

# Empirical Evidence that Modern Byte Pair Tokenizers are Zipfian

Jack Hanke      Daniel Plotkin  
Nicole Birova    David Demeter  
Northwestern University

## Abstract

A majority of large language models ingest word fragments produced by an algorithm known as byte pair encoding. This algorithm groups high-frequency element pairings in natural language into individual units called *tokens*. A natural question is whether the frequency distribution of the tokens produced by this pairing deviate significantly from the source language’s. Zipf showed that many natural language’s frequency distribution follow a power law, commonly known as Zipf’s Law. We examine two modern tokenizer’s adherence to Zipf’s law at the token level. We provide empirical evidence that these tokenizers are Zipfian on two corpuses, and speculate as to why this is. Additionally, we give evidence that Zipfness is preserved over many individual steps of byte pair encoding.

## 1 Introduction

George Zipf, in *The Psycho-Biology of Language* (Zipf, 1935) introduced the trend that the number of occurrences of a word, the word’s *frequency*, exhibits a power law. Specifically, Zipf said the word frequency in a corpus is proportional to the word’s frequency rank, that being

$$\text{word frequency} \propto \frac{1}{\text{word rank}}. \quad (1)$$

For example, Zipf’s law predicts that the second most frequent word will be half as common as the most frequent word. A more descriptive model of word frequency that is more commonly referenced in linguistics literature is the Zipf-Mandelbrot distribution, written

$$\text{word frequency} \propto \frac{1}{(\text{word rank} + b)^a}, \quad (2)$$

where  $a, b$  are fitted parameters. We say a distribution that follows the trend in Equation 2 with  $a \approx 1$  is *Zipf distributed*, or simply *Zipfian*.

The accuracy of the trend in Equation 2 has been examined in  $10^8$  English words in (Ferrer-i Cancho and Solé, 2000), over 50 languages in (Yu et al., 2018), and written-versus-spoken corpuses in (Lin et al., 2015). Each of these studies demonstrate that Zipf’s law is generally exhibited for common and somewhat-uncommon words, but rare words (words with high rank) appear less frequently than predicted. This deviation is shown to be statistically significant, and appears as two linear trends in log-log rank frequency plots. In English text, this transition is found around the  $10^4$ -th ranked word (Ferrer-i Cancho and Solé, 2000) (Yu et al., 2018). The authors further explore the linguistic relevance and universality of these multiple trends.

Nearly a century after Zipf’s discovery, large language models (LLMs) generate text comparable to human communication (?). LLMs digest text using a fixed vocabulary of word fragments called *tokens*. The mapping between natural language and tokens is most commonly computed using the *byte pair encoding algorithm* (BPE). BPE was introduced by Gage (Gage, 1994) for file compression, and later adapted for use in natural language processing by (Sennrich et al., 2016). The NLP version of BPE, instead of seeking to optimally compress a corpus of text into a vocabulary of any size, instead maps a corpus to a small fixed vocabulary. In this paper, we consider the natural language processing variant of the algorithm.

BPE also requires specifying the *elements* of a corpus, which are the portion of the text that is considered a indivisible chunk of the text. Elements can be words, phonemes, or letters. Modern tokenizers use the bytes of the Unicode representation as the elements, and so are called *byte-level* BPE tokenizers (Radford et al., 2019).

When considering these ideas in tandem, a natural question arises: given a corpus that appears Zipfian, how does the BPE tokenization process affect this trend?

## 2 Methods

To explore this question, we examine two modern byte-level BPE tokenizers. We choose the tokenizer for the RoBERTa language model (Liu et al., 2019), which has a vocabulary size of 50,265. As the training data is publicly known, this allows us to conduct frequency analyses on corpora that are on and off-distribution for the tokenizer. We also compare our results with the 200,019 vocab GPT-4o tokenizer (OpenAI et al., 2024), as a representative of industry-scale byte pair encoding.

We examine the 4.4GB bookscorpus dataset (Zhu et al., 2015), which is within RoBERTa’s training data, and the 5.6GB MiniPile dataset (Kaddour, 2023), which is not.

We compute the word and token frequency for both tokenizers on each corpus, where we define a word as anything separated by the space character (U+0020). We then compute the data’s  $\chi^2$  goodness-of-fit statistic for the fitted Zipf distribution, and compare the fit to exponential and log-normal distributions. We do not remove any words or tokens from consideration unless otherwise stated.

Additionally, we create a synthetic dataset of  $10^5$  samples from a zipf-distributed alphabet of size 100 to examine how the distribution deviates from Zipf over successive steps of BPE.

## 3 Results

For bookscorpus, we find 1.3M distinct words, 31,729 distinct RoBERTa tokens, and 12,859 distinct GPT-4o tokens. For MiniPile, we find 32M distinct words, 50,165 distinct RoBERTa tokens, and 178,416 distinct GPT-4o tokens. We plot the log-log rank frequency for both corpora in Figure 1.

We also report the least and most common words and tokens for each tokenizer and corpus. We report these words and tokens both including and excluding control tokens and punctuation, which are summarized in Table 1 and Table 2.

We report the  $\chi^2$  test statistics for fitted Zipf-Mandelbrot in Table 3, where we use Equation 2 with  $a = 1.1$  and  $b = 0$  for all Zipf tests.

Finally, we compute the mean squared error (MSE) between the Zipfian prediction and successive applications of a single step of BPE over the Zipfian synthetic dataset in Figure 2.

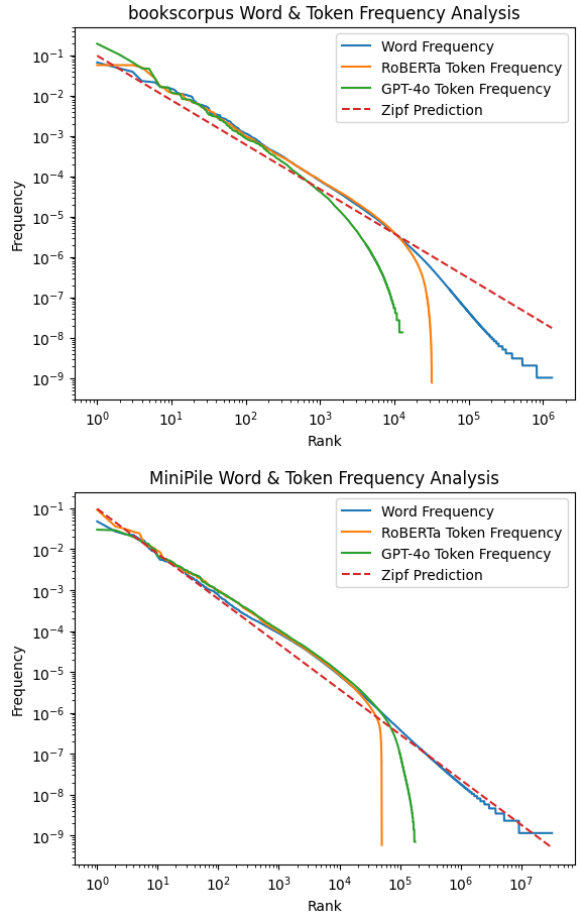


Figure 1: Log-log plots of word & token rank vs word & token frequency

## 4 Analysis

In general, we find strong evidence that BPE tokenizers are Zipfian. Deviations from the Zipfian trend only appear for the final few RoBERTa tokens for both corpora and the GPT-4o tokens for the MiniPile corpus. This includes in and out-of-distribution text. However, the GPT-4o tokenizer trend on the bookscorpus seems to deviate more dramatically, similar to the multiple regime studies mentioned in Section 1. We speculate that this may be due to MiniPile being a multilingual dataset, which may affect the distribution for a multilingual tokenizer such as GPT-4o. The  $\chi^2$  goodness-of-fit tests confirm a Zipf describes the data well, all having p-values  $< 0.0001$ .

Over the token vocabulary, the most common tokens tend to be punctuation, newline characters, and control tokens. Excluding punctuation and special tokens, we recover the known most common English words.

Additionally, applying the BPE pairing process

Rank	Word		RoBERTa		GPT-4o	
1	.	the	<s>	▯	“	,
2	,	of	\n	\n	i	▯the
3	the	and	<\s>	.	he	.
-1	restrain	RootDir,	seq	▯Archdemon	wares	aryny
-1	liarliar	homocystinemia	okemon	▯petertodd	slan	verlening
-1	shop-that	halfday	ython	▯councill	wier	▯myx

Table 1: Summary of most and least common words & tokens for each corpus, where the left subcolumn is bookcorpus and the right subcolumn is MiniPile, *including* control tokens and punctuation. A rank of -1 indicates a word or token that appears only once in the given corpus. The space character (Unicode U+0020) is denoted by the ‘▯’ character.

Rank	Word		RoBERTa		GPT-4o	
1	the	the	the	the	i	the
2	to	of	to	of	he	of
3	i	and	and	and	she	and

Table 2: Summary of most common words & tokens for each corpus, where the left subcolumn is bookcorpus and the right subcolumn is MiniPile, *excluding* control tokens and punctuation. All of these words are considered *function words* in linguistics literature.

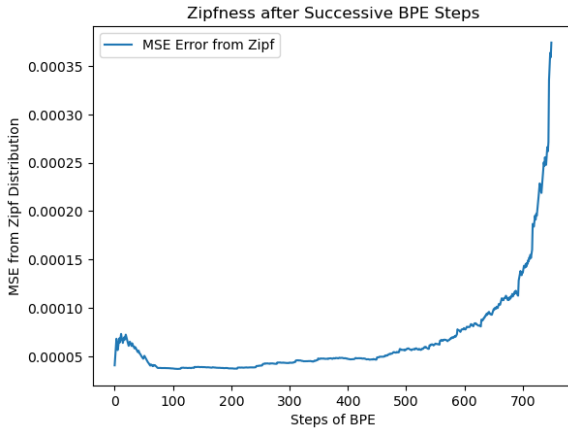


Figure 2: Mean Squared Error from Zipf prediction of BPE applied to  $10^5$  samples of a synthetic Zipf-distributed corpus of 100 elements

to the synthetic data shows almost no deviations from the original Zipf distribution for hundreds of iterations. However, after about 700 iterations the MSE from Zipf does begin to increase, though the error is still relatively small.

## 5 Conclusions & Future Work

The byte pair encoding algorithm, despite building vocabularies using the most frequent pairings of elements, generates vocabularies that are Zipfian, up to the rarest few tokens. We demonstrate this

over multiple tokenizers and corpuses, including both in and out-of-distribution text. We also provide empirical evidence that Zipfness is preserved over many iterations of BPE.

Noticeable theoretical work has been done to explain Zipf’s law for language (Li, 1992). Most significant was Belevitch’s *On the statistical laws of linguistic distributions* (Belevitch, 1959), in which the author shows the first order approximation of the rank ordering of many statistical distributions are all Zipfian. This indicates that Zipf’s law may be due to the rank ordering of words more than the underlying formation of language. Can Belevitch’s proof be shown to be invariant for some number of byte pair encoding transformations?

## References

- V. Belevitch. 1959. On the statistical laws of linguistic distributions. *Annales de la Société Scientifique de Bruxelles*, 73:310–326.
- Ramon Ferrer-i Cancho and Ricard Solé. 2000. Two regimes in the frequency of words and the origins of complex lexicons: Zipf’s law revisited. *Santa Fe Institute, Working Papers*.
- Philip Gage. 1994. A new algorithm for data compression. *C Users J.*, 12(2):23–38.
- Jean Kaddour. 2023. [The minipile challenge for data-efficient language models](#). *Preprint*, arXiv:2304.08442.
- W. Li. 1992. [Random texts exhibit zipf’s-law-like word frequency distribution](#). *IEEE Transactions on Information Theory*, 38(6):1842–1845.
- Ruokuang Lin, Qianli D. Y. Ma, and Chunhua Bian. 2015. [Scaling laws in human speech, decreasing emergence of new words and a generalized model](#). *Preprint*, arXiv:1412.4846.
- Yinhan Liu, Myle Ott, Naman Goyal, Jingfei Du, Mandar Joshi, Danqi Chen, Omer Levy, Mike Lewis,

Distribution	Words	RoBERTa	GPT-4o
bookscorpus	$2.5 \cdot 10^9$	$1.6 \cdot 10^9$	$1.4 \cdot 10^8$
MiniPile	$1.4 \cdot 10^9$	$4.1 \cdot 10^9$	$1.8 \cdot 10^9$

Table 3:  $\chi^2$  test statistics for distribution fit for each tokenizer and corpus, each of which have a p value of near 0

Luke Zettlemoyer, and Veselin Stoyanov. 2019. [Roberta: A robustly optimized bert pretraining approach](#). *Preprint*, arXiv:1907.11692.

OpenAI, Josh Achiam, Steven Adler, Sandhini Agarwal, Lama Ahmad, Ilge Akkaya, Florencia Leoni Aleman, Diogo Almeida, Janko Altenschmidt, Sam Altman, Shyamal Anadkat, Red Avila, Igor Babuschkin, Suchir Balaji, Valerie Balcom, Paul Baltescu, Haiming Bao, Mohammad Bavarian, Jeff Belgum, and 262 others. 2024. [Gpt-4 technical report](#). *Preprint*, arXiv:2303.08774.

Alec Radford, Jeffrey Wu, Rewon Child, David Luan, Dario Amodei, and Ilya Sutskever. 2019. [Language models are unsupervised multitask learners](#). *OpenAI*. Accessed: 2024-11-15.

Rico Sennrich, Barry Haddow, and Alexandra Birch. 2016. [Neural machine translation of rare words with subword units](#). *Preprint*, arXiv:1508.07909.

Shuiyuan Yu, Chunshan Xu, and Haitao Liu. 2018. [Zipf’s law in 50 languages: its structural pattern, linguistic interpretation, and cognitive motivation](#). *Preprint*, arXiv:1807.01855.

Yukun Zhu, Ryan Kiros, Richard Zemel, Ruslan Salakhutdinov, Raquel Urtasun, Antonio Torralba, and Sanja Fidler. 2015. [Aligning books and movies: Towards story-like visual explanations by watching movies and reading books](#). *Preprint*, arXiv:1506.06724.

G. K. Zipf. 1935. *The Psycho-Biology of Language*. Houghton Mifflin, Boston.