyset$  to leaf nodes in  $S_{\leq j^*}^a$ . From Lemma B.9,  $\sum_{j \leq j^*} |D_{\text{valid}} \cap S_j| \geq d_0/4n_D$ , where  $n_D = 1 - 2 \log(4Cd/\delta d_0)/\log(1-\delta)$ . Note that for sufficiently large  $d = \Omega(\log(1/\epsilon\delta)/\delta^5)$ , by Lemma B.7,  $\sum_{j \leq j^*/2} |S_j| = \sqrt{Cd/\delta}/\delta(1-\delta) \leq d_0/8n_D$ . Therefore,

$$\sum_{j^*/2 < j \leq j^*} |D_{\text{valid}} \cap S_j| \geq \frac{d_0}{8n_D}. \quad (21)$$

Define  $\Delta = \log(\delta)/\log(1-\delta)$ . Combining eq. (21) with eq. (20) and applying the probabilistic method, there exists an  $i^* \geq j^*/2$  such that,

$$\frac{|D_{\text{valid}} \cap (S_{i^*+1} \cup \dots \cup S_{i^*+\Delta})|}{|S_{i^*+1} \cup \dots \cup S_{i^*+\Delta}|} \geq \frac{\delta(1-\delta)d_0}{8Cd|\mathcal{A}|n_D}. \quad (22)$$

Note that  $\Delta$  is chosen to be sufficiently large, so that every infinite trajectory on  $\mathcal{T}_{a'}^*$  must intersect at least once with the band of vertices  $S_{i^*+\Delta+1}^{a'} \cup \dots \cup S_{i^*+2\Delta}^{a'}$ . Note that this band is different from the one considered in eq. (22). Define  $L_{a'}$  as the set of longest prefixes across infinite trajectories in  $\mathcal{T}_{a'}^*$  which belong to  $S_{i^*+\Delta+1}^{a'} \cup \dots \cup S_{i^*+2\Delta}^{a'}$ .

Note that our objective is to show that an infinite trajectory sampled on  $\mathcal{T}_{a'}^*$  where  $a' \sim P(\cdot|a)$ , has a long prefix in Dict. We can truncate this trajectory to lower bound this probability, and therefore, we assume that the infinite trajectories on  $\mathcal{T}_{a'}^*$  terminate once they reach a substring in  $L_{a'}$ . Furthermore, note that although  $\Delta$  is large, it is still a constant depending on  $\delta$ . Therefore, the band of states  $S_{i^*+\Delta+1}^{a'} \cup \dots \cup S_{i^*+2\Delta}^{a'}$  is not too wide, and all the substrings in  $L_{a'}$  have approximately similar probabilities to each other. In particular, for any character  $a \in \mathcal{A}$ , and for any  $a' \in \mathcal{A}$  and  $t \in L_{a'}$ , decomposing  $P(t|a)$  as  $P(t|t_1 = a')P(a'|a)$ ,

$$\delta^2(1-\delta) \cdot (1-\delta)^{i+\Delta} \stackrel{(i)}{\leq} P(t|a) \stackrel{(ii)}{\leq} (1-\delta)^{i+\Delta}. \quad (23)$$

Inequality (i) follows from the fact that all transition probabilities are at least  $\delta$ , so every leaf node in  $L_{a'}$  must have  $P(t|t_1 = a') \geq (1-\delta)^{i+2\Delta+1}$ , and the fact that  $P(a'|a) \geq \delta$ . Inequality (ii)

follows similarly from the fact that  $\mathbf{t}$  is a leaf node of  $L_{a'}$  and therefore  $P(\mathbf{t}|\mathbf{t}_1 = a') \leq (1 - \delta)^{i^* + \Delta}$ . Therefore, instead of bounding the probability of any event under the distribution over substrings in  $L_{a'}$  induced by truncating the infinite strings sampled on  $\mathcal{T}_{a'}^*$ , it suffices to count the fraction of substrings in  $L_{a'}$  satisfying the event (which are equivalent up to a  $\delta(1 - \delta)$  factor). Define,

$$\text{pre}(\mathbf{t}) = (\mathbf{t}_1, \mathbf{t}_{1:2}, \mathbf{t}_{1:3}, \dots, \mathbf{t}_{1:|\mathbf{t}|})$$

As the set of prefixes of  $\mathbf{t}$  (including  $\mathbf{t}$ ). Note that at most  $\Delta$  of the prefixes of any substring  $\mathbf{t}$  can intersect with  $S_{i^*+1}^a \cup \dots \cup S_{i^*+\Delta}^a$ . Therefore,

$$\begin{aligned} & \sum_{a' \in \mathcal{A}} \sum_{\mathbf{t} \in L_{a'}} \mathbf{1}(\text{pre}(\mathbf{t}) \cap D_{\text{valid}} \cap (S_{i^*+1}^{a'} \cup \dots \cup S_{i^*+\Delta}^{a'}) \neq \emptyset) \\ & \geq \sum_{a' \in \mathcal{A}} \sum_{\mathbf{t} \in L_{a'}} \frac{|\text{pre}(\mathbf{t}) \cap D_{\text{valid}} \cap (S_{i^*+1}^{a'} \cup \dots \cup S_{i^*+\Delta}^{a'})|}{\Delta} \\ & \stackrel{(i)}{\geq} \sum_{a' \in \mathcal{A}} \frac{|D_{\text{valid}} \cap (S_{i^*+1}^{a'} \cup \dots \cup S_{i^*+\Delta}^{a'})|}{\Delta} \\ & \stackrel{(ii)}{\geq} \frac{\delta d_0(1 - \delta)}{8Cd\Delta|\mathcal{A}|n_D} \sum_{a' \in \mathcal{A}} |S_{i^*+1}^{a'} \cup \dots \cup S_{i^*+\Delta}^{a'}| \\ & \stackrel{(iii)}{\geq} \frac{\delta^3 d_0(1 - \delta)}{8Cd\Delta|\mathcal{A}|n_D} \sum_{a' \in \mathcal{A}} |L_{a'}|, \end{aligned}$$

where (i) uses the fact that the prefixes of  $\mathbf{t} \in L_{a'}$  cover all the substrings in  $S_{\leq i^*+\Delta}^{a'}$ , and therefore  $\bigcup_{\mathbf{t} \in L_{a'}} \text{pre}(\mathbf{t}) \supset S_{i^*+1}^{a'} \cup \dots \cup S_{i^*+\Delta}^{a'}$ , and (ii) uses eq. (22). Finally, (iii) uses the fact that  $\Delta$  is not too large, and therefore, for any substring  $\mathbf{t}' \in S_{i^*+1}^{a'} \cup \dots \cup S_{i^*+\Delta}^{a'}$ , there are at most  $1/(1 - \delta)^{2\Delta} = 1/\delta^2$  substrings  $\mathbf{t} \in L_{a'}$  which contain it as a prefix. This means,  $|L_{a'}| \leq |S_{i^*+1}^{a'} \cup \dots \cup S_{i^*+\Delta}^{a'}|/\delta^2$ . After dividing by  $\sum_{a' \in \mathcal{A}} |L_{a'}|$  on both sides, this implies,

$$\mathbb{E}_{a' \sim \text{Unif}(\mathcal{A})} \left[ \Pr_{\mathbf{t} \sim \text{Unif}(L_{a'})} \left( \text{pre}(\mathbf{t}) \cap D_{\text{valid}} \cap (S_{i^*+1}^{a'} \cup \dots \cup S_{i^*+\Delta}^{a'}) \neq \emptyset \mid a' \right) \right] \geq \frac{\delta^3 d_0(1 - \delta)}{8Cd\Delta|\mathcal{A}|n_D}. \quad (24)$$

The event inside the inner probability term is the event that an infinitely long string (truncated at  $L_{a'}$ ) has a prefix which lies in  $D_{\text{valid}}$  and which intersects with  $S_{i^*+1}^{a'} \cup \dots \cup S_{i^*+\Delta}^{a'}$ , which implies that it has probability  $P(\mathbf{t}|a) \leq \sqrt{\delta/Cd}$ . Therefore, we have that for any  $a \in \mathcal{A}$ , sampling an  $a' \sim P(\cdot|a)$  and an infinite trajectory  $\text{traj} \sim \mathcal{T}_{a'}^*$ ,

$$\begin{aligned} & \mathbb{E}_{a' \sim P(\cdot|a)} \left[ \Pr_{\text{traj} \sim \mathcal{T}_{a'}^*} \left( \min_{\mathbf{t} \in \text{traj} \cap D_{\text{valid}}} P(\mathbf{t}|a) \leq \sqrt{\delta/Cd} \mid a' \right) \right] \\ & \stackrel{(i)}{\geq} \delta^2(1 - \delta) \cdot \mathbb{E}_{a' \sim P(\cdot|a)} \left[ \Pr_{\mathbf{t}' \sim \text{Unif}(L_{a'})} \left( \min_{\mathbf{t} \in \text{pre}(\mathbf{t}') \cap D_{\text{valid}}} P(\mathbf{t}|a) \leq \sqrt{\delta/Cd} \mid a' \right) \right] \\ & \stackrel{(ii)}{\geq} \delta^2(1 - \delta) \cdot \mathbb{E}_{a' \sim P(\cdot|a)} \left[ \Pr_{\mathbf{t}' \sim \text{Unif}(L_{a'})} \left( \text{pre}(\mathbf{t}') \cap D_{\text{valid}} \cap (S_{i^*+1}^{a'} \cup \dots \cup S_{i^*+\Delta}^{a'}) \neq \emptyset \mid a' \right) \right] \\ & \stackrel{(iii)}{\geq} \delta^3(1 - \delta) \cdot \mathbb{E}_{a' \sim \text{Unif}(\mathcal{A})} \left[ \Pr_{\mathbf{t}' \sim \text{Unif}(L_{a'})} \left( \text{pre}(\mathbf{t}') \cap D_{\text{valid}} \cap (S_{i^*+1}^a \cup \dots \cup S_{i^*+\Delta}^a) \neq \emptyset \mid a' \right) \right] \\ & \geq \delta^3(1 - \delta) \cdot \frac{\delta^3 d_0(1 - \delta)}{8Cd\Delta|\mathcal{A}|n_D}. \end{aligned}$$

Here (i) follows by truncating the trajectory  $\text{traj}$  to terminate at a node in  $\bigcup_{a' \in \mathcal{A}} L_{a'}$  and from eq. (23), (ii) follows by arguing that  $i^* \leq j^*/2$  and therefore if a prefix of  $\mathbf{t}'$  lies in  $S_{i^*+1}^{a'} \cup \dots \cup S_{i^*+\Delta}^{a'}$ , then it must have  $P(\mathbf{t}|a) \leq \sqrt{\delta/Cd}$ . Inequality (iii) follows by noting that all the transitions  $P(a'|a)$  have probability  $\geq \delta$ , and the last inequality follows from eq. (24).  $\square$

### Proof of Lemma B.5

Lemma B.10 concludes that given any previous sequence of tokens terminating in a character  $a$ , with constant probability, an infinite trajectory sampled from  $\mathcal{T}_a^*$  with  $a' \sim P(\cdot|a)$  has as prefix, a substring  $t$ , which not only has low probability, with  $P(t|a) \leq \sqrt{\delta/Cd}$ , but also belongs to the subset of tokens  $D_{\text{valid}}$ . Note that regardless of the previously sampled tokens, it is legal to sample any token in  $D_{\text{valid}}$  as the current token, since by definition, these tokens are not the suffixes of any other tokens in Dict. Moreover, if any trajectory on  $\mathcal{T}_a^*$  reaches a token in  $D_{\text{valid}}$ , then it must be largest token along that trajectory, since none of the tokens in  $D_{\text{valid}}$  are prefixes of another token in Dict.

Consider generating a new token by rejection sampling. Suppose the set of previous tokens  $t_1, \dots, t_i$  end in some character  $a$ . Sample the next character  $a' \sim P(\cdot|a)$  and an infinite trajectory on  $\mathcal{T}_a^*$ . If it reaches an illegal token  $t$  such that  $t_j t_{j+1} \dots t_i t$  already exists in Dict, this token is rejected and the trajectory is resampled. By the prefix-free property of these tokens, if this trajectory visits a token in  $D_{\text{valid}}$ , it must immediately be output as the next token. Note that this probability is lower bounded by,

$$\mathbb{E}_{a' \sim P(\cdot|a)} \left[ \Pr_{\text{traj} \sim \mathcal{T}_a^*} \left( \min_{t \in \text{traj} \cap D_{\text{valid}}} P(t|a) \leq \sqrt{\delta/Cd} \mid a' \right) \right]$$

which is lower bounded by  $\text{poly}(\epsilon, \delta)$ , the subject of Lemma B.10. Therefore with this probability, the process terminates in the first step with a token in  $D_{\text{valid}}$  being sampled.

### B.3 Analysis in the small dictionary case

In this section, we will prove Theorem B.3.2 and Theorem B.3.3. In particular we show that, either,

1. The dictionary is small with low probability. i.e.,  $\Pr(|\text{Dict}| < d_0) = e^{-\Omega(\epsilon^2 d / \log^2(1/\delta))}$ , or,
2. Or conditioned on the dictionary being small,  $|\text{Dict}| < d_0$ , with high probability  $\geq 1 - e^{-\Omega(\epsilon^2 d / \log^2(1/\delta))}$ ,

$$\min_{Q \in \mathcal{Q}_{1\text{-gram}}} \mathcal{L}(Q \circ \text{enc}(\cdot)) \leq 4 \left( 1 - \frac{2d_0}{d} + O\left(\frac{1}{\log(d)}\right) \right) H_\infty + \frac{2d_0}{d} \cdot \log(2|\mathcal{A}|).$$

For  $i \in [d]$ , define the indicator random variable,

$$X(s', \text{Dict}) = \mathbf{1}(\exists \text{ a pair of tokens in } \text{enc}_{\text{BPE}}(s') \text{ under Dict appears at least } \log(d) \text{ times}).$$

which captures the event that the string  $s'$  is compressed well by the dictionary Dict under the sequential encoder.

Let  $\text{Dict}_i$  denote the dictionary stored by Algorithm 1 right after  $\text{text}_i$  is processed. The key insight behind this lemma is the following statement, asserting that the sequential encoder satisfies a “monotonicity” property: for any  $j$  and string  $s'$ , if there exists a pair of tokens appearing more than  $\log(d)$  times consecutively in the sequential encoding of  $s'$  under  $\text{Dict}_j$ , then there must exist a pair of tokens appearing more than  $\log(d)$  times consecutively in the greedy encoding of  $s'$  under  $\text{Dict}_i$  for any  $i < j$ . This implies that  $X(s', \text{Dict}_j) \leq X(s', \text{Dict}_i)$  if  $i < j$  for any string  $s'$ . This monotonicity property implies that the last dictionary output by the learner,  $\text{Dict}_d$  sequentially encodes a  $1 - \epsilon$  fraction of the previously seen texts,  $\text{text}_i$  in a way where every pair of tokens appears at most  $\log(d)$  times. While  $\text{Dict}_d$  is correlated with these texts, we can circumvent this correlation by using a martingale argument to prove the statement of the lemma.

**Lemma B.11.** *Let Dict be the dictionary returned by Algorithm 1. Then,*

$$\min \left\{ \Pr \left( \mathbb{E}[X(s', \text{Dict}) | \text{Dict}] \geq 2d_0/d \mid |\text{Dict}| < d_0 \right), \Pr(|\text{Dict}| < d_0) \right\} \leq e^{-\epsilon^2 d / 8 \log^2(1/\delta)}.$$

where  $s'$  is a fresh substring of length  $d$  sampled from the stochastic source.

*Proof.* Let  $\text{Dict}_i$  denote the state of dictionary returned by Algorithm 1 right after  $\text{text}_i$  is processed. Then,  $\text{Dict}_d$  is the final dictionary returned by Algorithm 1. Suppose  $\mathbb{E}[X(s', \text{Dict}_d) | \text{Dict}_d] \geq 2d_0/d$ ,

where  $s'$  is a fresh substring of length  $d$  sampled from the stochastic source. Using monotonicity of the sequential encoder, almost surely for any string  $s'$ ,  $X(s', \text{Dict}_i) \leq X(s', \text{Dict}_j)$  for any  $j > i$ . Therefore,

$$\mathbb{E}[X(s', \text{Dict}_d) | \text{Dict}_d] \geq 2d_0/d \implies \sum_{i=1}^{d-1} \mathbb{E}[X(s', \text{Dict}_i) | \text{Dict}_i] \geq 2d_0 \cdot \frac{d-1}{d} \quad (25)$$

Note in this expectation,  $s'$  is an independent string of length  $d$  sampled from the stochastic source. Since  $\text{Dict}_i$  and  $\text{text}_{i+1}$  are independent, we may instead write,

$$\sum_{i=1}^{d-1} \mathbb{E}[X(\text{text}_{i+1}, \text{Dict}_i) | \text{Dict}_i, \text{text}_i, \text{Dict}_{i-1}, \dots, \text{Dict}_1, \text{text}_1] \geq 2d_0 \cdot \frac{d-1}{d}.$$

For brevity, denote  $X_i = X(\text{text}_{i+1}, \text{Dict}_i)$  and define the filtration  $\mathcal{F}_i = \sigma(\{\text{text}_1, \text{Dict}_1, \dots, \text{text}_i, \text{Dict}_i\})$ . Note that  $\sum_{j=1}^i X_j - \mathbb{E}[X_j | \mathcal{F}_i]$  forms a martingale sequence under the filtration  $\{\mathcal{F}_i : i \in [d]\}$ . Therefore, by the Azuma-Hoeffding inequality, for any  $\eta > 0$ ,

$$\Pr\left(\sum_{i=1}^{d-1} \mathbb{E}[X_i | \mathcal{F}_i] - X_i \leq -\eta\right) \leq e^{-\eta^2}. \quad (26)$$

Under Case I, we have that  $\sum_{i=1}^d X_i \leq d_0$ . Therefore, from eq. (25) and eq. (26),

$$\begin{aligned} \Pr\left(|\text{Dict}| < d_0; \mathbb{E}[X(s', \text{Dict}) | \text{Dict}] \geq 2d_0/d\right) &\leq \Pr\left(\sum_{i=1}^{d-1} X_i < d_0; \sum_{i=1}^{d-1} \mathbb{E}[X_i | \mathcal{F}_i] \geq 2d_0 \cdot \frac{d-1}{d}\right) \\ &\leq \Pr\left(\sum_{i=1}^{d-1} \mathbb{E}[X_i | \mathcal{F}_i] - X_i \geq d_0 \cdot \frac{d-2}{d}\right) \\ &\leq e^{-d_0^2(1-2/d)^2} \\ &\leq e^{-d_0^2/2} = e^{-\epsilon^2 d/8 \log^2(1/\delta)}. \end{aligned}$$

Finally, using the inequality  $\Pr(A, B) = \Pr(A|B) \Pr(B) \geq (\min\{\Pr(A), \Pr(B)\})^2$  completes the proof.  $\square$

**Proofs of Theorem B.3.2 and Theorem B.3.3** If  $\Pr(|\text{Dict}| < d_0) \leq e^{-\epsilon^2 d/8 \log^2(1/\delta)}$  the proof of Theorem B.3.2 concludes. Otherwise, consider the case  $\Pr(|\text{Dict}| < d_0) > e^{-\epsilon^2 d/8 \log^2(1/\delta)}$ , whereby,  $\mathbb{E}[X(s', \text{Dict}) | \text{Dict}] \leq 2d_0/d$  with probability  $\geq 1 - e^{-\epsilon^2 d/8 \log^2(1/\delta)}$  conditioned on  $|\text{Dict}| < d_0$  by Lemma B.11. Recall that when  $|\text{Dict}| < d_0$ , Algorithm 1 uses a parallel implementation of the sequential encoder which chunks a new string into pieces of length  $d$ , denoted  $\{\text{chunk}_i : i \in [d]\}$  and uses the sequential encoder under  $\text{Dict}_d$  to tokenize each chunk. Note that since the source is Markovian, the chunked process  $\{\text{chunk}_i = (X_{id+1}, X_{id+2}, \dots, X_{(i+1)d}) : i = 1, 2, \dots\}$  is also Markovian and ergodic. Therefore, by a similar limiting argument as in Lemma A.4, using the Krylov–Bogolyubov argument (cf. Proposition 4.2 in Chen (2018)) for Markov processes, we have that,

$$\lim_{\ell \rightarrow \infty} \frac{\sum_{i=1}^{\ell} X(\text{chunk}_i, \text{Dict})}{\ell} = \mathbb{E}[X(s', \text{Dict})] \leq \frac{2d_0}{d}.$$

where  $s'$  is a fresh string of length  $d$  sampled with initial state distribution as the stationary measure of the stochastic source. On the remaining (limiting)  $1 - 2d_0/d$  fraction of the chunks, their sequential encodings have every pair of tokens appearing at most  $\log(d)$  times consecutively. Using Theorem 1 of Navarro and Russo (2008), the number of tokens in the encoding of each of these chunks cannot be too large, and satisfies,

$$\begin{aligned} |\text{enc}_{\text{BPE}}(\text{chunk}_i)| \cdot \log |\text{enc}_{\text{BPE}}(\text{chunk}_i)| &\leq 2dH_{\infty} + O(d/\log(d)) \\ \implies |\text{enc}_{\text{BPE}}(\text{chunk}_i)| \cdot \log d &\leq 2dH_{\infty} + O(d/\log(d)) \end{aligned} \quad (27)$$

For the (limiting)  $2d_0/d$  fraction of the “bad” chunks, their sequential encodings may have one or more pairs of tokens which appear more than  $\log(d)$  times consecutively.

Define  $\mathcal{E}_i = \{X(\text{chunk}_i, \text{Dict}) = 1\}$  where  $\text{Dict} = \text{Dict}_d$  is the dictionary returned by Algorithm 1 and consider the unigram model  $Q_{\text{uni}}(\mathbf{t}) = \frac{1}{2}Q_1(\mathbf{t}) + \frac{1}{2}Q_2(\mathbf{t})$ , which is the uniform mixture of two models,

$$Q_1(\mathbf{t}) \propto \frac{1}{(2|\mathcal{A}|)^{|\mathbf{t}|}}, \quad \text{and} \quad Q_2(\mathbf{t}) = \mathbb{E} \left[ \frac{n_{\mathbf{t}}^1}{|\text{enc}_{\text{BPE.split}}(\text{chunk}_1)|} \middle| \mathcal{E}_1^c \right],$$

and let  $Q_{\text{uni}}(\mathbf{t}_1, \dots, \mathbf{t}_i) = Q_{\#}(j) \prod_{i=1}^i Q_{\text{uni}}(\mathbf{t}_i)$  for some distribution  $Q_{\#}(i)$  over the number of tokens to be chosen later. We will analyze the case where the total number of chunks  $\ell$  is finite and take the limit  $m \rightarrow \infty$  later. Then, the overall loss of the algorithm is,

$$\begin{aligned} \mathcal{L}_m(Q_{\text{uni}} \circ \text{enc}(\cdot)) &= -\mathbb{E}[\log Q_{\text{uni}}(\text{enc}_{\text{BPE.split}}(\mathbf{s}))] \\ &= -\sum_{\mathbf{t} \in \text{Dict}} \mathbb{E}[n_{\mathbf{t}} \log Q_{\text{uni}}(\mathbf{t}) + \log Q_{\text{uni}}(|\text{enc}_{\text{BPE.split}}(\mathbf{s})|)] \\ &\stackrel{(i)}{=} -\sum_{i=1}^{\ell} \mathbb{E} \left[ \sum_{\mathbf{t} \in \text{Dict}} n_{\mathbf{t}}^i \log Q_{\text{uni}}(\mathbf{t}) \right] + \log(m) \\ &= -\sum_{i=1}^{\ell} \mathbb{E} \left[ \sum_{\mathbf{t} \in \text{Dict}} n_{\mathbf{t}}^i \log Q_{\text{uni}}(\mathbf{t}) \middle| \mathcal{E}_i \right] \Pr(\mathcal{E}_i) + \mathbb{E} \left[ \sum_{\mathbf{t} \in \text{Dict}} n_{\mathbf{t}}^i \log Q_{\text{uni}}(\mathbf{t}) \middle| \mathcal{E}_i^c \right] \Pr(\mathcal{E}_i^c) + \log(m). \end{aligned} \tag{28}$$

where  $n_{\mathbf{t}}^i$  is the number of times  $\mathbf{t}$  is observed in the BPE encoding of  $\text{chunk}_i$  and (i) uses the fact that  $|\text{enc}_{\text{BPE.split}}(\mathbf{s})|$  follows some distribution supported on  $[m]$ , which implies its entropy is upper bounded by  $\log(m)$ . First observe that,

$$\begin{aligned} \sum_{i=1}^{\ell} \mathbb{E} \left[ \sum_{\mathbf{t} \in \text{Dict}} n_{\mathbf{t}}^i \log Q_{\text{uni}}(\mathbf{t}) \middle| \mathcal{E}_i^c \right] &\leq \sum_{i=1}^{\ell} \mathbb{E} \left[ |\text{enc}_{\text{BPE}}(\text{chunk}_i)| \cdot \sum_{\mathbf{t} \in \text{Dict}} \frac{n_{\mathbf{t}}^i}{|\text{enc}_{\text{BPE}}(\text{chunk}_i)|} \log Q_{\text{uni}}(\mathbf{t}) \middle| \mathcal{E}_i^c \right] \\ &\stackrel{(i)}{\leq} \ell \left( \frac{2dH_{\infty} + O(d/\log(d))}{\log(d)} \right) \sum_{\mathbf{t} \in \text{Dict}} Q_2(\mathbf{t}) \log Q_{\text{uni}}(\mathbf{t}) \end{aligned}$$

where (i) uses the upper bound on  $|\text{enc}_{\text{BPE.split}}(\text{chunk}_i)|$  under the event  $\mathcal{E}_i^c$  (eq. (27)). Since  $Q_{\text{uni}}(\mathbf{t}) = \frac{1}{2}Q_1(\mathbf{t}) + \frac{1}{2}Q_2(\mathbf{t}) \geq \frac{1}{2}Q_2(\mathbf{t})$  and  $Q_2$  is a distribution supported on at most  $d$  tokens, this term results in the upper bound,

$$\sum_{i=1}^{\ell} \mathbb{E} \left[ \sum_{\mathbf{t} \in \text{Dict}} n_{\mathbf{t}}^i \log Q_{\text{uni}}(\mathbf{t}) \middle| \mathcal{E}_i^c \right] \leq \ell \left( \frac{2dH_{\infty} + O(d/\log(d))}{\log(d)} \right) \log(2d). \tag{29}$$

On the other hand, since  $Q_{\text{uni}}(\mathbf{t}) \geq \frac{1}{2}Q_1(\mathbf{t})$ ,

$$\begin{aligned} \sum_{i=1}^{\ell} \mathbb{E} \left[ \sum_{\mathbf{t} \in \text{Dict}} n_{\mathbf{t}}^i \log(1/Q_{\text{uni}}(\mathbf{t})) \middle| \mathcal{E}_i \right] &\leq \sum_{i=1}^{\ell} \mathbb{E} \left[ \sum_{\mathbf{t} \in \text{Dict}} n_{\mathbf{t}}^i \log(2/Q_1(\mathbf{t})) \middle| \mathcal{E}_i \right] \\ &\leq \sum_{i=1}^{\ell} \mathbb{E} \left[ \sum_{\mathbf{t} \in \text{Dict}} n_{\mathbf{t}}^i (\log(2) + |\mathbf{t}| \log(2|\mathcal{A}|)) \middle| \mathcal{E}_i \right] \\ &\leq \ell d \log(2) + \ell d \log(2|\mathcal{A}|) \end{aligned} \tag{30}$$

where the last inequality uses the fact that  $\sum_{\mathbf{t} \in \text{Dict}} n_{\mathbf{t}}^i \leq d$  and  $\sum_{\mathbf{t} \in \text{Dict}} |\mathbf{t}| n_{\mathbf{t}}^i = d$  computes the length of  $\text{chunk}_i$ .

Overall, since  $\sum_{i=1}^{\ell} \Pr(\mathcal{E}_i) \leq 2d_0/d$  by eq. (27), combining this with eqs. (28) to (30),

$$\mathcal{L}_m(Q_{\text{uni}} \circ \text{enc}(\cdot)) \leq \left(1 - \frac{2d_0}{d}\right) \ell \left( \frac{2dH_{\infty} + O(d/\log(d))}{\log(d)} \right) \log(2d) + \frac{2d_0}{d} \ell d \log(4|\mathcal{A}|).$$

Dividing throughout by the length of the character sequence  $m \in [d(\ell - 1), d\ell]$  and letting  $\ell \rightarrow \infty$ ,

$$\min_{Q \in \mathcal{Q}_{1\text{-gram}}} \mathcal{L}(Q \circ \text{enc}(\cdot)) \leq \mathcal{L}(Q_{\text{uni}} \circ \text{enc}(\cdot)) \leq \left(1 - \frac{2d_0}{d}\right) \left(2H_\infty + O\left(\frac{1}{\log(d)}\right)\right) + \frac{2d_0}{d} \log(4|\mathcal{A}|).$$

## C Additional Theoretical Results II: Learning the likelihood model

The guarantees we prove in Theorems 3.1, 3.6 and B.2 on various tokenizers assume that the downstream model is trained optimally. In practice, these models are trained from a finite dataset and the sample complexity of learning this likelihood model scales with the number of tokens in the dictionary. In this section, we step away from the transformer architecture and focus on analyzing the performance of a simple estimator for the unigram model based on Laplace smoothing. We leave the problem of analyzing the finite-sample statistical error of simple transformer models trained with gradient descent as an interesting open direction for future research.

The result of Theorem 3.1 establishes that under appropriate assumptions on the Markov source, there exists a tokenizer  $\mathcal{T}$  and a unigram model over tokens  $Q^* \in \mathcal{Q}_{1\text{-gram}}$  such that,

$$\begin{aligned} \lim_{m \rightarrow \infty} \frac{1}{m} \mathbb{E} [\log(1/Q^*(\text{enc}(s)))] \\ \leq (1 + \varepsilon) \cdot \lim_{m \rightarrow \infty} \frac{1}{m} \mathbb{E} [\log(1/P(s))] \end{aligned}$$

Or in other words,

$$\lim_{m \rightarrow \infty} \frac{1}{m} \text{KL}(P, Q^*(\text{enc}(\cdot))) \leq \varepsilon \cdot \lim_{m \rightarrow \infty} \frac{1}{m} \mathbb{E} [\log(1/P(s))].$$

This implies that with the appropriate tokenization, the measure associated to the string by the best unigram model over tokens is close to that induced by the true Markov distribution over characters in KL divergence. In this section, we establish finite-sample guarantees on learning  $Q^*$  specifically for the LZW tokenizer. The approach we consider for distribution learning is a smoothed Laplace estimator described in more detail in Algorithm 2.

For any constant  $\theta \in (0, 1)$ , define  $\mathcal{E}_\theta$  as the event that every maximal token  $t$  (Definition A.5) in the LZW dictionary satisfies  $1/d^{1-\theta} \geq \max_a P(t|a) \geq \delta/d^{1+\theta}$ . By Lemmas A.10 and A.11 when the LZW tokenizer is trained on a dataset of size  $\tilde{\Omega}_\delta(d)$  drawn from a stochastic source satisfying Assumption 3.2,  $\mathcal{E}_\theta$  occurs with probability  $\geq 1 - d^{-\Omega_{\theta,\delta}(\log(d))}$ .

**Theorem C.1.** *Consider any constant  $\theta \in (0, 1)$ , failure probability  $\eta \in (0, 1)$  and approximation error  $\xi \in (0, 1)$ . Assume that the learnt LZW tokenizer  $\mathcal{T}_{\text{LZW}}$  satisfies the event  $\mathcal{E}_\theta$ , which occurs with probability  $\geq 1 - d^{-\Omega_{\theta,\delta}(\log(d))}$ . Assume that  $d^{1-3\theta} \geq 1 + \delta^{-2}$  and that the stochastic source satisfies Assumption 3.2. For an absolute constant  $C > 0$ , assume that the size of the training dataset is at least  $n_{lm}^*(\xi)$ , where,*

$$n_{lm}^* \triangleq \frac{Cd^{1+\theta} \log^3(d/\eta\delta) \log \log(d/\eta)}{\delta\xi^2}$$

*Then, Algorithm 2 learns a unigram model  $\hat{Q}$  such that,*

$$\mathcal{L}(\hat{Q} \circ \text{enc}_{\text{gre}}(\cdot)) \leq (1 + \xi) \min_{Q \in \mathcal{Q}_{1\text{-gram}}} \mathcal{L}(Q \circ \text{enc}_{\text{gre}}(\cdot))$$

*with probability  $\geq 1 - \eta$ .*

In conjunction with Theorem 3.6, this gives end-to-end guarantees on the cross-entropy loss of the LZW tokenizer (with vocabulary size  $\leq d$ ) with the Laplace estimator as the downstream unigram model. We instantiate this result choosing  $\theta = 0.01$  in Theorem C.1.

**Corollary C.2.** *Choose any  $\xi \in (0, 1)$ . Suppose the data source satisfies Assumption 3.2. On a dataset of size  $\tilde{\Omega}_\delta(d)$  drawn from the source, train an LZW tokenizer  $\mathcal{T}_{\text{LZW}}$  with  $d$  tokens. Subsequently, using Algorithm 2, learn a unigram model  $\hat{Q}$  using a dataset of size at least  $\tilde{\Omega}(d^{1.01}/\delta\xi^2)$  drawn from the source. Then, with probability  $\geq 1 - d^{-\Omega_\delta(\log(d))}$ ,*

$$\mathcal{L}(\hat{Q} \circ \text{enc}_{\text{gre}}(\cdot)) \leq \frac{1 + \xi}{1 - \varepsilon} \min_Q \mathcal{L}(Q),$$

*where  $\varepsilon = \frac{\log(1/\delta)}{0.99 \log(d)}$ .*

The analysis of Theorem C.1 relies on showing that the distribution over tokens induced when a string sampled from the data source is encoded into tokens by the greedy encoder and the LZW dictionary is a Markov process. In general, given a set of previously sampled tokens  $t_1, \dots, t_i$ , the next token  $t_{i+1}$  is sampled from the distribution  $P(t_{i+1}|t_i; \forall j \in [i], t_{i-j+1} \dots t_i t_{i+1} \notin \text{Dict})$ . The conditioning is to simply guarantee that the previous tokens which were sampled were indeed maximal, since if  $t_i t_{i+1} \in \text{Dict}$ , then the previous token returned would in fact have been this longer token and not  $t_i$  (and likewise for  $t_{i-1} t_i t_{i+1}$  and so on). While in general, this process is complicated and depends on all the previous tokens sampled, for the LZW dictionary, we show that the conditioning  $\{\forall j \in [i], t_{i-j+1} \dots t_i t_{i+1} \notin \text{Dict}\}$  can be removed, thereby resulting in a simple Markov process over tokens.

Furthermore, we establish that this Markov process has a relatively large spectral gap. The optimal unigram model ends up being the stationary distribution over tokens induced by greedy encoder. Given the large spectral gap of the Markov process over tokens, estimating the stationary distribution of this process in KL divergence ends up being closely related to estimating a distribution from i.i.d. samples in the same metric. For this problem, the de-facto choice of estimator is the Laplace estimator, and several existing results provide finite-sample bounds on the KL divergence (Braess and Sauer, 2004; Han et al., 2021; Mourtada and Gaïffas, 2022). The Laplace estimator (Line 6 of Algorithm 2) is simply a smoothed empirical estimate to account for the degeneracy of the KL divergence in its second argument as any coordinate approaches 0. The non-i.i.d.ness of the Markov process is circumvented by using concentration inequalities which are a function of the spectral gap (Naor et al., 2020).

---

**Algorithm 2** Training likelihood model on tokens

---

**Input:** A training dataset of size  $n_{\text{lm}}$ , likelihood model class  $\mathcal{Q}$ , likelihood model training algorithm `TrainLM`

**Output:** Likelihood model  $Q \in \mathcal{Q}$ .

- 1: Tokenize the training dataset into a sequence of tokens  $\mathcal{T} = (t_1, \dots, t_i)$ .
- 2: Train a likelihood model  $Q$  on the tokenized dataset  $\mathcal{T}$  using the `TrainLM( $\mathcal{T}, Q$ )` subroutine.

*// In the case of  $Q = Q_{1\text{-gram}}$  use the Laplace estimator*

**def** `TrainLM( $\mathcal{T}, Q_{1\text{-gram}}$ )`:

- 3: Truncate the dataset to the first  $n' = \lfloor n_{\text{lm}} / \ell_{\text{max}} \rfloor$  tokens where  $\ell_{\text{max}} = 4 \log(d|\mathcal{A}|) / \delta$ . Let the truncated dataset be  $\mathcal{T}_{\text{trunc}}$
- 4: Construct the unigram model  $\hat{Q}$  with  $\hat{Q}_{\#} = \text{Unif}([m])$  and  $\hat{Q}_{\text{tok}}(t) = \frac{n_t + 1}{n_t + |\text{Dict}|}$ .

*//  $n_t$  is the number of times  $t$  appears in  $\mathcal{T}_{\text{trunc}}$ .*

*// Test sequences are assumed to be of length  $m$ .*

---

### C.1 Proof of Theorem C.1

Since  $\mathcal{T}_{\text{LZW}}$  uses the greedy encoder, the cross-entropy loss of the unigram model learnt by Algorithm 2 is,

$$\begin{aligned}
& \mathcal{L}(\hat{Q} \circ \text{enc}_{\text{gre}}(\cdot)) - \min_{Q \in \mathcal{Q}_{1\text{-gram}}} \mathcal{L}(Q \circ \text{enc}_{\text{gre}}(\cdot)) \\
&= \max_{Q \in \mathcal{Q}_{1\text{-gram}}} \lim_{m \rightarrow \infty} \frac{1}{m} \mathbb{E}[\log(Q(\text{enc}_{\text{gre}}(s)) / \hat{Q}(\text{enc}_{\text{gre}}(s)))] \\
&\stackrel{(i)}{=} \max_{Q \in \mathcal{Q}_{1\text{-gram}}} \lim_{m \rightarrow \infty} \frac{1}{m} \mathbb{E} \left[ |\text{enc}_{\text{gre}}(s)| \sum_{t \in \text{Dict}} \frac{n_t}{|\text{enc}_{\text{gre}}(s)|} \log(Q_{\text{tok}}(t) / \hat{Q}_{\text{tok}}(t)) \right] + \frac{\log(m)}{m} \\
&\stackrel{(ii)}{\leq} \lim_{m \rightarrow \infty} \frac{1}{m} \mathbb{E} \left[ |\text{enc}_{\text{gre}}(s)| \sum_{t \in \text{Dict}} \frac{n_t}{|\text{enc}_{\text{gre}}(s)|} \log \left( \frac{n_t / |\text{enc}_{\text{gre}}(s)|}{\hat{Q}_{\text{tok}}(t)} \right) \right] + \frac{\log(m)}{m}
\end{aligned}$$

where in (i) we use the fact that  $\hat{Q}_{\#} = \text{Unif}([m])$  and in (ii) we take the  $\max\{\cdot\}$  inside the limit and the expectation (Fatou's lemma and Jensen's inequality) and plug in the maximizer of the negative cross-entropy,  $Q_{\text{tok}}(t) = \frac{n_t}{|\text{enc}_{\text{gre}}(s)|}$ . Note that  $\lim_{m \rightarrow \infty} \frac{n_t}{|\text{enc}_{\text{gre}}(s)|} \stackrel{\text{a.s.}}{=} Q_{\text{MLE}}(t)$  by Lemma A.4.

Moreover, since  $|\text{enc}(s)|/m \leq 1$  and  $\widehat{Q}_{\text{tok}}(\mathbf{t}) > 0$  surely, by the Dominated Convergence Theorem,

$$\mathcal{L}(\widehat{Q} \circ \text{enc}_{\text{gre}}(\cdot)) - \min_{Q \in \mathcal{Q}_{1\text{-gram}}} \mathcal{L}(Q \circ \text{enc}_{\text{gre}}(\cdot)) \leq \lim_{m \rightarrow \infty} \frac{1}{m} \mathbb{E}[|\text{enc}_{\text{gre}}(s)|] \cdot \text{KL}(Q_{\text{MLE}}, \widehat{Q}_{\text{tok}}) \quad (31)$$

By eq. (6), we have that for any tokenizer using the greedy encoder,

$$\lim_{m \rightarrow \infty} \frac{|\text{enc}_{\text{gre}}(s)| (H(Q_{\text{MLE}}, P) - \log(1/\delta))}{m} \stackrel{\text{a.s.}}{\leq} H_{\infty}.$$

Furthermore under the event  $\mathcal{E}_{\theta}$  which implies that the learnt dictionary is  $(1 - \theta)$ -heavy hitting (cf. Definition A.5), which implies that,

$$H(Q_{\text{MLE}}, P) \geq (1 - \theta) \log(d).$$

Therefore, by almost sure boundedness, we have that,

$$\lim_{m \rightarrow \infty} \frac{1}{m} \mathbb{E}[|\text{enc}_{\text{gre}}(s)|] \leq \frac{H_{\infty}}{(1 - \theta) \log(d) - \log(1/\delta)} \leq \frac{\min_{Q \in \mathcal{Q}_{1\text{-gram}}} \mathcal{L}(Q \circ \text{enc}_{\text{gre}}(\cdot))}{(1 - \theta) \log(d) - \log(1/\delta)}$$

Putting this together with eq. (31), we have that,

$$\mathcal{L}(\widehat{Q} \circ \text{enc}_{\text{gre}}(\cdot)) \leq \left(1 + \text{KL}(Q_{\text{MLE}}, \widehat{Q}_{\text{tok}})\right) \min_{Q \in \mathcal{Q}_{1\text{-gram}}} \mathcal{L}(Q \circ \text{enc}_{\text{gre}}(\cdot)), \quad (32)$$

which uses the assumption  $(1 - \theta) \log(d) \geq 1 + \log(1/\delta)$ . In the remainder of the proof we upper bound the KL term.

By the law of large numbers established in eq. (34) and the fact that  $\frac{n_{\mathbf{t}}}{\sum_{\mathbf{t}'} n_{\mathbf{t}'}} \in [0, 1]$ , we have that,

$$Q_{\text{MLE}}(\mathbf{t}) = \lim_{m \rightarrow \infty} \mathbb{E} \left[ \frac{n_{\mathbf{t}}}{\sum_{\mathbf{t}'} n_{\mathbf{t}'}} \right] = \lim_{m \rightarrow \infty} \frac{\mathbb{E}[n_{\mathbf{t}}]}{\mathbb{E}[\sum_{\mathbf{t}'} n_{\mathbf{t}'}]} = \pi(\mathbf{t}),$$

where  $\pi(\mathbf{t})$  denote the stationary distribution over tokens induced by the greedy encoding process, which exists for the LZW tokenizer. This distribution is in fact an ergodic Markov process, as we discuss next.

By Lemmas A.10 and A.11, for any constant  $\theta \in (0, 1)$ , with probability  $\geq 1 - d^{-\Omega_{\theta, \delta}(\log(d))}$ , every maximal token in the the LZW dictionary satisfies  $1/d^{1-\theta} \geq \max_a P(\mathbf{t}|a) \geq \delta/d^{1+\theta}$ . Let  $S_{\text{gre}}$  denote the set of tokens which have a non-zero probability (over a string drawn from the Markov source) of being chosen by the greedy encoder while encoding the string. More importantly, note that for any sequence of tokens  $t_1, \dots, t_i$ , the next token is necessarily in  $S_{\text{gre}}$  and can be any token in this set. The reason for this is that for any  $t_i, \mathbf{t} \in S_{\text{gre}}$ , the concatenation  $t_i \mathbf{t} \notin S_{\text{gre}}$  since  $\max_{a \in \mathcal{A}} P(t_i \mathbf{t}|a) \leq 1/\delta d^{2(1-\theta)}$ , which is smaller than the  $\max_{a \in \mathcal{A}} P(\mathbf{t}'|a) \geq \delta/d^{1+\theta}$  for any token  $\mathbf{t}' \in S_{\text{gre}}$  as long as  $d^{1-3\theta} \geq 1/\delta^2$ . This constraint implies that in the sampling procedure in Figure 7, it suffices to drop the conditioning on the event  $t_j t_{j+1} \dots t_i \mathbf{t} \notin \text{Dict}$  while sampling the next token  $\mathbf{t}$ . This condition automatically implies that the sequence of tokens conditionally follows a Markov process with  $\Pr(t_{i+1} = \mathbf{t} | t_1, \dots, t_i) = P(\mathbf{t} | \text{last}(t_i))$ . Since the probability of every transition is lower bounded, this means that the Markov chain is ergodic. Moreover, the pseudo-spectral gap (Naor et al., 2020),  $1 - \lambda$  can be lower bounded by the Dobrushin contraction coefficient,  $\kappa$ ,

$$\begin{aligned} 1 - \lambda &\leq \kappa \triangleq \max_{(\mathbf{t}, \mathbf{t}') \in \text{Dict}^2} \|\Pr(\cdot | \mathbf{t}) - \Pr(\cdot | \mathbf{t}')\|_{\text{TV}} \\ &= \max_{(\mathbf{t}, \mathbf{t}') \in \text{Dict}^2} 1 - \sum_{\mathbf{t}'' \in \text{Dict}} \min\{\Pr(\mathbf{t}'' | \mathbf{t}), \Pr(\mathbf{t}'' | \mathbf{t}')\} \\ &\leq 1 - \delta d/d^{1+\theta} \\ &= 1 - \delta d^{-\theta}. \end{aligned} \quad (33)$$

Recall that the learner is given a training dataset of  $n_{\text{lm}}$  characters to train the likelihood model. By Lemma A.8, with probability  $\geq 1 - d^{-\Omega(\log(d/\delta)/\delta)}$ , in the run of the LZW tokenization algorithm, every token in the dictionary has length at most  $\ell_{\text{max}} = 4 \log(d|\mathcal{A}|)/\delta$ . Therefore, suppose the learner



always truncates the dataset to the first  $n' = \lfloor n_{\text{lm}}/\ell_{\text{max}} \rfloor$  tokens and runs the Laplace estimator on this truncated dataset. With this, we move onto upper bounding,

$$\text{KL}(Q_{\text{MLE}}, \hat{Q}_{\text{tok}}) = \sum_{\mathbf{t} \in \text{Dict}} \pi(\mathbf{t}) \log \left( \pi(\mathbf{t}) / \hat{Q}_{\text{tok}}(\mathbf{t}) \right)$$

which necessitates lower bounding  $\hat{Q}_{\text{tok}}(\mathbf{t})$  for every  $\mathbf{t}$ . Recall that the learner's estimate  $\hat{Q}(\mathbf{t})$  in Algorithm 2 is the Laplace estimator,  $\frac{n_{\mathbf{t}}+1}{\sum_{\mathbf{t}'} n_{\mathbf{t}'} + |\text{Dict}|}$ , where  $\{n_{\mathbf{t}} : \mathbf{t} \in \text{Dict}\}$  is computed by truncating the dataset to the first  $n'$  tokens. Firstly, by invoking Corollary 1.3 of Naor et al. (2020) for the function  $n_{\mathbf{t}} = \sum_{i=1}^{n'} \mathbb{I}(\mathbf{t}_i = \mathbf{t})$ ,

$$\Pr \left( |n_{\mathbf{t}} - \mathbb{E}[n_{\mathbf{t}}]| \geq c \sqrt{\frac{\mathbb{E}[n_{\mathbf{t}}]}{1-\lambda} \cdot \log(1/\eta)} \right) \leq \eta \quad (34)$$

for a universal constant  $c > 0$ . In particular, this implies that with probability  $\geq 1 - \eta$ , simultaneously for all  $\mathbf{t}$ ,

$$|n_{\mathbf{t}} - \mathbb{E}[n_{\mathbf{t}}]| \leq \Delta_{\mathbf{t}} \triangleq \sqrt{\frac{d^\theta}{\delta} \mathbb{E}[n_{\mathbf{t}}] \cdot \log(|\text{Dict}|/\eta)}, \text{ and, } \mathbb{E}[n_{\mathbf{t}}] - n_{\mathbf{t}} \geq \mathbb{E}[n_{\mathbf{t}}].$$

Under this event, for any  $\mathbf{t}$ , the estimate is lower bounded by,

$$\begin{aligned} \hat{Q}_{\text{tok}}(\mathbf{t}) &= \frac{n_{\mathbf{t}} + 1}{n' + |\text{Dict}|} \geq \frac{\mathbb{E}[n_{\mathbf{t}}] + 1 - \min\{\mathbb{E}[n_{\mathbf{t}}], \Delta_{\mathbf{t}}\}}{n' + |\text{Dict}|} \\ &\geq \max \left\{ \pi(\mathbf{t}) - \frac{(\Delta_{\mathbf{t}} - 1) n' + |\text{Dict}| \mathbb{E}[n_{\mathbf{t}}]}{(n')^2 + n' |\text{Dict}|}, \frac{1}{n' + |\text{Dict}|} \right\} \\ &\geq \max \left\{ \pi(\mathbf{t}) - \frac{\Delta_{\mathbf{t}} n' + |\text{Dict}| \mathbb{E}[n_{\mathbf{t}}]}{(n')^2}, \frac{1}{n' + |\text{Dict}|} \right\} \end{aligned}$$

Suppose the following condition is satisfied,

$$n' = \frac{4r d^\theta |\text{Dict}| \log(|\text{Dict}|/\eta)}{\delta} \quad (\text{C1})$$

for some  $r \geq 4$ . Under this condition, we have that  $n' \geq 2\sqrt{r}\Delta$  and  $n' \geq 4r|\text{Dict}|$ .

**Case I.**  $\Delta_{\mathbf{t}} n' \geq |\text{Dict}| \mathbb{E}[n_{\mathbf{t}}]$ .

In this case, we have the upper bound,

$$\begin{aligned} \hat{Q}_{\text{tok}}(\mathbf{t}) &\geq \max \left\{ \pi(\mathbf{t}) - \frac{2\Delta_{\mathbf{t}}}{n'}, \frac{1}{n' + |\text{Dict}|} \right\} \\ &= \max \left\{ \pi(\mathbf{t}) - 2 \frac{\sqrt{\frac{d^\theta}{\delta} \mathbb{E}[n_{\mathbf{t}}] \cdot \log(|\text{Dict}|/\eta)}}{n'}, \frac{1}{n' + |\text{Dict}|} \right\} \\ &\geq \max \left\{ \pi(\mathbf{t}) - \sqrt{\frac{\pi(\mathbf{t})}{r|\text{Dict}|}}, \frac{1}{2n'} \right\}. \end{aligned}$$

where the last inequality uses eq. (C1).

Consider two sub-cases,

**Sub-case I.**  $\pi(\mathbf{t}) \geq 2/r|\text{Dict}|$ . Define this event  $\mathcal{C}_1$ .

Here,

$$\pi(\mathbf{t}) \log(\pi(\mathbf{t}) / \hat{Q}_{\text{tok}}(\mathbf{t})) \leq -\pi(\mathbf{t}) \log \left( 1 - \sqrt{\frac{1}{\pi(\mathbf{t}) r |\text{Dict}|}} \right) \leq \frac{3}{2} \sqrt{\frac{\pi(\mathbf{t})}{r |\text{Dict}|}}. \quad (35)$$

**Sub-case II.**  $\pi(\mathbf{t}) \leq 2/r|\text{Dict}|$ . Define this event  $\mathcal{C}_{\text{II}}$ .

Here,

$$\begin{aligned}\pi(\mathbf{t}) \log(\pi(\mathbf{t})/\hat{Q}_{\text{tok}}(\mathbf{t})) &\leq \pi(\mathbf{t}) \log(2n'\pi(\mathbf{t})) \leq \max \left\{ 0, \frac{2}{r|\text{Dict}|} \log \left( \frac{4n'}{r|\text{Dict}|} \right) \right\} \\ &\leq \frac{2}{r|\text{Dict}|} \log(16d^\theta \log(|\text{Dict}|/\eta))\end{aligned}\quad (36)$$

**Case II.**  $\Delta_{\mathbf{t}} n' < |\text{Dict}| \mathbb{E}[n_{\mathbf{t}}]$ . Define this event  $\mathcal{C}_{\text{III}}$ .

In this case we have the upper bound,

$$\hat{Q}_{\text{tok}}(\mathbf{t}) \geq \pi(\mathbf{t}) - \frac{2|\text{Dict}| \mathbb{E}[n_{\mathbf{t}}]}{(n')^2} \geq \pi(\mathbf{t}) - \frac{\pi(\mathbf{t})}{2r}$$

where the last inequality follows from eq. (C1). This implies that,

$$\pi(\mathbf{t}) \log(\pi(\mathbf{t})/\hat{Q}_{\text{tok}}(\mathbf{t})) \leq -\pi(\mathbf{t}) \log(1 - 1/2r) \leq \frac{\pi(\mathbf{t})}{r}. \quad (37)$$

By using the geometric ergodicity of this Markov process (eq. (33)), when  $n'$  tokens are sampled from an arbitrary initial distribution,

$$(1 - \kappa^{n'}) \pi(\mathbf{t}) \leq \frac{\mathbb{E}[n_{\mathbf{t}}]}{n'} \leq \kappa^{n'} + (1 - \kappa^{n'}) \pi(\mathbf{t}) \implies \pi(\mathbf{t}) \leq \frac{\hat{Q}_{\text{tok}}(\mathbf{t})}{1 - e^{-4r|\text{Dict}| \log(|\text{Dict}|/\eta)}} = \frac{\hat{Q}_{\text{tok}}(\mathbf{t})}{1 - d^{-r}}$$

where in the implication, we use the condition on  $n'$  in eq. (C1) and the bound on the contraction coefficient  $\kappa$  in eq. (33).

$$\begin{aligned}\text{KL}(Q_{\text{MLE}}, \hat{Q}_{\text{tok}}) &= \sum_{\mathbf{t} \in \text{Dict}} \pi(\mathbf{t}) \log(\pi(\mathbf{t})/\hat{Q}_{\text{tok}}(\mathbf{t})) \\ &\leq \sum_{\mathbf{t} \in \text{Dict}} \frac{\pi(\mathbf{t})}{1 - d^{-r}} \log(\pi(\mathbf{t})/\hat{Q}_{\text{tok}}(\mathbf{t})) - \log(1 - d^{-r}) \\ &\leq \frac{1}{1 - d^{-r}} \sum_{\mathbf{t} \in \text{Dict}} \mathbb{I}(\mathcal{C}_{\text{I}}) \pi(\mathbf{t}) \log(\pi(\mathbf{t})/\hat{Q}_{\text{tok}}(\mathbf{t})) + \mathbb{I}(\mathcal{C}_{\text{II}}) \pi(\mathbf{t}) \log(\pi(\mathbf{t})/\hat{Q}_{\text{tok}}(\mathbf{t})) + \mathbb{I}(\mathcal{C}_{\text{III}}) \pi(\mathbf{t}) \log(\pi(\mathbf{t})/\hat{Q}_{\text{tok}}(\mathbf{t})) + 2d^{-r} \\ &\leq \frac{1}{1 - d^{-r}} \sum_{\mathbf{t} \in \text{Dict}} \underbrace{\mathbb{I}(\mathcal{C}_{\text{I}}) \frac{3}{2} \sqrt{\frac{\pi(\mathbf{t})}{r|\text{Dict}|}}}_{\text{eq. (35)}} + \underbrace{\mathbb{I}(\mathcal{C}_{\text{II}}) \frac{2}{r|\text{Dict}|} \log(16d^\theta \log(|\text{Dict}|/\eta))}_{\text{eq. (36)}} + \underbrace{\mathbb{I}(\mathcal{C}_{\text{III}}) \frac{\pi(\mathbf{t})}{r}}_{\text{eq. (37)}} + 2d^{-r} \\ &\leq \frac{1}{1 - d^{-r}} \left( \frac{3}{2} \sqrt{\frac{|\text{Dict}|}{r|\text{Dict}|}} + \frac{2}{r} \log(16d^\theta \log(|\text{Dict}|/\eta)) + \frac{1}{r} \right) + 2d^{-r} \\ &\leq \frac{5}{\sqrt{r}} \log(16d^\theta \log(|\text{Dict}|/\eta))\end{aligned}$$

Combining with eq. (32), we get the bound,

$$\begin{aligned}\mathcal{L}(\hat{Q} \circ \text{enc}_{\text{gre}}(\cdot)) &\leq \left( 1 + \text{KL}(Q_{\text{MLE}}, \hat{Q}) \right) \min_{Q \in \mathcal{Q}_{1\text{-gram}}} \mathcal{L}(Q \circ \text{enc}_{\text{gre}}(\cdot)) \\ &\leq \left( 1 + \frac{5}{\sqrt{r}} \log(16d^\theta \log(d/\eta)) \right) \min_{Q \in \mathcal{Q}_{1\text{-gram}}} \mathcal{L}(Q \circ \text{enc}_{\text{gre}}(\cdot)).\end{aligned}$$

Rescaling  $r$  to be  $r(\log(16d^\theta \log(d/\eta)))^2$  completes the proof.

## D Additional Theoretical Results III: The generalization ability of tokenizers

The proofs of the upper bounds in the paper (Theorems 3.6 and B.2) relied on showing that the entropy  $H(Q_{\text{MLE}}, P)$  is large, or in other words, the algorithm typically encodes new strings into long length (i.e. low probability under  $P$ ) tokens. This statement about generalization to new strings is fundamentally different from having a tokenizer which compresses the training dataset well. In other words, consider the following modification: the measure  $Q_{\text{MLE}}$  is defined as the expected empirical distribution over tokens when a new string is encoded into tokens, and not on the source dataset used to construct the dictionary. Suppose the definition of  $Q_{\text{MLE}}$  is changed to the empirical distribution over tokens in the source dataset. Under this new definition of the MLE unigram model, the largeness of the  $H(Q_{\text{MLE}}, P)$  metric, in a sense, captures compressing the source dataset well. However, we show that in general, this does not result in good tokenizers that minimize the population cross-entropy loss, suffering from  $\min_{Q \in \mathcal{Q}_{1\text{-gram}}} \mathcal{L}(Q \circ \text{enc}(\cdot)) \approx H(\pi) \gg H_\infty$ .

**Theorem D.1.** *Consider the stochastic source in example A.1 having entropy rate  $H_\infty = \delta \log(1/\delta) + (1 - \delta) \log(1/(1 - \delta))$ . Consider a training dataset of size  $n$ . For a dictionary  $\text{Dict}$  and  $\mathbf{t} \in \text{Dict}$ , define  $\hat{Q}_{\text{MLE}}(\mathbf{t}) = \frac{n_{\mathbf{t}}(\mathbf{s}_{\text{src}})}{|\text{enc}(\mathbf{s}_{\text{src}})|}$  as the empirical distribution over tokens induced by the greedy encoder when encoding the training dataset,  $\mathbf{s}_{\text{src}}$ . There exists a dictionary  $\text{Dict}$  such that with probability  $\geq 1 - e^{-\Omega(\sqrt{n})}$  over the training dataset,*

$$H(\hat{Q}_{\text{MLE}}, P_\gamma) \geq nH_\infty(1 - O(n^{-1/4}))$$

is large. However, for this dictionary, for any encoding algorithm (including the greedy encoder), the resulting tokenizer  $\mathcal{T} = (\text{Dict}, \emptyset, \text{enc}(\cdot), \text{dec}(\cdot))$  satisfies,

$$\min_{Q \in \mathcal{Q}_{1\text{-gram}}} \mathcal{L}(Q \circ \text{enc}(\cdot)) \geq (1 - \varepsilon)H(\pi)$$

where  $\varepsilon = 2ne^{-nH_\infty(1 - O(n^{-1/4}))}$ .

*Proof.* Suppose the entire training dataset was compressed into a single token,  $\mathbf{t}_{\text{src}}$ . The dictionary is  $\mathcal{A} \cup \mathbf{t}_{\text{src}}$ . In the following argument, we show that the number of occurrences,  $n_{\mathbf{t}_{\text{src}}}$ , of the entire training dataset  $\mathbf{t}_{\text{src}}$  in a new string of length  $m$  generated from the stochastic source,  $\mathbf{s}$ , converges to its expectation as  $m \rightarrow \infty$ . Let  $\pi_n^{(i)}$  denote the stationary distribution of the Markov process induced by the stochastic source over length- $n$  strings with a shift of  $i$  from the starting position, and let  $n_{\mathbf{t}}^{(i)}$  denote the number of times  $\mathbf{t}$  appears in the training dataset starting at the position  $i + rn$  for some  $r > 0$ . Then,

$$\lim_{m \rightarrow \infty} \frac{n_{\mathbf{t}_{\text{src}}}}{m} = \frac{1}{n} \lim_{m \rightarrow \infty} \sum_{i=0}^{n-1} \frac{n_{\mathbf{t}_{\text{src}}}^{(i)}}{m/n} \stackrel{\text{a.s.}}{=} \frac{1}{n} \sum_{i=0}^{n-1} \mathbb{E}_{\mathbf{t}' \sim \pi_n^{(i)}} [P(\mathbf{t}_{\text{src}} | \mathbf{t}')] \leq \max_{a \in \mathcal{A}} P(\mathbf{t}_{\text{src}} | a). \quad (38)$$

The second equation follows by considering the Markov process induced over length  $n$  strings and applying the Krylov–Bogolyubov argument for ergodic and homogeneous Markov processes.

In Lemma D.2, we show that with probability  $\geq 1 - e^{-\Omega(\sqrt{n})}$ , the token  $\mathbf{t}_{\text{src}}$  constructed from the source dataset satisfies,  $\max_{a \in \mathcal{A}} P(\mathbf{t} | a) \leq e^{-nH_\infty(1 - O(n^{-1/4}))}$ . In other words, the source string has exponentially small probability. Combining this with eq. (38), with probability  $\geq 1 - e^{-\Omega(\sqrt{n})}$  over the source dataset, the number of occurrences of the substring  $\mathbf{t}_{\text{src}}$  in a new string  $\mathbf{s}$  is upper bounded by,

$$\lim_{m \rightarrow \infty} \frac{n_{\mathbf{t}_{\text{src}}}}{m} \stackrel{\text{a.s.}}{\leq} e^{-nH_\infty(1 - O(n^{-1/4}))} \triangleq \varepsilon/2n.$$

By the Krylov–Bogolyubov argument, for each  $a \in \mathcal{A} = \{0, 1\}$ ,  $\lim_{m \rightarrow \infty} \frac{n_a}{m} \stackrel{\text{a.s.}}{=} \pi(a)$ . More importantly, the number of times  $a$  is made as a token is upper bounded by  $n_a$  and lower bounded by  $n_a - nn_{\mathbf{t}_{\text{src}}}$ . Therefore,

$$(1 - \varepsilon)\pi(a) = \pi(a) - \frac{\varepsilon}{2} \stackrel{\text{a.s.}}{\leq} \lim_{m \rightarrow \infty} \frac{n_a}{m} \stackrel{\text{a.s.}}{\leq} \pi(a) = \frac{1}{2} \quad (39)$$

Finally, putting everything together,

$$\begin{aligned}
\min_{Q \in \mathcal{Q}_{1\text{-gram}}} \lim_{m \rightarrow \infty} \frac{1}{m} \mathcal{L}_m(Q \circ \text{enc}(\cdot)) &= \min_{Q \in \mathcal{Q}_{1\text{-gram}}} \lim_{m \rightarrow \infty} -\frac{1}{m} \mathbb{E} \left[ \log(Q_{\#}(|\text{enc}(\mathbf{s})|)) + \sum_{t \in \text{Dict}} n_t \log Q_{\text{tok}}(t) \right] \\
&\geq \min_{Q \in \mathcal{Q}_{1\text{-gram}}} \lim_{m \rightarrow \infty} -\frac{1}{m} \mathbb{E} \left[ \sum_{a \in \mathcal{A}} n_a \log Q_{\text{tok}}(a) \right] \\
&\stackrel{(i)}{\geq} \min_{Q \in \mathcal{Q}_{1\text{-gram}}} -(1 - \varepsilon) \sum_{a \in \mathcal{A}} \pi(a) \log Q_{\text{tok}}(a) \\
&\geq (1 - \varepsilon) H(\pi).
\end{aligned}$$

where (i) follows from the lower bound on  $n_a/m$  in eq. (39). This completes the proof.  $\square$

**Lemma D.2.** *With probability  $\geq 1 - e^{-\Omega(\sqrt{n})}$  over the source dataset,*

$$\max_{a \in \mathcal{A}} P(\mathbf{t}_{\text{src}}|a) \leq e^{-nH(\delta)(1-O(n^{-1/4}))}.$$

*Proof.* Let  $X$  denote the number of  $i \in [n-1]$  such that  $\mathbf{s}_i \neq \mathbf{s}_{i+1}$  in  $\mathbf{s}$ , the stochastic source. Since the transition of the Markov process only depends on whether the next character is the same as the previous character, we can write down,

$$\max_{a \in \mathcal{A}} \log P(\mathbf{t}_{\text{src}}|a) = -(X+1) \log(\delta) - (n-1-X) \log(1-\delta).$$

Note that  $X$  is a sum of  $n-1$  i.i.d. random variables, since  $\mathbb{I}(\mathbf{s}_i \neq \mathbf{s}_{i+1}) \sim \text{Ber}(\delta)$  does not depend on whether  $\mathbf{s}_i = 0$  or  $1$ . In particular, by Hoeffding's inequality, we have that with probability  $\geq 1 - e^{-\Omega(\sqrt{n})}$ ,

$$\left| \frac{1}{n} \max_{a \in \mathcal{A}} \log P(\mathbf{t}_{\text{src}}|a) - H(\delta) \right| \leq O(n^{-1/4}),$$

which uses the fact that  $\mathbb{E}[X] = \delta(n-1)$  and  $H_{\infty} = \delta \log(1/\delta) + (1-\delta) \log(1/(1-\delta))$ . Taking an exponential on both sides proves the statement of the lemma.  $\square$

## E Additional Theoretical Results IV: Interaction between the dictionary and encoding algorithm

In this section, we show another kind of barrier to generalization, which brings out the relationship between the encoding algorithm and the dictionary. We show that there exist dictionaries which generalize under the minimal encoder, i.e. the encoding algorithm which encodes a string into the shortest number of possible tokens, but at the same time, completely fail to generalize under the greedy encoder. This means that in the process of constructing good tokenizers, it does not suffice to think about the dictionary in isolation. Its interaction with the encoding algorithm is pertinent.

**Definition E.1** (minimal encoder). The minimal encoder parses a new string into the fewest possible number of tokens from the dictionary as possible. Ties are broken arbitrarily.

**Theorem E.2.** *There exists a stochastic source parameterized by  $\delta \in (0, 0.5)$  and a dictionary  $\text{Dict}$  such that under the minimal encoder/decoder pair, the resulting tokenizer,  $\mathcal{T} = (\text{Dict}, \emptyset, \text{enc}_{\min}(\cdot), \text{dec}_{\min}(\cdot))$  generalizes near-optimally,*

$$\min_{Q \in \mathcal{Q}_{1\text{-gram}}} \mathcal{L}(Q \circ \text{enc}_{\min}(\cdot)) \leq 1.273 H_{\infty}. \quad (40)$$

Here the entropy rate of the source,  $H_{\infty}$ , is  $\delta \log(\sqrt{2}/\delta) + (1-\delta) \log(1/(1-\delta))$ . However, the same dictionary  $\text{Dict}$  under the greedy encoder/decoder pair, i.e.  $\mathcal{T}' = (\text{Dict}, \emptyset, \text{enc}_{\text{gre}}(\cdot), \text{dec}_{\text{gre}}(\cdot))$ , generalizes poorly, suffering from cross-entropy scaling as,

$$\min_{Q \in \mathcal{Q}_{1\text{-gram}}} \mathcal{L}(Q \circ \text{enc}_{\text{gre}}(\cdot)) \geq \frac{1 - o_{\delta}(1)}{3} H(\pi). \quad (41)$$

where the entropy of the stationary distribution of the source is  $H(\pi) = \frac{1}{2} \log(8)$  and the  $1 - o_{\delta}(1)$  term is  $(1-\delta)^2(1+\delta)^{-1}$ .

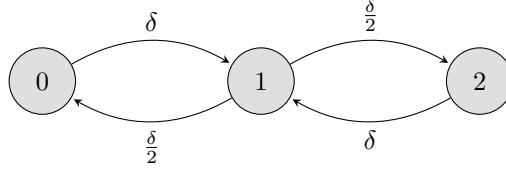


Figure 10: order-1 Markov source used in the proof of Theorem E.2

This means that the greedy encoder is not really compatible with the dictionary in the sense that the cross-entropy loss of the tokenizer is a constant multiple away from that achieved by the character-level tokenizer. The separation between eq. (40), and eq. (41) only manifests as  $\delta$  becomes smaller and smaller.

In this section, we prove that generalization of a dictionary is a function of the underlying tokenization algorithm used. In particular, the greedy encoder is not universal, and there exists dictionaries under the minimum-length encoder/decoder which achieve small cross-entropy loss, which do not generalize under the greedy encoder/decoder.

We split the proof of Theorem E.2 into two parts. We first define the stochastic source and dictionary we consider. Then we show that under the minimum-length encoder, the asymptotic cross-entropy loss is upper bounded by  $H_\infty$  up to a constant. Finally, we show that under the greedy-encoder, the same dictionary suffers from high cross-entropy loss, which is a constant factor away from that of the character encoder.

### E.1 Stochastic source and dictionary.

Consider an extension of the switching Markov source in example A.1 to  $\mathcal{A} = \{0, 1, 2\}$ . The Markov chain is described in Figure 10. The transition of the Markov chain is  $P(0|0) = P(1|1) = P(2|2) = 1 - \delta$ , and  $P(1|0) = P(2|1) = \delta$  and  $P(2|1) = P(0|1) = \delta/2$ , with the remaining transitions being 0-probability. For a parameter  $\ell > 0$  to be instantiated later, define  $S_1$  (resp.  $S_0, S_2$ ) as the set of all-1 (resp. all-0, all-2) strings of length  $\leq \ell - 1$ , including the empty string. Consider a dictionary composed of the following set of tokens,  $\{1s : s \in S_0 \cup S_1 \cup S_2\}$ . Therefore, the tokens follow the template  $10 \dots 0, 11 \dots 1$  or  $12 \dots 2$  and are of length at most  $\ell$ .  $\ell$  is chosen to be  $1 + 2 \log(1/\delta)/\delta$ .

Although we use the minimal encoder in the statement of Theorem E.2, for the purpose of analysis, define the following encoding algorithm: if the new string is prefixed by  $10 \dots 0$  or  $12 \dots 2$ , select the largest prefix which exists in dictionary and assign it as a token. If the new string starts with a sequence  $11 \dots 1$  of length  $x$ , consider the first  $\max\{\ell, x - 1\}$  length prefix and assign it as a token. Finally, if the string starts with 0 or 2, assign that character as token. Once the first token has been assigned, remove it and repeat.

### E.2 Minimal encoder achieves the optimal cross-entropy loss up to a constant.

First consider a simplification of the overall cross-entropy loss,

$$\begin{aligned} \min_{Q \in \mathcal{Q}_{1\text{-gram}}} \lim_{m \rightarrow \infty} \frac{1}{m} \mathcal{L}_m(Q \circ \text{enc}(\cdot)) \\ = \min_{Q \in \mathcal{Q}_{1\text{-gram}}} \lim_{m \rightarrow \infty} -\frac{1}{m} \mathbb{E} \left[ \log Q_{\#}(|\text{enc}_{\min}(s)|) + \sum_{t \in \text{Dict}} n_t \log Q_{\text{tok}}(t) \right] \end{aligned} \quad (42)$$

$$\leq \lim_{m \rightarrow \infty} \frac{1}{m} \mathbb{E} [\log(m) + |\text{enc}_{\min}(s)| \log |\text{Dict}|], \quad (43)$$

where in the last inequality we upper bound by choosing  $Q_{\#} = \text{Unif}([m])$  and  $Q_{\text{tok}}(t) = 1/|\text{Dict}|$ . Note that  $|\text{Dict}| \leq 2\ell + 1$  and letting  $\lim_{m \rightarrow \infty} \log(m)/m = 0$ ,

$$\begin{aligned} \min_{Q \in \mathcal{Q}_{1\text{-gram}}} \lim_{m \rightarrow \infty} \frac{1}{m} \mathcal{L}_m(Q \circ \text{enc}(\cdot)) &\leq \lim_{m \rightarrow \infty} \frac{1}{m} \mathbb{E} [|\text{enc}_{\min}(s)| \log(2\ell + 1)] \\ &\leq \lim_{m \rightarrow \infty} \frac{1}{m} \mathbb{E} [|\text{enc}(s)| \log(2\ell + 1)], \end{aligned} \quad (44)$$

where in (i), we replace  $|\text{enc}_{\min}(s)|$  by  $|\text{enc}(s)|$ , which is the encoder we define in Appendix E.1. By definition of the minimal encoder,  $|\text{enc}_{\min}(s)| \leq |\text{enc}(s)|$  surely. Recall that the encoder  $\text{enc}(\cdot)$  processes strings in a sequential (left-to-right) manner. In particular, by a similar argument as Lemma A.4, we can show that under this encoder, the limit  $n_t / \sum_{t'} n_{t'}$  almost surely converges to its expectation. More importantly, since,  $\sum_{t \in \text{Dict}} |t| n_t = m$ , we have that,

$$\lim_{m \rightarrow \infty} \frac{|\text{enc}(s)|}{m} \stackrel{\text{a.s.}}{=} \frac{1}{\mathbb{E}_{t \sim Q_{\text{MLE}}} [|t|]}.$$

converges to some limit almost surely. Therefore, from eq. (44),

$$\min_{Q \in \mathcal{Q}_{1\text{-gram}}} \lim_{m \rightarrow \infty} \frac{1}{m} \mathcal{L}_m(Q \circ \text{enc}(\cdot)) \leq \text{ess limsup}_{m \rightarrow \infty} \frac{|\text{enc}(s)|}{m} \log(2\ell + 1). \quad (45)$$

where the essential lim-sup captures the almost sure limit  $1/\mathbb{E}_{t \sim Q_{\text{MLE}}} [|t|]$ . The almost sure convergence of  $|\text{enc}(s)|/m$  also implies that we can let the limit  $m$  go to  $\infty$  in any manner, and the limit will remain the same. In particular, consider a process parameterized by  $i^*$  for generating the source string, such that surely  $m \geq i^*$ , where the total number of characters,  $m$ , is a random variable. As  $i^* \rightarrow \infty$ , we will also have  $m \rightarrow \infty$  surely, and so the limit of  $|\text{enc}(s)|/m$  under this modified stochastic process should also converge to the same limit.

Rather than sampling a string of a fixed length  $m$  from the source, consider the following sampling model: for  $i^* \rightarrow \infty$ , sample  $i^*$  geometric random variables  $X_1, \dots, X_{i^*} \stackrel{\text{i.i.d.}}{\sim} \text{Geo}(\delta)$  and construct the source string as the concatenation of  $i^*$  strings alternating between successive 1's and successive 0's or 2's (with the choice between the two made uniformly at random), with the  $i^{\text{th}}$  string of length  $X_i + 1$ . The overall number of characters sampled,  $m$ , is surely at least  $i^*$ .

Under this stochastic process, the size of the encoding of the string is upper bounded by,

$$|\text{enc}(s)| \leq |X_1 + 1| + \sum_{i=2}^{i^*} \left(1 + (X_i + 1 - \ell)_+\right)$$

This bound follows from the fact that in any substring  $s'$  of successive 1's followed by a substring  $s''$  of successive 0's or 2's, the encoder tokenizes the first  $\max\{\ell, |s'| - 1\}$  length prefix of  $s'$  as a token, and the remaining characters in  $s'$  into individual tokens except the last. Then, the last character of  $s'$  and the first  $\max\{\ell - 1, |s''|\}$  characters of  $s''$  are assigned as token. The remainder of  $s''$  is assigned as individual tokens. Each of  $s'$  or  $s''$  of length  $x$ , is allocated into at most  $1 + (x + 1 - \ell)_+$  tokens.

For any  $i$ ,  $\Pr(X_i \geq u) = (1 - \delta)^u$ , and therefore, summing over  $u \geq \ell$ , we get that  $\mathbb{E}[(X_i + 1 - \ell)_+] = \frac{(1 - \delta)^{\ell-1}}{\delta}$ . With  $\ell = 1 + 2 \log(1/\delta)/\delta$ , this expectation is upper bounded by  $\delta$ . Therefore,

$$\lim_{i^* \rightarrow \infty} \frac{\mathbb{E}[|\text{enc}(s)|]}{i^*} \leq \lim_{i^* \rightarrow \infty} \frac{1}{i^*} \mathbb{E} \left[ |X_1 + 1| + \sum_{i=2}^{i^*} \left(1 + (X_i + 1 - \ell)_+\right) \right] \leq 1 + \delta$$

More importantly, by the strong law of large numbers for a sum of independent random variables,  $(|X_1 + 1| + \sum_{i=2}^{i^*} (1 + (X_i + 1 - \ell)_+))/i^*$ , and therefore  $|\text{enc}(s)|/i^*$  is asymptotically almost surely upper bounded as,

$$\lim_{i^* \rightarrow \infty} \frac{|\text{enc}(s)|}{i^*} \stackrel{\text{a.s.}}{\leq} 1 + \delta, \quad (46)$$

On the other hand, the number of characters generated,  $m$ , equals  $\sum_{i=1}^{i^*} (X_i + 1)$ , and satisfies,  $\lim_{i^* \rightarrow \infty} \mathbb{E}[m]/i^* = 1 + \delta^{-1}$ . By another application of the strong law of large numbers for a sum of independent random variables,

$$\lim_{i^* \rightarrow \infty} \frac{m}{i^*} \stackrel{\text{a.s.}}{=} 1 + \delta^{-1}. \quad (47)$$

By combining eqs. (46) and (47), we have that,

$$\lim_{i^* \rightarrow \infty} \frac{|\text{enc}(s)|}{m} \stackrel{\text{a.s.}}{\leq} \frac{1 + \delta}{1 + \delta^{-1}} = \frac{1}{\delta}.$$

Finally, combining with eq. (45) and the ensuing discussion, we may upper bound the limiting cross-entropy loss by,

$$\min_{Q \in \mathcal{Q}_{1\text{-gram}}} \lim_{m \rightarrow \infty} \frac{1}{m} \mathcal{L}_m(Q \circ \text{enc}(\cdot)) \leq \delta \log(2\ell + 1) = \delta \log(3 + 4 \log(1/\delta)/\delta).$$

Note for this Markovian source, it is a short calculation to see that,

$$H_\infty = \mathbb{E}_{x \sim \pi} [H(P(\cdot|x))] = \delta \log(\sqrt{2}/\delta) + (1 - \delta) \log(1/(1 - \delta))$$

Note that for any  $\delta \leq 1/2$ , numerical evaluation gives the inequality,

$$1 \leq \frac{\delta \log(3 + 4 \log(1/\delta)/\delta)}{H_\infty} \leq 1.273$$

with the approximation factor improving as  $\delta$  becomes smaller. Therefore, this tokenizer achieves a normalized cross-entropy loss which asymptotically scales as a constant multiple of the entropy rate of the source.

### E.3 Greedy-encoder achieves poor cross-entropy loss

Note that the greedy encoder picks the largest prefix of the string which is a token, assigns and removes it, and iterates on the rest of the string. The greedy encoder's behavior is easy to analyze - every string of consecutive 1's in the new string is broken into chunks of length  $\ell$  (save potentially the last chunk) and each chunk is assigned as a token in  $\{1s : s \in S_1\} \subset \text{Dict}$ . If the length of this substring of successive 1's is not  $1, \ell + 1, 2\ell + 1, \dots$ , or in general,  $\equiv 1 \pmod{\ell}$ , every character in the next sequence, composed of 0's or 2's is tokenized into individual characters.

Similar to eq. (42) to eq. (43), consider a simplification of the overall cross-entropy loss,

$$\begin{aligned} & \min_{Q \in \mathcal{Q}_{1\text{-gram}}} \lim_{m \rightarrow \infty} \frac{1}{m} \mathcal{L}_m(Q \circ \text{enc}_{\text{gre}}(\cdot)) \\ &= \min_{Q \in \mathcal{Q}_{1\text{-gram}}} \lim_{m \rightarrow \infty} -\frac{1}{m} \mathbb{E} \left[ \log Q_{\#}(|\text{enc}_{\text{gre}}(s)|) + |\text{enc}_{\text{gre}}(s)| \sum_{t \in \text{Dict}} \frac{n_t}{|\text{enc}_{\text{gre}}(s)|} \log Q_{\text{tok}}(t) \right] \\ &\geq \min_{Q \in \mathcal{Q}_{1\text{-gram}}} \lim_{m \rightarrow \infty} -\frac{1}{m} \mathbb{E} \left[ |\text{enc}_{\text{gre}}(s)| \sum_{\substack{t \in \text{Dict} \\ Q_{\text{MLE}}(t) > 0}} Q_{\text{MLE}}(t) \log Q_{\text{tok}}(t) \right], \end{aligned}$$

where the last equation uses the fact that by Lemma A.4, for the greedy encoder,  $\lim_{m \rightarrow \infty} \frac{n_t}{|\text{enc}_{\text{gre}}(s)|} \stackrel{\text{a.s.}}{=} Q_{\text{MLE}}(t)$ . The minimizer of this objective subject to  $\sum_{t \in \text{Dict}: Q_{\text{MLE}}(t) > 0} Q_{\text{tok}}(t) \leq 1$  is  $Q_{\text{tok}}(t) = Q_{\text{MLE}}(t)$  resulting in the inequality,

$$\min_{Q \in \mathcal{Q}_{1\text{-gram}}} \lim_{m \rightarrow \infty} \frac{1}{m} \mathcal{L}_m(Q \circ \text{enc}_{\text{gre}}(\cdot)) \geq \lim_{m \rightarrow \infty} \frac{1}{m} \mathbb{E} [|\text{enc}_{\text{gre}}(s)| H(Q_{\text{MLE}})], \quad (48)$$

where we use the convention  $0 \log(1/0) \triangleq \lim_{P \rightarrow 0} P \log(1/P) = 0$  and therefore we may sum over tokens such that  $Q_{\text{MLE}}(t) = 0$  for free.

Considering the same geometric sampling model as in Appendix E.2, and Lemma A.4, we may study the almost sure limit  $Q_{\text{MLE}}(t) = \lim_{m \rightarrow \infty} n_t/|\text{enc}_{\text{gre}}(s)|$  by computing  $\lim_{i^* \rightarrow \infty} n_t/|\text{enc}_{\text{gre}}(s)|$  under the geometric sampling model since the almost sure limit exists. Recall that in the geometric sampling model, we generate the overall source string by concatenating  $i^*$  strings of length  $X_1 + 1, \dots, X_{i^*} + 1$  where  $X_i \sim \text{Geo}(\delta)$ , with the strings alternating between successive 1's and successive 0's or 2's (with the choice between the two made by the flip of a fair coin). For  $x \in \{0, 1, 2\}$ , let  $\mathcal{E}_i(x)$  denote the event that  $X_i$  is a string composed only of all  $x$ 's. The length of the greedy encoding of  $s$  is lower bounded by,

$$|\text{enc}_{\text{gre}}(s)| \geq \sum_{i=1}^{i^*} X_i \cdot \mathbb{I}(X_{i-1} \not\equiv 1 \pmod{\ell}) \mathbb{I}(\mathcal{E}_i(0) \cup \mathcal{E}_i(2)). \quad (49)$$

Which captures for the fact that all 0's and 2's are encoded into singular tokens unless the previous string of 1's was of length  $\equiv 1 \pmod{\ell}$ . By the law of large numbers of the RHS of eq. (49), the following a.s. lower bound is satisfied,

$$\lim_{i^* \rightarrow \infty} \frac{|\text{enc}_{\text{gre}}(\mathbf{s})|}{i^*} \stackrel{\text{a.s.}}{\geq} \frac{1}{2\delta} \left( 1 - \sum_{u=0}^{\infty} \delta(1-\delta)^{\ell u+1} \right) = \frac{1}{2\delta} \left( 1 - \frac{\delta(1-\delta)}{1-(1-\delta)^\ell} \right) \geq \frac{1-\delta}{2\delta}, \quad (50)$$

where the last inequality uses the fact that  $\ell = 1 + 2 \log(1/\delta)/\delta$ . Likewise, observe that,  $|\text{enc}_{\text{gre}}(\mathbf{s})| \leq m$  surely, and following the analysis in Appendix E.2 of eq. (47), we have that,

$$\lim_{i^* \rightarrow \infty} \frac{|\text{enc}_{\text{gre}}(\mathbf{s})|}{i^*} \leq \lim_{i^* \rightarrow \infty} \frac{m}{i^*} \stackrel{\text{a.s.}}{=} 1 + \delta^{-1}. \quad (51)$$

For  $x \in \{0, 2\}$ , observe that the expected number of times the token  $x$  is observed in the encoding of  $\mathbf{s}$ ,  $n_x$  can be written as,

$$n_x \geq \sum_{i=1}^{i^*} ((X_i + 1) \cdot \mathbb{I}(X_{i-1} \not\equiv 1 \pmod{\ell})) \mathbb{I}(\mathcal{E}_i(x)). \quad (52)$$

In particular, taking the expectation of eq. (52),

$$\mathbb{E}[n_x | \mathcal{E}_1(0) \cup \mathcal{E}_1(2)], \mathbb{E}[n_x | \mathcal{E}_1(1)] \geq \frac{i^* - 1}{4} (1 + \delta^{-1}) \left( 1 - \sum_{u=0}^{\infty} \delta(1-\delta)^{\ell u+1} \right) \geq \frac{i^* - 1}{4} \cdot \frac{1 - \delta^2}{\delta}. \quad (53)$$

Note that in any realization of the geometric sampling process, in eq. (52), either the odd indexed substrings are all-1's or the even indexed substrings are all-1's. Therefore, surely, all the non-zero terms in the above summation are of the same parity. Moreover, since the  $i^{\text{th}}$  term in the sum only depends on  $X_i$  and  $X_{i-1}$ , conditioned on whether the non-zero parities are even or odd,  $n_x$  can be written as a sum of  $\approx i^*/2$  mutually independent terms. By the strong law of large numbers on each of the conditional processes, eqs. (52) and (53) implies that for  $x \in \{0, 2\}$ ,

$$\lim_{i^* \rightarrow \infty} \frac{n_x}{i^*} \stackrel{\text{a.s.}}{\geq} \frac{1 - \delta^2}{4\delta}.$$

To upper bound  $n_x$ , note that it is upper bounded by the number of times the character  $x$  appears in the source string, which by the strong law of large numbers a.s (after normalizing by  $i^*$ ), scales as  $1/4\delta$ . Finally, to bound  $Q_{\text{MLE}}(\mathbf{t})$  which is the sequential nature of the encoder, using a similar proof as Lemma A.4, we can show that  $n_{\mathbf{t}} / \sum_{\mathbf{t}'} n_{\mathbf{t}'}$  converges to the unigram MLE model for this tokenizer. For the token  $x \in \{0, 2\}$ ,

$$\lim_{i^* \rightarrow \infty} \frac{n_x}{|\text{enc}(\mathbf{s})|} = Q_{\text{MLE}}(x) \leq \mathbb{E} \left[ \lim_{i^* \rightarrow \infty} \frac{n_x}{n_2 + n_0} \right] \quad (54)$$

Using the a.s. upper and lower bounds on  $|\text{enc}(\mathbf{s})|$ ,  $n_0$  and  $n_2$  derived in eqs. (51) and (54), we arrive at lower and upper bounds on  $Q_{\text{MLE}}(x)$  for  $x \in \{0, 2\}$ ,

$$\frac{1}{4} \approx \frac{1-\delta}{4} = \frac{(1-\delta^2)}{4\delta(1+\delta^{-1})} \leq Q_{\text{MLE}}(x) \leq \frac{1}{2(1-\delta^2)} \approx \frac{1}{2}.$$

Since there are at least two tokens having probability bounded away from 0 and 1 by a constant under the MLE unigram model, the entropy of  $Q_{\text{MLE}}$  must also be lower bounded by a constant. Indeed,

$$H(Q_{\text{MLE}}) \geq 2 \min_{\frac{1-\delta}{4} \leq y \leq \frac{1}{2(1-\delta^2)}} y \log(1/y).$$

It is easy to verify that for  $\delta \leq 0.5$ , the minimizer is achieved at  $y = \frac{1-\delta}{4}$ , which leads to the lower bound,

$$H(Q_{\text{MLE}}) \geq \left( \frac{1-\delta}{2} \right) \log \left( \frac{4}{1-\delta} \right)$$



Architecture	GPT-2
Batch size	Grid-searched in {8, 16, 32}
Gradient acc. steps	1
Tokenizer dictionary size	{10, 20}
Tokenizer dataset size	10,000
Optimizer	AdamW ( $\beta_1 = 0.9, \beta_2 = 0.95$ )
Learning rate	0.002
Scheduler	Cosine
# Iterations	8000
Weight decay	$1 \times 10^{-3}$
Dropout	0
Sequence length	512
Embedding dimension	Grid-searched in {10, 20, 30, 40}
# layers	Grid-searched in {1, 2, 4, 8}
# heads	Grid-searched in {1, 2, 4, 8, 16}
Repetitions	5

Table 3: Hyperparameter choices

Finally, combining this lower bound on  $H(Q_{\text{MLE}})$  with eq. (48), we have that,

$$\begin{aligned}
\min_{Q \in \mathcal{Q}_{1\text{-gram}}} \lim_{m \rightarrow \infty} \frac{1}{m} \mathcal{L}_m(Q \circ \text{enc}(\cdot)) &= \lim_{i^* \rightarrow \infty} \mathbb{E} \left[ \frac{|\text{enc}_{\text{gre}}(s)|}{m} H(Q_{\text{MLE}}) \right] \\
&\geq \lim_{i^* \rightarrow \infty} \mathbb{E} \left[ \frac{|\text{enc}_{\text{gre}}(s)|}{m} \right] \cdot \left( \frac{1-\delta}{2} \right) \log \left( \frac{4}{1-\delta} \right) \\
&\stackrel{(i)}{\geq} \frac{1-\delta}{2\delta(1+\delta^{-1})} \cdot \left( \frac{1-\delta}{2} \right) \log \left( \frac{4}{1-\delta} \right) \\
&\geq \frac{(1-\delta)^2}{3(1+\delta)} H(\pi)
\end{aligned}$$

where (i) follows from the lower bound on  $|\text{enc}_{\text{gre}}(s)|$  in eq. (50) with the almost sure limit of  $m$  in eq. (47) and noting that  $|\text{enc}_{\text{gre}}(s)|/m \leq 1$  surely. The last inequality follows by simplifying using  $\pi = (1/4, 1/2, 1/4)$  and  $H(\pi) = \frac{1}{2} \log(8)$ .

## F Experiment details

**Experiment 1 (Figures 4a and 4b).** In this and previous experiments (Figures 2, 3a and 3b), we train the transformers on a single GPU on an  $8 \times$  A100 node. The wall-clock time measured does not count time spent in validation loss evaluations. The hyperparameter choices are listed in Table 3.

**Experiment 2 (Table 1).** We evaluate pre-trained tokenizers on various datasets. In this experiment, we do not evaluate the likelihood model on test sequences, rather, we estimate the cross-entropy of the best unigram model by using the approximation,

$$-\mathbb{E} \left[ \sum_{t \in \text{Dict}} n_t \log Q_{\text{MLE}}(t) \right] \approx - \sum_{t \in \text{Dict}} \hat{n}_t \log(\hat{Q}(t)) \quad (55)$$

where  $\hat{Q}(t) = \frac{\hat{n}_t}{\sum_t \hat{n}_t}$  is the MLE unigram model learnt from a finite dataset, which we choose here as GLUE (Wang et al., 2019), and  $\hat{n}_t$  is the number of times the token  $t$  is observed in the encoding of the dataset. This approximation allows us to separate the error stemming from learning a suboptimal likelihood model which tends to have higher sample complexity requirements and focus on the asymptotic error of the tokenizer.

We use Monte-carlo sampling to approximate the cross-entropy loss estimator in eq. (55). These approximations tends to underestimate the true cross-entropy loss due to the concavity of  $x \log(1/x)$  close to 0. In general, the gap between the approximation and the true error is expected to grow with  $k$ . Therefore, the true difference between the estimate of the best unigram model on a tokenizer and the best  $k$ -gram model for  $k \geq 2$  on the character level tokenizer is likely to be larger than the reported figures.

**Experiment 3 (Figure 5).** We train the LZW, BPE, Unigram and Wordpiece tokenizers with dictionary sizes  $\{5000, 6000, 8000, 12000, 20000, 32000, 50000, 80000\}$ . The cross-entropy loss incurred by the best 1-gram model is estimated using eq. (55) while for  $k$ -gram models for  $k \geq 2$ , we use Monte-carlo sampling to estimate the cross-entropy of the empirical  $k$ -gram model computed using the GLUE dataset. For the  $k$ -gram models trained on the character level tokenizer, since the vocabulary size is fixed, we instead plot the number of distinct  $k$ -grams on the  $x$ -axis. While this is not a true measure of the number of parameters in the underlying  $k$ -gram model, we use this as a proxy for the same.

## G NeurIPS Paper Checklist

The checklist is designed to encourage best practices for responsible machine learning research, addressing issues of reproducibility, transparency, research ethics, and societal impact. Do not remove the checklist: **The papers not including the checklist will be desk rejected.** The checklist should follow the references and precede the (optional) supplemental material. The checklist does NOT count towards the page limit.

Please read the checklist guidelines carefully for information on how to answer these questions. For each question in the checklist:

- You should answer [Yes], [No], or [NA].
- [NA] means either that the question is Not Applicable for that particular paper or the relevant information is Not Available.
- Please provide a short (1–2 sentence) justification right after your answer (even for NA).

### 1. Claims

Question: Do the main claims made in the abstract and introduction accurately reflect the paper's contributions and scope?

Answer: [Yes]

Justification: The paper lists an empirical phenomenon (justified in Fig. 2) and theoretical contributions justified in Theorems 3.1, 3.3 and 3.5

Guidelines:

- The answer NA means that the abstract and introduction do not include the claims made in the paper.
- The abstract and/or introduction should clearly state the claims made, including the contributions made in the paper and important assumptions and limitations. A No or NA answer to this question will not be perceived well by the reviewers.
- The claims made should match theoretical and experimental results, and reflect how much the results can be expected to generalize to other settings.
- It is fine to include aspirational goals as motivation as long as it is clear that these goals are not attained by the paper.

### 2. Limitations

Question: Does the paper discuss the limitations of the work performed by the authors?

Answer: [Yes]

Justification: Remark 3.3

Guidelines:

- The answer NA means that the paper has no limitation while the answer No means that the paper has limitations, but those are not discussed in the paper.
- The authors are encouraged to create a separate "Limitations" section in their paper.
- The paper should point out any strong assumptions and how robust the results are to violations of these assumptions (e.g., independence assumptions, noiseless settings, model well-specification, asymptotic approximations only holding locally). The authors should reflect on how these assumptions might be violated in practice and what the implications would be.
- The authors should reflect on the scope of the claims made, e.g., if the approach was only tested on a few datasets or with a few runs. In general, empirical results often depend on implicit assumptions, which should be articulated.
- The authors should reflect on the factors that influence the performance of the approach. For example, a facial recognition algorithm may perform poorly when image resolution is low or images are taken in low lighting. Or a speech-to-text system might not be used reliably to provide closed captions for online lectures because it fails to handle technical jargon.

- The authors should discuss the computational efficiency of the proposed algorithms and how they scale with dataset size.
- If applicable, the authors should discuss possible limitations of their approach to address problems of privacy and fairness.
- While the authors might fear that complete honesty about limitations might be used by reviewers as grounds for rejection, a worse outcome might be that reviewers discover limitations that aren't acknowledged in the paper. The authors should use their best judgment and recognize that individual actions in favor of transparency play an important role in developing norms that preserve the integrity of the community. Reviewers will be specifically instructed to not penalize honesty concerning limitations.

### 3. Theory Assumptions and Proofs

Question: For each theoretical result, does the paper provide the full set of assumptions and a complete (and correct) proof?

Answer: [\[Yes\]](#)

Justification: Assumption 3.2

Guidelines:

- The answer NA means that the paper does not include theoretical results.
- All the theorems, formulas, and proofs in the paper should be numbered and cross-referenced.
- All assumptions should be clearly stated or referenced in the statement of any theorems.
- The proofs can either appear in the main paper or the supplemental material, but if they appear in the supplemental material, the authors are encouraged to provide a short proof sketch to provide intuition.
- Inversely, any informal proof provided in the core of the paper should be complemented by formal proofs provided in appendix or supplemental material.
- Theorems and Lemmas that the proof relies upon should be properly referenced.

### 4. Experimental Result Reproducibility

Question: Does the paper fully disclose all the information needed to reproduce the main experimental results of the paper to the extent that it affects the main claims and/or conclusions of the paper (regardless of whether the code and data are provided or not)?

Answer: [\[Yes\]](#)

Justification: The code has been released along with the rest of the submission.

Guidelines:

- The answer NA means that the paper does not include experiments.
- If the paper includes experiments, a No answer to this question will not be perceived well by the reviewers: Making the paper reproducible is important, regardless of whether the code and data are provided or not.
- If the contribution is a dataset and/or model, the authors should describe the steps taken to make their results reproducible or verifiable.
- Depending on the contribution, reproducibility can be accomplished in various ways. For example, if the contribution is a novel architecture, describing the architecture fully might suffice, or if the contribution is a specific model and empirical evaluation, it may be necessary to either make it possible for others to replicate the model with the same dataset, or provide access to the model. In general, releasing code and data is often one good way to accomplish this, but reproducibility can also be provided via detailed instructions for how to replicate the results, access to a hosted model (e.g., in the case of a large language model), releasing of a model checkpoint, or other means that are appropriate to the research performed.
- While NeurIPS does not require releasing code, the conference does require all submissions to provide some reasonable avenue for reproducibility, which may depend on the nature of the contribution. For example
  - (a) If the contribution is primarily a new algorithm, the paper should make it clear how to reproduce that algorithm.

- (b) If the contribution is primarily a new model architecture, the paper should describe the architecture clearly and fully.
- (c) If the contribution is a new model (e.g., a large language model), then there should either be a way to access this model for reproducing the results or a way to reproduce the model (e.g., with an open-source dataset or instructions for how to construct the dataset).
- (d) We recognize that reproducibility may be tricky in some cases, in which case authors are welcome to describe the particular way they provide for reproducibility. In the case of closed-source models, it may be that access to the model is limited in some way (e.g., to registered users), but it should be possible for other researchers to have some path to reproducing or verifying the results.

## 5. Open access to data and code

Question: Does the paper provide open access to the data and code, with sufficient instructions to faithfully reproduce the main experimental results, as described in supplemental material?

Answer: [Yes]

Justification: Instructions provided in the jupyter notebook.

Guidelines:

- The answer NA means that paper does not include experiments requiring code.
- Please see the NeurIPS code and data submission guidelines (<https://nips.cc/public/guides/CodeSubmissionPolicy>) for more details.
- While we encourage the release of code and data, we understand that this might not be possible, so “No” is an acceptable answer. Papers cannot be rejected simply for not including code, unless this is central to the contribution (e.g., for a new open-source benchmark).
- The instructions should contain the exact command and environment needed to run to reproduce the results. See the NeurIPS code and data submission guidelines (<https://nips.cc/public/guides/CodeSubmissionPolicy>) for more details.
- The authors should provide instructions on data access and preparation, including how to access the raw data, preprocessed data, intermediate data, and generated data, etc.
- The authors should provide scripts to reproduce all experimental results for the new proposed method and baselines. If only a subset of experiments are reproducible, they should state which ones are omitted from the script and why.
- At submission time, to preserve anonymity, the authors should release anonymized versions (if applicable).
- Providing as much information as possible in supplemental material (appended to the paper) is recommended, but including URLs to data and code is permitted.

## 6. Experimental Setting/Details

Question: Does the paper specify all the training and test details (e.g., data splits, hyperparameters, how they were chosen, type of optimizer, etc.) necessary to understand the results?

Answer: [Yes]

Justification: Table 3

Guidelines:

- The answer NA means that the paper does not include experiments.
- The experimental setting should be presented in the core of the paper to a level of detail that is necessary to appreciate the results and make sense of them.
- The full details can be provided either with the code, in appendix, or as supplemental material.

## 7. Experiment Statistical Significance

Question: Does the paper report error bars suitably and correctly defined or other appropriate information about the statistical significance of the experiments?

Answer: [Yes]

Justification: All plots which allow for it, contain standard error bars.

Guidelines:

- The answer NA means that the paper does not include experiments.
- The authors should answer "Yes" if the results are accompanied by error bars, confidence intervals, or statistical significance tests, at least for the experiments that support the main claims of the paper.
- The factors of variability that the error bars are capturing should be clearly stated (for example, train/test split, initialization, random drawing of some parameter, or overall run with given experimental conditions).
- The method for calculating the error bars should be explained (closed form formula, call to a library function, bootstrap, etc.)
- The assumptions made should be given (e.g., Normally distributed errors).
- It should be clear whether the error bar is the standard deviation or the standard error of the mean.
- It is OK to report 1-sigma error bars, but one should state it. The authors should preferably report a 2-sigma error bar than state that they have a 96% CI, if the hypothesis of Normality of errors is not verified.
- For asymmetric distributions, the authors should be careful not to show in tables or figures symmetric error bars that would yield results that are out of range (e.g. negative error rates).
- If error bars are reported in tables or plots, The authors should explain in the text how they were calculated and reference the corresponding figures or tables in the text.

## 8. Experiments Compute Resources

Question: For each experiment, does the paper provide sufficient information on the computer resources (type of compute workers, memory, time of execution) needed to reproduce the experiments?

Answer: [Yes]

Justification: Appendix F contains this information.

Guidelines:

- The answer NA means that the paper does not include experiments.
- The paper should indicate the type of compute workers CPU or GPU, internal cluster, or cloud provider, including relevant memory and storage.
- The paper should provide the amount of compute required for each of the individual experimental runs as well as estimate the total compute.
- The paper should disclose whether the full research project required more compute than the experiments reported in the paper (e.g., preliminary or failed experiments that didn't make it into the paper).

## 9. Code Of Ethics

Question: Does the research conducted in the paper conform, in every respect, with the NeurIPS Code of Ethics <https://neurips.cc/public/EthicsGuidelines?>

Answer: [Yes]

Justification: No NeurIPS code of ethics were violated.

Guidelines:

- The answer NA means that the authors have not reviewed the NeurIPS Code of Ethics.
- If the authors answer No, they should explain the special circumstances that require a deviation from the Code of Ethics.
- The authors should make sure to preserve anonymity (e.g., if there is a special consideration due to laws or regulations in their jurisdiction).

## 10. Broader Impacts

Question: Does the paper discuss both potential positive societal impacts and negative societal impacts of the work performed?

Answer: [NA]

Justification: This work is a primarily theoretical study on the behavior of tokenization on toy problems (learning Markov chains). The societal impact of this research is not likely to be significant.

Guidelines:

- The answer NA means that there is no societal impact of the work performed.
- If the authors answer NA or No, they should explain why their work has no societal impact or why the paper does not address societal impact.
- Examples of negative societal impacts include potential malicious or unintended uses (e.g., disinformation, generating fake profiles, surveillance), fairness considerations (e.g., deployment of technologies that could make decisions that unfairly impact specific groups), privacy considerations, and security considerations.
- The conference expects that many papers will be foundational research and not tied to particular applications, let alone deployments. However, if there is a direct path to any negative applications, the authors should point it out. For example, it is legitimate to point out that an improvement in the quality of generative models could be used to generate deepfakes for disinformation. On the other hand, it is not needed to point out that a generic algorithm for optimizing neural networks could enable people to train models that generate Deepfakes faster.
- The authors should consider possible harms that could arise when the technology is being used as intended and functioning correctly, harms that could arise when the technology is being used as intended but gives incorrect results, and harms following from (intentional or unintentional) misuse of the technology.
- If there are negative societal impacts, the authors could also discuss possible mitigation strategies (e.g., gated release of models, providing defenses in addition to attacks, mechanisms for monitoring misuse, mechanisms to monitor how a system learns from feedback over time, improving the efficiency and accessibility of ML).

## 11. Safeguards

Question: Does the paper describe safeguards that have been put in place for responsible release of data or models that have a high risk for misuse (e.g., pretrained language models, image generators, or scraped datasets)?

Answer: [NA]

Justification: No models with a high risk for misuse were trained or released.

Guidelines:

- The answer NA means that the paper poses no such risks.
- Released models that have a high risk for misuse or dual-use should be released with necessary safeguards to allow for controlled use of the model, for example by requiring that users adhere to usage guidelines or restrictions to access the model or implementing safety filters.
- Datasets that have been scraped from the Internet could pose safety risks. The authors should describe how they avoided releasing unsafe images.
- We recognize that providing effective safeguards is challenging, and many papers do not require this, but we encourage authors to take this into account and make a best faith effort.

## 12. Licenses for existing assets

Question: Are the creators or original owners of assets (e.g., code, data, models), used in the paper, properly credited and are the license and terms of use explicitly mentioned and properly respected?

Answer: [Yes]

Justification: Code has been properly credited, via citing the relevant paper.

Guidelines:

- The answer NA means that the paper does not use existing assets.
- The authors should cite the original paper that produced the code package or dataset.
- The authors should state which version of the asset is used and, if possible, include a URL.
- The name of the license (e.g., CC-BY 4.0) should be included for each asset.
- For scraped data from a particular source (e.g., website), the copyright and terms of service of that source should be provided.
- If assets are released, the license, copyright information, and terms of use in the package should be provided. For popular datasets, [paperswithcode.com/datasets](https://paperswithcode.com/datasets) has curated licenses for some datasets. Their licensing guide can help determine the license of a dataset.
- For existing datasets that are re-packaged, both the original license and the license of the derived asset (if it has changed) should be provided.
- If this information is not available online, the authors are encouraged to reach out to the asset's creators.

### 13. New Assets

Question: Are new assets introduced in the paper well documented and is the documentation provided alongside the assets?

Answer: [NA]

Justification: No new assets released.

Guidelines:

- The answer NA means that the paper does not release new assets.
- Researchers should communicate the details of the dataset/code/model as part of their submissions via structured templates. This includes details about training, license, limitations, etc.
- The paper should discuss whether and how consent was obtained from people whose asset is used.
- At submission time, remember to anonymize your assets (if applicable). You can either create an anonymized URL or include an anonymized zip file.

### 14. Crowdsourcing and Research with Human Subjects

Question: For crowdsourcing experiments and research with human subjects, does the paper include the full text of instructions given to participants and screenshots, if applicable, as well as details about compensation (if any)?

Answer: [NA]

Justification: No crowdsourcing or research with human subjects.

Guidelines:

- The answer NA means that the paper does not involve crowdsourcing nor research with human subjects.
- Including this information in the supplemental material is fine, but if the main contribution of the paper involves human subjects, then as much detail as possible should be included in the main paper.
- According to the NeurIPS Code of Ethics, workers involved in data collection, curation, or other labor should be paid at least the minimum wage in the country of the data collector.

### 15. Institutional Review Board (IRB) Approvals or Equivalent for Research with Human Subjects

Question: Does the paper describe potential risks incurred by study participants, whether such risks were disclosed to the subjects, and whether Institutional Review Board (IRB) approvals (or an equivalent approval/review based on the requirements of your country or institution) were obtained?

Answer: [NA]

Justification: No crowdsourcing or research with human subjects.



Guidelines:

- The answer NA means that the paper does not involve crowdsourcing nor research with human subjects.
- Depending on the country in which research is conducted, IRB approval (or equivalent) may be required for any human subjects research. If you obtained IRB approval, you should clearly state this in the paper.
- We recognize that the procedures for this may vary significantly between institutions and locations, and we expect authors to adhere to the NeurIPS Code of Ethics and the guidelines for their institution.
- For initial submissions, do not include any information that would break anonymity (if applicable), such as the institution conducting the review.