

Dotless Representation of Arabic Text: Analysis and Modeling

Maged S. Al-Shaibani*

King Fahd University of Petroleum and Minerals

Irfan Ahmad ** †

King Fahd University of Petroleum and Minerals

This paper presents a novel dotless representation of Arabic text as an alternative to the standard Arabic text representation. We delve into its implications through comprehensive analysis across five diverse corpora and four different tokenization techniques. We explore the impact of dotless representation on the relationships between tokenization granularity and vocabulary size and compare them with standard text representation. Moreover, we analyze the information density of dotless versus standard text using text entropy calculations. To delve deeper into the implications of the dotless representation, statistical and neural language models are constructed using the various text corpora and tokenization techniques. A comparative assessment is then made against language models developed using the standard Arabic text representation. This multifaceted analysis provides valuable insights into the potential advantages and challenges associated with the dotless representation. Last but not the least, utilizing parallel corpora, we draw comparisons between the text analysis of Arabic and English to gain further insights. Our findings shed light on the potential benefits of dotless representation for various NLP tasks, paving the way for further exploration for Arabic natural language processing.

Keywords: Arabic text representation, dotless text, tokenization, vocabulary size, text entropy, language modeling, cross-lingual analysis.

1. Introduction

Arabic serves as the official language across 22 nations within the Arab world, wielding a native speaker base exceeding 400 million individuals (Guellil et al. 2021). Moreover, it functions as a secondary language for numerous non-Arabic adherents within the Muslim community. Recognized as one of the United Nation’s six official languages (Boudad et al. 2018), Arabic occupies a significant linguistic sphere. Notably, it has ascended to become the fourth most prevalent language on the internet, experiencing rapid growth in user engagement between 2013 and 2018 (Boudad et al. 2018).

The language itself exhibits a tripartite classification. Classical Arabic (CA) represents the linguistic embodiment within sacred texts, tightly interwoven with Islamic cultural and religious heritage, as well as ancient Arabic literature. Modern Standard Arabic (MSA) constitutes the formal register employed in official documentation and news dissemination. Finally, Dialectical or Colloquial Arabic (DA) encompasses a mul-

* Information & Computer Science Department, KFUPM, E-mail: g201381710@kfupm.edu.sa.

** Information & Computer Science Department, KFUPM, E-mail: irfan.ahmad@kfupm.edu.sa.

† SDAIA–KFUPM Joint Research Center for AI, Dhahran 31261, Saudi Arabia.

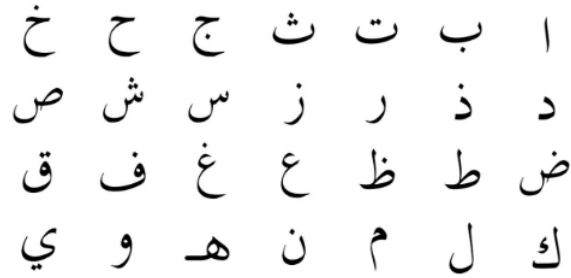


Figure 1: Characters in the Arabic script.

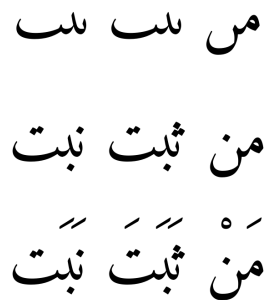


Figure 2: Arabic text with and without dots and diacritics

titude of regionally specific vernaculars derived from MSA, utilized extensively in day-to-day interpersonal communication.

Arabic, akin to various Semitic languages, is distinguished for its dense morphology, characterized by a non-concatenative structure. In this morphological framework, words are crafted by interweaving vowel letters amidst the root consonants, coupled with the incorporation of prefixes and affixes. This intricate morphology poses formidable challenges for prevailing natural language processing (NLP) tools. For instance, the expansion of the language's lexicon surpasses that of languages such as English. Research by Alotaiby, Alkharashi, and Foda (2009) underscores this, revealing that the vocabulary size within an Arabic corpus is twice that of an English corpus from a commensurate domain, despite containing a roughly equivalent number of tokens.

The Arabic script, a cursive writing system, follows a right-to-left orientation. Each letter within this script exhibits diverse glyphs contingent upon its placement within a word. Comprising a total of 28 letters, 25 of these represent consonants. Notably, nearly half of these letters share fundamental base glyphs, known as *rasm*, but are differentiated by the presence of dots either atop, within, or below their *rasm*, a visual distinction elucidated in Figure 1. Additionally, the script incorporates 8 supplementary symbols known as diacritics. These diacritics serve the purpose of offering phonetic guidance, thereby mitigating potential ambiguity.

During the initial stages of the language writing system development, the utilization of dots lacked prominence and formalization (Bentouati 2020). This stemmed from the limited popularity of writing practices during that era, coupled with the relative ease with which native speakers could disambiguate dotless text within its contextual

framework. The inception of standardized dot usage can be attributed to Abu Al-Aswad Al-Dualli, a renowned Arabic scholar, who introduced dots in the sacred text of the Holy Quran. These dots functioned akin to diacritics, with particular emphasis placed on the final letter in a word, representing the word's grammatical function within a sentence. However, this proposal was subsequently supplanted by the utilization of dots to differentiate letters sharing identical *rasm* and diacritics to emulate phonetic assistance.

The primary motivation behind this endeavor was to facilitate Arabic comprehension for non-native speakers thereby safeguarding against foreign influences on subsequent generations of native speakers. Figure 2 showcases an Arabic sentence depicted in three distinct formats. The first is the dotless sentence, followed by the sentence with dots, and finally, the sentence incorporating both dots and diacritics. Notably, the second word illustrates the potential for multiple meanings, a disparity resolved through the use of dots.

Considerable attention has been devoted to the automatic diacritization of Arabic text, as evidenced by extensive research efforts documented in works such as (Diab, Ghoneim, and Habash 2007; Abandah et al. 2015; Metwally, Rashwan, and Atiya 2016; Mubarak et al. 2019). This emphasis is primarily due to the potential utility of diacritization in facilitating downstream tasks like machine translation (Diab, Ghoneim, and Habash 2007), word sense disambiguation and part-of-speech tagging (Darwish, Mubarak, and Abdelali 2017). Comprehensive surveys conducted by Hamed and Zesch (2017) and Azmi and Almajed (2015) aim to delineate the advancements achieved in this domain.

Conversely, in the realm of Optical Character Recognition (OCR), the treatment of dotless text has garnered attention. Certain studies have explored techniques focused on the separation of dots from their respective *rasms* during the OCR process (Ahmad and Fink 2015, 2019). However, the utilization of dotless text in natural language processing tasks, such as sentiment analysis or language modeling, remains, to the best of our knowledge, an unexplored avenue. This presents a promising area warranting further investigation and potential exploration.

Before delving into an experimental analysis to evaluate the applicability and utility of dotless Arabic text as a representation for NLP tasks, conducting a preliminary statistical analysis would offer valuable insights. This statistical examination encompasses token counts across multiple levels of granularity, namely words, subwords, and characters.

Additionally, employing quantitative analytical tools such as Zipf's law (Zipf and Behavior 1949) would be beneficial. Zipf's law elucidates an exponential relationship between the frequency of vocabulary items v_i and their respective rank. Extensively studied across numerous languages (Yu, Xu, and Liu 2018), literary styles, and varying degrees of morphological abstraction (Moreno-Sánchez, Font-Clos, and Corral 2016), Zipf's law has demonstrated applicability across diverse linguistic units, encompassing phrases, pragmatics (Sicilia-Garcia et al. 2002), spoken dialogues (Linders and Louwerse 2022), and even in domains seemingly unrelated to natural languages, including music (Serrà et al. 2012), users' telephone call patterns (Newman 2005), programming languages (Sano, Takayasu, and Takayasu 2012), constructed languages like Esperanto (Manaris et al. 2006), and even in chess openings (Blasius and Tönjes 2009).

Furthermore, Heap's law (Heaps 1978) serves as another valuable analytical tool. This law explores the growth of vocabulary within a given corpus relative to its size. According to Heap's law, the expansion of vocabulary exhibits exponential growth proportional to the corpus size, characterized by an exponential parameter β typically

less than 1. This implies that the growth rate of vocabulary lags behind the rate of corpus size expansion, elucidating a slower-than-linear increase in vocabulary size. Employing these statistical tools can provide foundational insights essential for subsequent experimental analyses.

Employing dotless Arabic text presents a promising avenue to address various challenges encountered in Arabic natural language processing (ANLP). The inherent density of Arabic morphology often results in a considerably expansive language vocabulary. However, dotless text adoption could aid in mitigating this issue by consolidating many dotted words into a singular dotless homographic word. This consolidation contributes to reducing the overall vocabulary size, offering a potential solution within ANLP.

Moreover, this line of inquiry holds significant relevance in the realm of automatic recognition of ancient parchments where the script used was predominantly dotless. By leveraging dotless text methodologies, researchers and practitioners can enhance their capabilities in deciphering and analyzing historical documents etched in dotless Arabic script.

Furthermore, this study's findings can be extended to augmenting processing tools for social media platforms, especially amid a recent trend encouraging users to employ dotless text instead of dotted text in posts that contravene platform regulations (Rom and Bar 2021; Alimardani and Elswah 2021). These studies have highlighted that despite the absence of dots, native speakers retain a substantial ability to recognize and interpret dotless text, underscoring the resilience and adaptability of native speakers within this context. Incorporating dotless text methodologies could serve as a means to address platform regulation breaches while acknowledging the inherent readability of dotless text by native speakers.

Our research interest is deeply rooted in the historical context wherein ancient Arabs adeptly interpreted dotless text, motivating us to undertake a comprehensive analysis encompassing both dotted and dotless textual dimensions. Our primary objective is to assess the efficacy of context-aware deep learning models specifically when applied to dotless text, contrasting their performance with the dotted variant. To explore this, we propose leveraging language modeling as a fundamental task within NLP, wherein we meticulously analyze the perplexity of text generated from both dotted and dotless formats.

We also hypothesize that through harnessing sufficiently extensive contextual information, these sophisticated algorithms hold the potential to achieve a level of proficiency in disambiguating dotless text that rivals human-level comprehension. This pursuit seeks to elucidate the capacity of modern computational models to grasp and interpret dotless Arabic text, drawing parallels with the interpretive capabilities of historical Arab scholars.

To summarize the contribution of this paper, it presents:

- A novel method to represent Arabic text by removing dots from the Arabic letters.
- A comprehensive analysis of Arabic text rigorously comparing dotted and dotless text conducted on five different datasets, varying in size and domains. The analysis encompasses multiple tokenization levels, specifically character, word, and subword tokenizations. Additionally, this study supplements its findings by conducting a comparative analysis

between Arabic dotted and dotless text and English text, enhancing the depth and breadth of the research insights.

- An extensive examination and comparison of dotted and dotless text conducted with two primary language modeling techniques – statistical n-grams and neural language models – utilizing corpora sourced from diverse domains. This analysis encompasses various tokenization schemes. To the best of our knowledge, no existing literature comprehensively covers this specific analysis, making it a valuable contribution to the field of language modeling techniques.

The remaining sections of this paper are organized as follows: Section 2 delves into the advancements and developments pertinent to this topic in the existing literature. Section 3 provides an overview of the corpora utilized in this study, detailing the preprocessing steps undertaken and elucidating the employed tokenization methods. In Section 4, we present the various text analyses conducted to compare dotted and dotless text. Section 5 outlines our experiments involving language modeling. Section 6 offers a comparative analysis between dotted, dotless text, and English text to maintain a sense of broad applicability and generalization. Finally, Section 7 concludes this study, encapsulating the findings and envisaging future directions for research.

2. Related Work

As far as our research indicates, this is the first work to analyze dotless Arabic text in the domain of natural language processing tasks. However, the utilization of dotless text in social media to circumvent content filtering algorithms has attracted recent attention (Alimardani and Elswah 2021). In response to this trend, Rom and Bar (2021) delved into methodologies to accommodate undotted text within the prevailing settings of language models trained predominantly on dotted text, particularly in the context of social media platforms. Their exploration yielded two primary approaches. The first approach involved reconfiguring the tokenizer either by removing dots from its tokens or expanding it with undotted variants of these tokens. These undotted versions were then mapped to the same identifier as their dotted counterparts, effectively enabling the language model to accept both dotted and dotless text. The second approach proposed a dotting algorithm aimed at restoring dots to the undotted text. Through this method, they demonstrated that the restored dotting approach produced results closely aligned with the original experiments conducted on dotted text when evaluated across various downstream tasks.

Alhathloul and Ahmad (2022) presented a deep learning-based approach aimed at storing dots for a given dotless text. Employing an extensive experimentation strategy across various datasets, they conducted an in-depth error analysis to assess the efficacy of their proposed method. Their approach relied on a Gated Recurrent Networks (GRUs) based sequence-tagging model, designed to generate dotted text as an output from dotless text input. In order to facilitate a more equitable comparison and evaluation, the researchers introduced a new metric called "dottization error rate." Across their experimentation involving four distinct datasets, their method exhibited promising results, achieving character error rates ranging between 2% and 5.5%, along with dottization error rates spanning from 4.2% to 11%.

Zipf's and Heap's laws serve as pivotal statistical tools in linguistics, offering profound insights when applied to language corpora, extensively examined in languages

like English, Chinese, Norwegian, and Old German (Li 1992; Moreno-Sánchez, Font-Clos, and Corral 2016; Sicilia-Garcia et al. 2003; Welsh 1988). However, their application to Arabic remains relatively limited, potentially influenced by Arabic's perception as a low-resource language (Magueresse, Carles, and Heetderks 2020). Within Arabic language studies, two primary research trajectories emerge. Firstly, analytical comparisons of Arabic with other languages have been explored (Alotaiby, Alkharashi, and Foda 2009; Alotaiby, Foda, and Alkharashi 2014; Al-Kadi 1998). Secondly, these analyses have been utilized to assess the quality of proposed corpora (Alarifi et al. 2012; Almeman and Lee 2013; Selab and Guessoum 2015; Khreisat 2009). Notably, much of this research has focused on evaluating extensive corpora within a singular domain, predominantly newswire. Consequently, a notable research gap exists in conducting a comprehensive analysis encompassing diverse tokenization levels, linguistic units, and datasets spanning various domains within Arabic language studies.

Language models constitute mathematical models designed to predict the most probable word within a sequence of preceding words, representing a critical component in numerous natural language processing endeavors. Their utility extends across various tasks, including generative applications like machine translation (Diab, Ghoneim, and Habash 2007), optical character recognition enhancement (Smith 2011), and the augmentation of automatic speech recognition systems (Abushariah et al. 2010). Leveraging their capacity to encapsulate implicit linguistic information, neural language models find applications in sequence labeling tasks such as part-of-speech tagging and named entity recognition through a transfer learning style. Within these contexts, language models play a pivotal role in generating high-quality embeddings to enhance subsequent classification processes.

The evolution of language modeling boasts a substantial developmental history. Initially, these models sought to learn the joint probabilities of sequences within expansive training corpora. Employing the Markov principle, these sequences were pruned to encompass orders 2, 3, 4, or higher n tokens, commonly known as n -grams, aiming to achieve an acceptable level of generalization. However, this approach grappled with an inherent challenge to dealing with Out-Of-Vocabulary (OOV) tokens, wherein the model encountered difficulty computing the probability of words absent from the training corpus. To address this hurdle, various smoothing techniques emerged, including Laplace (additive) smoothing, Katz smoothing (Katz 1987), and Kneser-Ney smoothing (Kneser and Ney 1995). Subsequently, specialized toolkits emerged as standardized solutions for constructing these language models using sufficiently extensive corpora. Prominent examples of such toolkits include KenLM (Heafield 2011) and SriLM (Stolcke 2002). These toolkits have streamlined the development process by providing efficient means to build robust language models, marking significant strides in the field of language modeling.

Neural language models surpassed statistical counterparts (Bengio, Ducharme, and Vincent 2000), with Recurrent Neural Networks (RNNs) initially dominating sequence modeling. However, they grappled with issues like vanishing or exploding gradients, particularly evident in lengthy sequences. To address this, Long Short-Term Memory (LSTM) architectures were introduced (Hochreiter and Schmidhuber 1997), albeit slower and more challenging to train. Gated Recurrent Units (GRUs) (Chung et al. 2014) emerged as a compromise between Vanilla RNNs and LSTMs. Recently, transformer-based models (Vaswani et al. 2017) sparked a revolution in Natural Language Processing (NLP), overcoming the long recurrence limitation of previous recurrent networks. Pretrained models like BERT (Devlin et al. 2018), Wav2Vec2 (Baevski et al. 2020), and GPT-3 (Brown et al. 2020) attained state-of-the-art performance across various

NLP tasks using this architecture. Nonetheless, the attention mechanism, intrinsic to transformers, exhibited limitations in mastering simpler tasks easily tackled by earlier models (Dehghani et al. 2018; Chernyavskiy, Ilvovsky, and Nakov 2021).

3. Text Corpora, preprocessing and Tokenization

In this section, we introduce the corpora utilized in this research, outlining the preprocessing steps employed in their preparation. Additionally, an overview of the tokenization methods applied is presented for comprehensive understanding and context within this study.

3.1 Text Corpora

We selected corpora from different domains with different sizes to evaluate dotless text as an alternative representation for Arabic NLP. The following are the selected corpora:

- **Quran** (Aloufi 2019): a dataset compiled from Quran verses. In Quran, most chapters, called Surah, start with a special sentence, called Basmala. In order not to confuse the training with this repetitive sentence, we removed it from the beginning of each chapter except the first one. The total number of samples is 6236.
- **Sanadset dataset** (Mghari, Bouras, and El Hibaoui 2022): This dataset is a collection of Sacred narrates by the prophet Muhammed called Hadeeth. Each Hadeeth contains a Matn which is the sacred text and a Sanad which is the chain of narrators. We selected only the Matn. This dataset is interesting because its text is coherent with a relatively small vocabulary size compared to its size.
- **Ashaar (Poems)** (Alyafeai and Al-Shaibani 2020): This is a poetry dataset collected from various online sources. Samples extracted from this dataset comprises a set of 6 verses. Throughout this research, we will refer to this dataset as "Ashaar" or "Poems" dataset.
- **Newsware** dataset (El-Khair 2016): This dataset is a massive collection of news collected from various news wires. We only selected one newspaper called Alittihad.
- **Wikipedia** (Foundation): a dataset collected from Wikipedia dumps of Arabic articles up to 20-10-2022.

Table 1 presents detailed statistics with respect to dotted and dotless tokens for each dataset. In this table, N represents the running text, i.e. all words, V is the unique vocabulary, and V' is the unique dotless vocabulary. As illustrated in the table, there is a noticeable reduction of the dotless vocabulary as compared to dotted. It can also be noticed that the Quran dataset is the smallest dataset compared to others. For Sanadset dataset, it contains very coherent text with many repetitive phrases. This can be deduced from the ratio of its unique vocabulary to its running text. This measure highlights the abundant richness of the poetry dataset, which aligns with our expectations, given the inherent creativity associated with poetry. Furthermore, classical poetry adheres to strict metrical and rhyming constraints, which compel poets to seek words that harmonize with the poem's rhyme scheme. Regarding news and Wikipedia datasets, they are huge

Dataset	N_w	V_w	V'_w	N_c	V_c	V'_c
Quran	77,797	14,748	13,229	330,709		
Sanadset	28,880,818	317,331	227,590	112,879,152		
Poems	34,940,290	1,007,279	631,487	145,574,240	31	19
News	134,862,816	892,583	654,982	663,796,116		
Wikipedia	177,422,512	1,811,244	1,345,853	851,699,872		
Aggregated	376,184,233	2,739,172	1,865,126	1,774,280,089		

Table 1: Datasets vocabulary and characters statistics, dotted and dotless

Letter	Mapping	Letter	Mapping	Letter	Mapping	Letter	Mapping
ا	ا	د	د	ض	ص	ك	ك
ب	ب	ذ	ر	ط	ط	ل	ل
ت		ر		ظ		م	م
ث		ز		ع		ن	ن
ج	ح	س	س	غ	ع	ه	ه
ح		ش		ف		و	و
خ		ص		ق		ي	ي

Table 2: Arabic dotted letters mapped to their dotless variant

corpora capturing various aspects and structures of the language. However, the news dataset has a much narrower domain as compared to Wikipedia.

From the above statistics, it is clear that the datasets we chose are of different sizes and can be grouped into three categories. From the size of the running text N , the Quran text is a significantly small corpus, Sanadset and Ashaar are medium-sized corpora, Wikipedia and news are larger corpora.

3.2 Text pre-processing and undotting

In this subsection, we present the procedure we followed to pre-process and undot text. For pre-processing, we removed any non-Arabic characters including numeral characters and punctuation symbols. Diacritics are also omitted. Further, as the letter Hamza can appear adjoined with other characters, usually vowels, we removed it keeping only the adjoined letter. However, if it comes alone, it is kept unchanged.

For undotting, table 2 illustrates Arabic letters and their dotless mapped letter. Due to the cursive nature of the Arabic script, we tried to mimic the same writing shape for a few of the letters to its closest dotless letter if it comes either at the beginning or the middle of the word. For the letters ن and ي, if they come either at the beginning or the middle of the word, they are mapped to the dotless version of ب. Similarly, for the letter, ق, if it comes at the beginning or middle of the word, it is mapped to the dotless version of ف.

3.3 Tokenization

Tokenization stands as a crucial process essential for converting text into integer identifiers, serving as inputs for deep learning systems. This transformation involves segmenting text into various levels of granularity, including words, characters, or in-between units known as subwords, achieved through a meticulous segmentation procedure. Notably, subword tokenization emerges as an intriguing approach due to its adaptability, being either language-dependent or language-agnostic. This method integrates benefits from different ends of the spectrum, presenting a versatile and robust means to represent input text. The proliferation of large language models such as BERT (Devlin et al. 2018) and GPT-3 (Brown et al. 2020) significantly contributed to the popularity of subword tokenization. Data-driven approaches like WordPiece (Song et al. 2020) and SentencePiece (Kudo and Richardson 2018) exemplify such tokenization methodologies, underscoring their significance in modern language modeling techniques.

In this research, we studied the behavior of dotless text compared to dotted text in the context of these levels of granularity. We used word tokenization, character tokenization, morphology-based, and language-specific subword tokenizations. These tokenizations are described as follows.

- **Word tokenization:** This tokenization splits text into words.
- **Farasa Morphological tokenization:** This tokenization splits text into morphemes. That is roots and affixes. The tool used to deliver this tokenization is called farasa (Abdelali et al. 2016), hence the name. We used this tokenization to represent morphological splits of the Arabic text.
- **Disjoint-Letters tokenization:** This tokenization is a language dependent tokenization introduced by (Alyafeai et al. 2023). Although the Arabic script is cursive, some of its letters are written disconnected from their subsequent letter in the word. Based on this property, this tokenization splits the text into subwords according to this property of the letters. That is, each subword in this tokenization is a sequence of fully connected letters.
- **Character tokenization:** This tokenization splits text into character.

Figure 3 shows an example of an Arabic sentence being tokenized with the aforementioned tokenizations.

(Alyafeai et al. 2023) conducted an in-depth comparison of various types of tokenization for Arabic where this tokenization was first introduced for NLP tasks. To perform the tokenization task, we extended their tokenizers package, named tkseem¹, to fit our needs while utilizing their implemented functionalities.

4. Text Analysis

This section delineates the experimental framework established to compare dotless text with its dotted counterpart. Initially, it provides a comprehensive exposition of detailed

¹ <https://github.com/ARBML/tkseem>

العربية إحدى أكثر اللغات انتشارا عالميا																													
Normal Text																													
العربية	إحدى	أكثر	اللغات	انتشارا	عالميا																								
Word Tokenization																													
ال	عربي	ة	إحدى	أكثر	ال	لغ	ات	انتشار	ا	عالمي	ا																		
Farasa Morphological Tokenization																													
ا	لعربية	إ	حدى	أ	كثر	ا	للغات	ا	نتشا	ر	ا	عا	لميا																
Disjoint-Letters Tokenization																													
ا	ل	ع	ر	ب	ي	ة	إ	ح	د	ى	أ	ك	ث	ر	ا	ل	ل	غ	ا	ت	ا	ن	ت	ش	ا	ر	ا	...	
Character Tokenization																													

Figure 3: An Arabic sentence tokenized with different tokenizations

vocabulary statistics spanning diverse tokenization levels. Subsequently, it delves into an examination of dotless and dotted text within the purview of Zipf's and Heap's laws, offering insights into their respective applicability and manifestations. Furthermore, it encompasses a subsection dedicated to elucidating the outcomes derived from conducting text entropy analysis over both dotted and dotless text, providing a nuanced understanding of their respective linguistic properties and information content.

4.1 Vocabulary Statistics

As this work presents a new representation for writing Arabic text based on the removal of the dots, it is a valuable endeavor to discuss dotless vocabulary as compared to dotted ones on all our datasets and across multiple tokenizers. Tables 3, 4, and 5 present such statistics for words, farasa, and disjoint tokenizations, respectively. It also presents the average tokens' length for each dataset. In these tables, V' represents dotless vocabulary, S is the average vocabulary length for the dotted text, and S' represents average vocabulary' length for dotless text.

Table 3 serves to elucidate the vocabulary statistics pertaining to words tokenization across various datasets. This table reveals notable insights, particularly in gauging

Dataset	V	N	V/N	V'	V'/N	V'/V (%)	S	S'
Quran	14,748	77,797	18.96	13,229	17	89.70	5.39	5.49
Sanadset	317,331	28.88M	1.10	227,590	0.79	71.72	5.95	6.23
Ashaar	1,007,279	34.94M	2.88	631,487	1.81	62.69	6.13	6.49
Wikipedia	1,811,244	177.42M	1.02	1,345,853	0.76	74.31	7.33	7.82
News	892,583	134.86M	0.66	654,982	0.49	73.38	6.73	7.09
Aggregated	2,739,172	376.18M	0.73	1,865,126	0.5	68.09	7.12	7.66

Table 3: Datasets statistics with words tokenization covering dotted and dotless vocabulary

Dataset	V	N	V/N	V'	V'/N	V'/V (%)	S	S'
Quran	7,677	129,788	5.92	6,074	4.68	79.04	4.46	4.65
Sanadset	105,469	44,069,593	2.39	74,674	0.17	70.85	5.50	5.89
Ashaar	361,739	57,350,212	0.63	242,735	0.42	67.09	5.79	6.22
Wikipedia	1,002,148	303,831,223	0.33	791,084	0.26	78.91	7.32	7.81
News	344,627	241,222,485	0.14	263,143	0.11	76.36	6.45	6.88
Aggregated	1,417,437	646,603,301	0.22	1,061,481	0.16	74.92	7.14	7.66

Table 4: Datasets statistics with farasa tokenization covering dotted and dotless vocabulary

Dataset	V	N	V/N	V'	V'/N	V'/V (%)	S	S'
Quran	6,484	164,609	3.94	4,303	2.61	66.36	4.06	4.29
Sanadset	76,495	55,892,178	0.14	37,412	0.07	48.91	4.95	5.31
Ashaar	191,393	72,437,612	0.26	79,968	0.11	41.78	5.16	5.56
Wikipedia	266,822	412,843,347	0.06	118,150	0.03	44.28	5.54	6.07
News	151,557	326,024,821	0.05	68,943	0.02	45.49	5.31	5.75
Aggregated	402,019	867,375,567	0.05	162,389	0.02	40.39	5.59	6.11

Table 5: Datasets statistics with disjoint-letters tokenization covering dotted and dotless vocabulary

vocabulary richness, denoted by the size of V in comparison to N. In smaller datasets like Quran, the value of this metric appears notably high, contrasting with sufficiently large datasets where this value may even fall below 1. This reduction, expounded by Zipf's and Heap's laws, illustrates that within expansive datasets, a small subset of V accounts for a substantial portion of the running text N. Delving deeper, disparities in vocabulary richness among datasets become apparent, such as Ashaar being richer than Sanadset. This discrepancy arises from the inherent constraints imposed by the creativity, classical metric, and rhythmic structure in classical Arabic poetry, necessitating a broader vocabulary to adhere to these constraints. Conversely, Sanadset exhibits higher consistency and lesser diversity in topics, resulting in a less rich vocabulary. Furthermore, in the comparison between Wikipedia and news datasets, the former displays higher richness owing to its diverse domains and topic variety, whereas the

latter, sourced from fewer domains with more consistent article structures, manifests a lower richness. These observations underscore the influence of content diversity and structural consistency within datasets on their respective vocabulary richness, offering valuable insights into their distinct linguistic characteristics.

The comparison of V with respect to N for subword tokenizations is detailed in Table 4 for farasa tokenization and Table 5 for disjoint tokenization. Evident from both tables is the finer granularity of disjoint tokenization compared to farasa, indicating a higher prevalence of disjoint letters within the language vocabulary in contrast to morphologically segmented components. Moreover, the discrepancy observed between V for farasa and disjoint tokenizations is influenced by two key factors: the extent of running text (N) and the inherent richness of the dataset. Notably, this discrepancy is more pronounced when datasets exhibit richness and consequently higher N values. This trend is particularly observable in datasets like Wikipedia and Ashaar within the dotted text domain, where the difference in V between farasa and disjoint tokenizations is notably elevated. This delineates the influence of dataset richness and the volume of running text on the variance between farasa and disjoint tokenizations, emphasizing their impact on the vocabulary distinctions observed across different tokenization schemes.

A clear trend observable across all tokenizations from these tables is that S' —representing the average vocabulary length for dotless text—is generally marginally higher than S , denoting the average vocabulary length for dotted text. This recurrent pattern underscores the characteristic distinction between least and most frequent vocabulary items across various tokenizations and the subtle disparity in average vocabulary lengths between dotless and dotted text.

From these tables, it can be noticed that rich datasets characterized by richness tend to exhibit smaller V'/V ratios across almost all tokenizations. For example, the Wikipedia dataset showcases lower V'/V values for words and disjoint tokenization compared to the News dataset, and a comparable ratio in farasa tokenization, despite having a considerably larger running text N . This pattern is even more apparent in the poetry dataset, which boasts the smallest V'/V ratio owing to its richness in vocabulary. This observation leads to a hypothesis that, asymptotically, V' grows slower concurrently with N in the long run. In essence, it suggests that, generally, the least frequent vocabulary items tend to share more *rasms* compared to the more frequent ones, leading to a proportional growth of dotless vocabulary as the volume of running text expands.

The observation that farasa's V'/V ratio is higher than Word's V'/V for datasets with high vocabulary as can be shown in table 4 is of a particular interest. This discrepancy suggests that subwords segmented by farasa might not share as many *rasms* when dots are removed, in comparison to words tokenization. This disparity in the V'/V ratio between farasa and words tokenizations highlights a potential difference in the underlying structure of segmented subwords, indicating that the removal of dots might affect these segmented subwords differently in terms of shared *rasms*.

An other interesting finding is the consistently smaller V'/V ratio observed for disjoint tokenization compared to other tokenizations. This ratio is even smaller than the characters' V'/V ratio in the three largest datasets. This discrepancy suggests that these subwords segmented by disjoint tokenization are notably rich in dots and, consequently, exhibit a higher likelihood of sharing the same *rasms* upon the removal of dots. This observation underscores the distinct characteristic of disjoint subwords, implying a stronger association or shared *rasms* among these segments when dots are excluded.

Another final notable observation is the consistently lowest V'/V ratio observed in the aggregated dataset across all tokenizations. This observation implies that with larger

datasets, a greater proportion of *rasms* are being shared. This trend suggests that as the dataset size increases, there is a higher tendency for shared *rasms* among vocabulary.

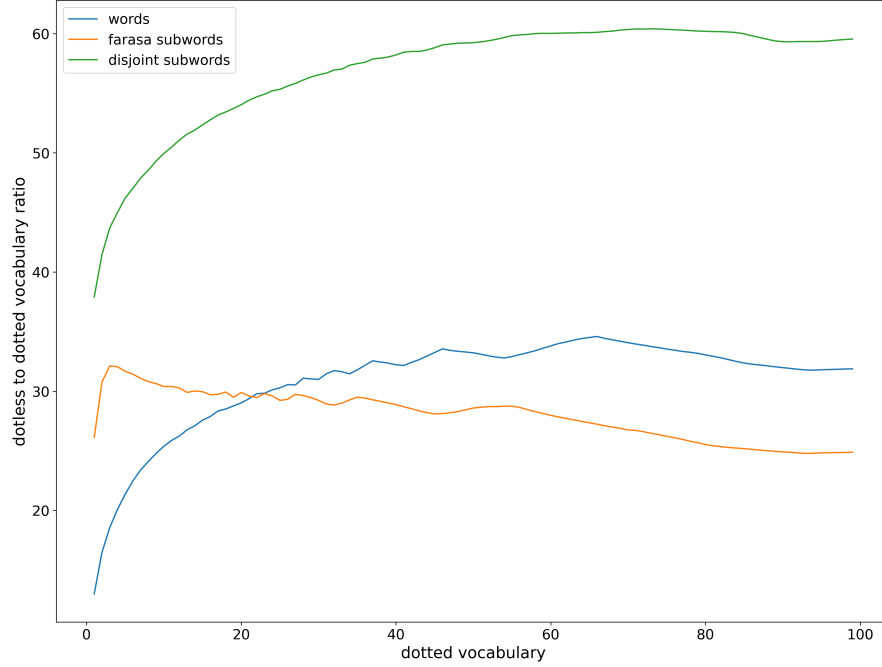


Figure 4: Dotless vocabulary reduction ratio on different tokenization levels

To explore the correlation between dotted and dotless vocabulary, we conducted an analysis on the top i frequently occurring vocabularies, ranging from 1 to 100. Subsequently, we transformed these vocabularies into their dotless forms and calculated the ratio between dotless and dotted vocabulary counts. This analysis was conducted across aggregated datasets employing words, farasa, and disjoint tokenizations. The findings of this analysis are represented in 4.

Upon examination of the plot, a noteworthy observation emerges regarding vocabulary reduction. The most significant reduction occurs within the disjoint tokenization scheme, where the ratio of dotless to dotted vocabulary size experiences the highest reduction, starting from approximately a 40% decrease. Additionally, it becomes evident that the reduction in vocabulary through words tokenization initially starts at a lower rate but increases as the vocabulary expands, surpassing the reduction observed in the farasa tokenization.

Moreover, across all three tokenization types, the ratio exhibits fluctuations, displaying intermittent rises and falls, forming plateaus. As a result of this analysis, we hypothesize that while many common vocabularies share identical *rasms*, this ratio stabilizes as the vocabulary expands, corroborating our previous observations.

Furthermore, it is observed that within the farasa tokenization, the most frequently occurring vocabularies, anticipated to represent affixes, exhibit the highest number of shared *rasms* compared to other tokenizations, whereas its least common vocabularies display the lowest number of shared *rasms*. This inference suggests that affixes are notably abundant in dots compared to other vocabularies in this specific tokenization scheme.

Additionally, the reduction pattern observed in disjoint tokenization closely resembles that of words but appears smoother. Notably, the decrease in the ratio for the least common vocabularies in words exhibits a slightly faster descent compared to those in disjoint tokenization. This discrepancy implies distinct structural characteristics within these tokenization schemes.

4.2 Zipf's and Heap's laws

We applied Zipf's law to examine the correlation between vocabulary and its frequency. This empirical law delineates a seemingly straightforward exponential relationship between the rank of vocabulary and its frequency when the vocabulary is arranged in descending order based on frequency, particularly within a considerably large corpus. Our objective in this analysis was to visually explore the frequency distributions of both dotted and undotted vocabularies across varying corpus sizes.

Zipf's law can be mathematically expressed as shown in equation 1, where r signifies the rank of the vocabulary, $F(r)$ represents its frequency, and α stands for a corpus-specific parameter typically approximating unity. The primary aim of our investigation was to scrutinize how this law manifests itself in the context of the frequency distributions of dotted and undotted vocabularies across diverse corpus sizes.

$$F(r) \propto r^{-\alpha} \quad (1)$$

To fit a corpus with Zipf's law and determine the value of α , we applied a log transformation to both sides of the equation, as illustrated in equation 2. By setting the constant C as the frequency of the most frequent vocabulary within each dataset, we aimed to prevent potential biases when fitting the regression line. This transformation led to the formulation of an equation that exhibits near-linearity on a log-log scale, allowing us to employ linear regression to estimate the slope (α).

$$\log F(r) = -C\alpha \log r \implies \alpha = -\frac{\log F(r)}{C \log r} \quad (2)$$

In investigating the expansion of vocabulary alongside the running text within our corpus, we employed Heap's law (Heaps 1978). This law similarly demonstrates an exponential correlation between the size of the running text denoted by n and the corresponding vocabulary size denoted by V , governed by an exponent parameter β typically falling below 1. The essence of this relationship is articulated in equation 3.

$$V \propto n^\beta \implies V = kn^\beta \quad (3)$$

Following the same procedure conducted for Zipf's law to estimate the parameters, we transform both sides of the equation using the log function as in equation 4. We fit this equation using linear regression to estimate both values k and β .

$$\log V = \log k + \beta \log n \quad (4)$$

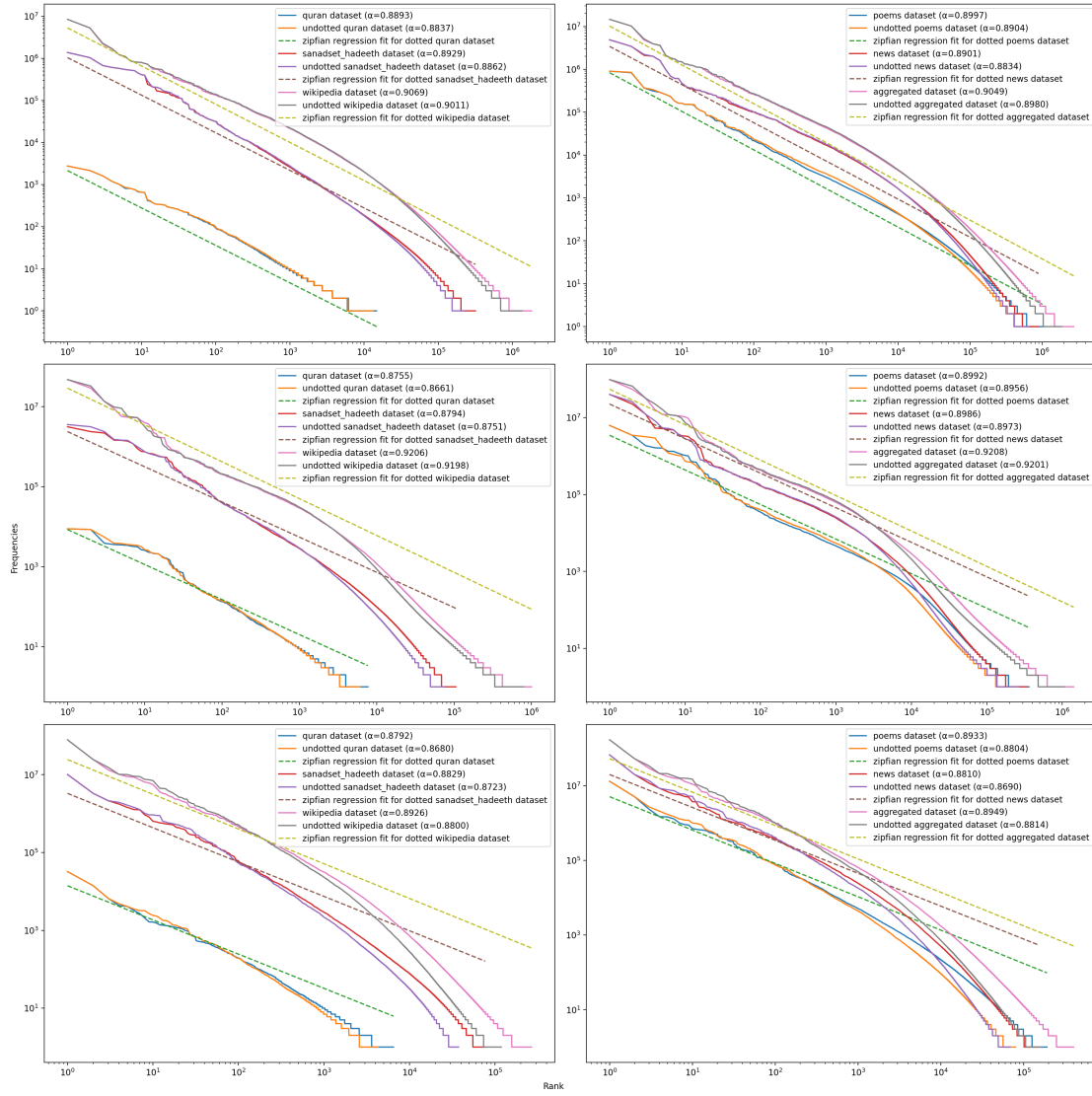


Figure 5: Zipf's Law plots for various datasets with different tokenizations

Figure 5 plots Zipf's law applied to our datasets. Each row represents a tokenization level. The first row is for word tokenization, the second row is for farasa tokenization, and the third row is for disjoint-letters tokenization. In each of these rows, datasets are split into two figures as it is difficult to plot all of them in one graph.

From an overarching perspective gathered from these graphs, it is evident that all tokenization methods adhere to Zipf's law. Nevertheless, the coarse-grained tokenizations exhibit a more pronounced adherence to Zipf's law, showcasing a stronger Zipfian distribution compared to other finer-grained tokenization approaches.

The depicted figure reveals a nearly identical pattern between dotless and dotted frequencies, with marginal distinctions apparent primarily within the most frequent and least frequent vocabularies. Notably, across various datasets, a consistent trend

emerges where, with finer-grained tokenization methods, the dotless frequencies exhibit a slightly lower count. This discrepancy becomes particularly conspicuous in the disjoint plots. Remarkably, this observation remains consistent across all datasets, suggesting that this behavior is independent of the specific dataset characteristics.

Additionally, a distinct pattern emerges within the farasa vocabulary. Across all datasets, a consistent decline in the frequencies of vocabulary—both dotted and dotless—can be observed, occurring between what we can term as the most frequent and least frequent vocabularies. This drop in the frequencies of the least frequent vocabulary exhibits a consistent behavior, irrespective of the dataset or tokenization method employed. However, an intriguing observation arises concerning the decline in the frequencies of what we designate as the mid-frequent vocabulary. It appears that the most frequent vocabularies, presumed to represent affixes due to their high frequency resulting from their attachment to multiple roots, initiate this decline from the mid-frequent vocabulary and continue downward. This particular pattern seems to manifest across various datasets and other tokenization methods, suggesting a consistent relationship between these mid-level frequencies and the prevalence of affixes within the farasa vocabulary.

Figure 6 illustrates Heap’s law applied to various tokenizations across all datasets. Each row within the figure corresponds to a distinct tokenization type, while the first column depicts the plots for the Quran datasets, notably smaller in terms of running text compared to the other datasets. The subsequent column showcases the plots for the remaining datasets, exhibiting more substantial textual content in comparison.

The depicted plots demonstrate the conformity of both dotted and dotless vocabularies to Heap’s law. Notably, the growth of dotless vocabulary consistently trails behind that of the dotted vocabulary across various tokenizations. However, a significant observation arises particularly in the case of disjoint tokenization where the growth rate appears notably slower. This phenomenon suggests a substantial sharing of *rasms* among the dotted vocabularies within this specific tokenization scheme, contributing to a slower expansion of the vocabulary.

Another notable observation is the influence of dataset richness on the growth disparity between dotted and dotless vocabularies. Specifically, a more substantial gap between the growth rates of dotted and dotless vocabularies is noticeable in datasets with higher richness, such as Ashaar and Wikipedia. In contrast, this gap notably narrows in datasets characterized by lower richness, such as news and Sanadset. This divergence suggests a correlation between dataset diversity or size and the distinctiveness in growth rates between the dotted and dotless vocabularies.

4.3 Text Entropy

In the realm of information theory, entropy serves as a metric quantifying the degree of disorder within a given system. Its purpose lies in computing the number of bits necessary to transmit information, i.e., the information content of text in our context. In our investigation, we try to assess the extent of information loss resulting from the elimination of dots across various tokenization levels for all token strings, denoted as T , within our dataset. Equation 5 was employed to compute the entropy of tokens within the corpus, a dimension that has not been extensively explored in Arabic literature. Previous work by (Al-Kadi 1998) delved into entropy calculations on Arabic text, reporting an entropy of 9.98 bits/words and 2.39 bits/letters. This study, considered a corpus comprising 3025 vocabularies, was conducted nearly two decades ago, predating recent advancements in computational resources. Surprisingly, no recent studies have delved

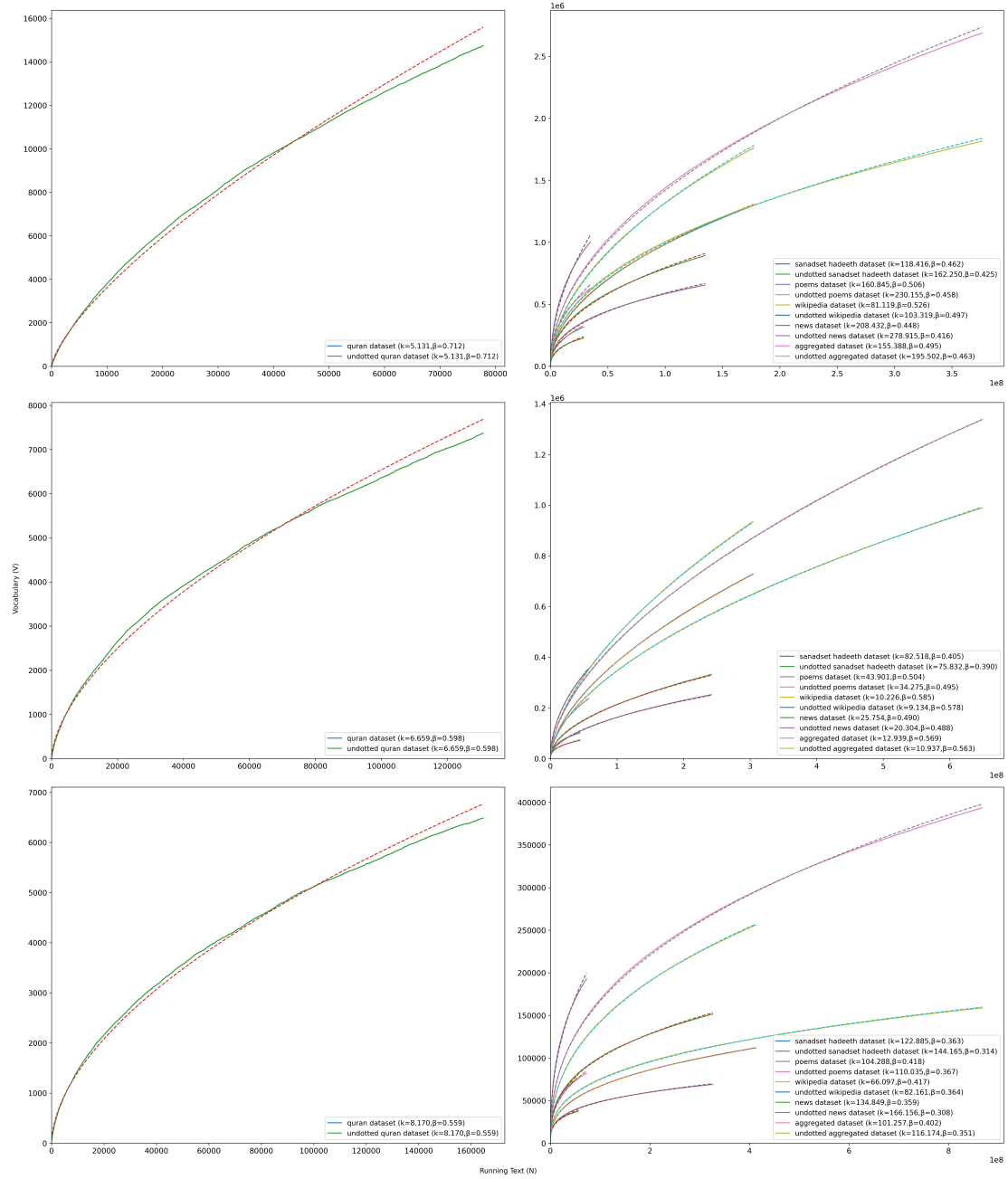


Figure 6: Heap's Law plots for various datasets with different tokenizations

into this particular aspect of Arabic text analysis. Table 6 elucidates the entropy metrics for each tokenization scheme, where H denotes dotted text entropy and H' signifies dotless text entropy.

Dataset	Words		Characters		Farasa		Disjoint	
	H	H'	H	H'	H	H'	H	H'
Quran	11.02	10.87	4.15	3.83	8.240	7.968	7.451	6.958
Sanadset	10.92	10.67	4.24	3.86	8.364	8.004	7.644	7.098
Ashaar	14.33	13.70	4.3	3.9	9.659	9.148	8.480	7.687
Wikipedia	13.20	12.94	4.27	3.87	8.892	8.599	8.153	7.493
News	13.10	12.87	4.25	3.85	8.563	8.286	7.981	7.363
Aggregated	13.6	13.26	4.27	3.87	9.063	8.711	8.210	7.523

Table 6: Entropy results for each tokenization across the proposed datasets

$$H(t) = \sum_t^T -P(t) \log_2 P(t) \quad (5)$$

For character tokenization, we can see that the number of dotless characters drops by a margin of 12 letters with a ratio of 61.3% of the dotted characters. However, we can also notice that the highest entropy for dotted characters is 4.3 and the lowest is 4.15 while the highest for the undotted are 3.9 and the lowest is 3.83. If we consider the aggregated dataset as a representation of the language, we can see that the reduction of entropy due to undotting is 0.4 with a ratio of 9.1%. That is, the reduction in entropy is less compared to the reduction in the characters. From this, we can see that the undotted text still holds entropy close to the dotted one indicating only a little information being lost with this dotless representation compared to the reduction in the character set. This gives an indication that this type of text can be learned with a decent performance close to the dotted text.

The upper bound of the entropy for characters can be calculated using equation 5 when we have a string of 1 character to be $-\sum_i^{31} P(i) * \log_2 P(i) = \log_2 31 = 4.95$ for dotted characters and $\log_2 19 = 4.24$ for undotted ones. Based on this, if we consider the aggregated dataset as a representative set of the language, the language has a characters redundancy of $1 - \frac{4.27}{4.95} = 13.74\%$ for dotted and $1 - \frac{3.87}{4.24} = 8.76\%$ for undotted characters. The redundancy is less for undotted text because dots were helpful in resolving confusion. It can be concluded from this analysis that, dots, at most, decrease the text redundancy by 5%. This also indicates that there are other sources of redundancy shared between dotted and undotted text.

Entropy, as discussed previously, measures the surprise in a sequence given a prior history of sequences. We can see that Ashaar has the highest entropy in word tokens. It is even higher than Wikipedia which is larger by magnitude in terms of running text N. This confirms our hypothesis about the vocabulary richness of Ashaar being due to creativity and the use of unusual phrases and sentences. We can also note that Sanadset has the lowest entropy. It is even lower than Quran dataset which is smaller by magnitude. This indicates that Sanadset samples share a lot of similarities and it is easily predictable.

As can be noticed in Table 6, the entropy results for subwords tokenizations are in consensus with words and characters. This is expected as this tokenization still preserves the datasets' characteristics. It is also interesting to see that the reduction in

entropy is low compared to the reduction in vocabulary (V). This indicates that a little deeper model is required to capture the loss in entropy but with the advantage of much fewer vocabulary. This also confirms the choice of subwords tokenization as a standard method for language modeling in the current practice.

5. Language Modeling

Language modeling stands as a core area of research within NLP, focusing on generating text that mimics human language patterns. Its fundamental function involves learning the associations between tokens present in extensive text corpora. These models play a vital role in diverse downstream applications such as optical character recognition (OCR), speech-to-text (STT), and text grammar correction, among others.

Language models are constructed through either statistical methodologies or deep learning approaches. In statistical methods, text generation relies on determining the most likely token given a preceding sequence of tokens, computed from the corpus to establish conditional probabilities of tokens. These models, often termed n -gram models where n denotes the sequence length of prior tokens considered, operate within a probabilistic framework. Conversely, neural language models employ neural networks to predict the most probable token, diverging from probabilistic strategies. Notably, statistical language models are typically shallower in structure compared to neural models, offering lower language comprehension abilities in their linguistic analyses.

Statistical language models encounter significant challenges in predicting unseen or infrequently used words due to their reliance on word probabilities computed from corpus counts. To mitigate this issue, smoothing techniques have been developed, including Laplace smoothing, discounting, interpolation, backoff, and notably, Kneser-Ney smoothing (Ney, Essen, and Kneser 1994). Acknowledged as one of the most effective smoothing methods in literature (Zhang and Chiang 2014), Kneser-Ney smoothing operates by generating the subsequent token using Equation 6. In this equation, $c(w_{i-1}, w_i)$ represents the count of bi-grams in the training corpus, δ denotes a discounting factor, $\sum_w c(w_{i-1}, w)$ signifies the sum of counts of all words w following the word w_{i-1} in the training corpus, and λ_{w-1} stands for a normalization factor.

$$p_{KN}(w_i|w_{i-1}) = \frac{\max(c(w_{i-1}, w_i) - \delta, 0)}{\sum_{w'} c(w_{i-1}, w')} + \lambda_{w_{i-1}} p_{KN}(w_i) \quad (6)$$

In this study, we experimented with both of these types of language modeling. For statistical language models, we used KenLM toolkit (Heafield 2011). KenLM is an n -grams language model toolkit that implements the modified Kneser-Ney smoothing (Heafield et al. 2013). The toolkit implements disk-based streaming algorithms allowing it to be 7 times faster than its counterpart SriLM (Stolcke 2002), a popular language models toolkit, in estimation with efficient memory management (Heafield 2013). We experimented with different orders of grams ranging from bigrams to 6 grams. However, other hyper-parameters are set to defaults. For datasets' splits, each dataset has been split into 90%, 10% subsets for training and testing respectively.

Language models are subject to evaluation through two distinct approaches: extrinsic and intrinsic. Extrinsic evaluation involves assessing the model's performance within a downstream task, such as speech recognition or optical character recognition. Conversely, intrinsic evaluation entails self-assessment of the model's competence without reliance on external tasks. Perplexity PPL serves as a metric for intrinsic evaluation,

Dataset	Text Type	Train Vocab	Val Vocab	Test Vocab
Quran	Dotted	13,435	1,720	3,161
	Dotless	12,102	1,658	3,018
Sanadset	Dotted	299,162	81,763	121,488
	Dotless	215,758	65,941	95,132
Poems	Dotted	881,669	208,236	321,650
	Dotless	560,069	151,916	225,293
Wikipedia	Dotted	970,567	211,025	323,331
	Dotless	730,885	173,672	259,353
News	Dotted	837,244	231,783	333,842
	Dotless	617,593	187,819	263,524

Table 7: Train, validation, and test statistics for datasets used in language modeling tasks

quantifying the model’s surprise when presented with a given sentence, leveraging the probabilities learned during training. The calculation of perplexity is expressed in Equation 7 as follows:

$$\text{PPL}(D, \mathcal{M}) = \exp \left(-\frac{1}{N} \sum_{i=1}^N \log \mathcal{M}(w_i) \right) \quad (7)$$

Where $\text{PPL}(D, \mathcal{M})$ represents the perplexity of the language model \mathcal{M} on a test dataset D , N is the total number of words in the dataset, w_i is the i -th word in the dataset, and $\mathcal{M}(w_i)$ is the probability assigned by the language model \mathcal{M} to the i -th word.

A lower perplexity value indicates that the language model is better at predicting the given sequence of words.

5.1 Statistical n-grams

In this subsection, we present the results of building language models on dotted and dotless text using an n-grams-based model.

Before discussing the results of the proposed language models, we present our preparation procedure for our datasets. Table 7 shows the statistics of these datasets for each split in terms of vocabulary size.

It should be noted that we applied some selection criteria for our datasets where some samples were eliminated. For instance, in the poetry dataset, we were interested in studying classical poetry. Hence, we selected samples where the meter is known to be out of the classical poetry meters.

Furthermore, the Wikipedia dataset contains samples with various lengths. To capture meaningful samples, we only selected those with lengths greater than or equal to 30 tokens.

Figure 7 presents the n-grams counts for each order we investigated in this work. In this figure, the dotted line is the n-grams counts for the undotted dataset sharing the same color. Each row in the figure presents the results of a tokenization type with two groups of datasets presented in two plots as it is confusing to plot all in one plot.

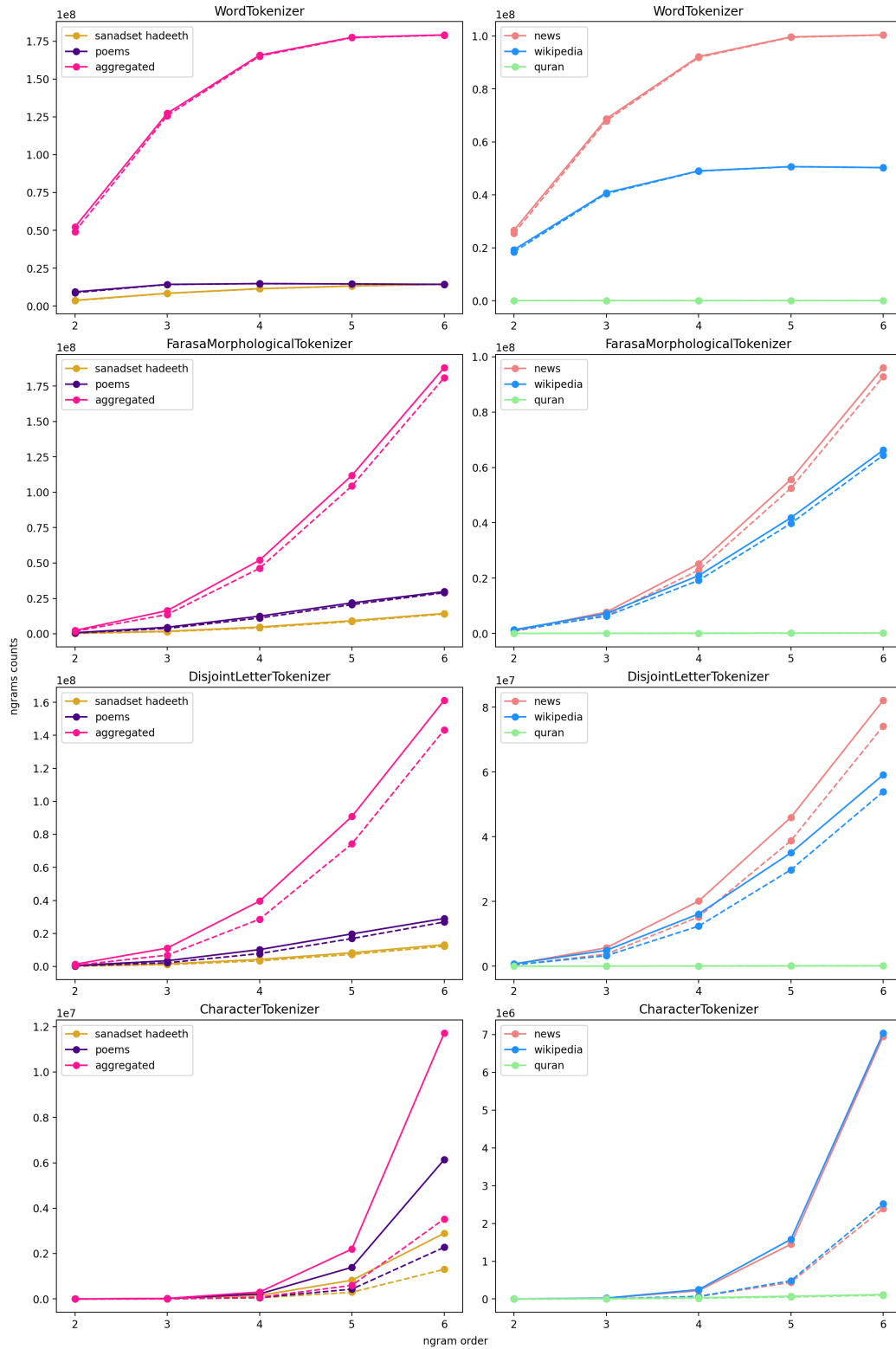


Figure 7: n-grams counts of various tokenizations applied to all proposed datasets

A general observation from this graph is that the difference between dotted and dotless counts increases along with the granularity level of the tokenization. That is, as the tokenization is more coarse-grained, the dotless n-grams counts are closer to the dotted ones.

Another interesting observation is that as the n-grams order grows, the rich datasets' counts tend to decrease faster than the non-rich datasets in both dotted and dotless texts. This can be noticed clearly in poems and Wikipedia datasets as compared to the Sanadset and the news datasets. This behavior is more apparent in coarse-grained tokenization. For instance, the poetry dataset's sixth-order counts are less than its fifth-order. This indicates that although rich datasets have more vocabulary, many of their vocabularies and, subsequently, sequences are rare as compared to non-rich datasets.

Table 8 presents the Out-Of-Vocabulary (OOVs) statistics of dotted and dotless text after splitting our dataset to training and testing splits for each tokenization scheme. The *ratio* row presents the dotless to dotted OOVs ratio to measure the reduction in OOVs of dotless text as compared to dotted. It should be noted that the last column represents an aggregated dataset with samples collected from all the datasets we used.

Tokenization		quran	sanadset	poems	wikipedia	news	aggregated
Words	<i>Dotted</i>	910	13,273	34,577	57,329	39,177	94,216
	<i>Dotless</i>	786	8,619	20,080	40,901	26,507	60,093
	<i>Ratio(%)</i>	86.37	64.94	58.07	71.34	67.66	63.78
Farasa	<i>Dotted</i>	390	4,253	12,514	32,153	17,997	51,268
	<i>Dotless</i>	300	2,948	8,112	25,104	13,680	37,430
	<i>Ratio(%)</i>	76.92	69.32	64.82	78.08	76.01	73.01
Disjoint	<i>Dotted</i>	287	2,534	5,985	7,747	5,619	12,196
	<i>Dotless</i>	185	1,078	2,216	3,116	2,252	4,417
	<i>Ratio(%)</i>	64.46	42.54	37.03	40.22	40.08	36.22

Table 8: OOVs statistics with dotted and dotless texts across different tokenizations and datasets

It can be noticed from this table that the dotless to dotted ratio resembles similar properties for rich datasets as in Tables 3, 4, and 5. For instance, It is noticed that the poems dataset has a low ratio across all tokenizations as compared to the Sanadset dataset. Wikipedia, although it has a large running text size as compared to news, has a similar ratio to news ratio. We can also note that farasa tokenization ratios are higher than other tokenizations. The results in this table seem to follow from our conclusions and hypotheses drawn in our discussion in section 4.

Figure 8 plots the perplexity of our statistical models trained on all the proposed datasets across different tokenizations with various n-grams orders. As in the n-grams counts figure, the dashed lines represent the dotless version of the dataset.

A general pattern to note in this figure is that as the n-gram order increases, the difference in perplexity between dotted and dotless text decreases. This pattern becomes more apparent as the tokenization is more fine-grained. We can also note that, as the dataset is not rich, i.e. has a smaller V/N, the dotted text perplexity is almost the same as dotless with a higher order of n-grams in almost all tokenizations. This indicates that dotless text may yield the same performance as dotted one in this task.

Unsurprisingly, as the dataset is rich, its perplexity is high. However, it is interesting to note that the perplexity of these rich datasets is notably higher than the aggregated

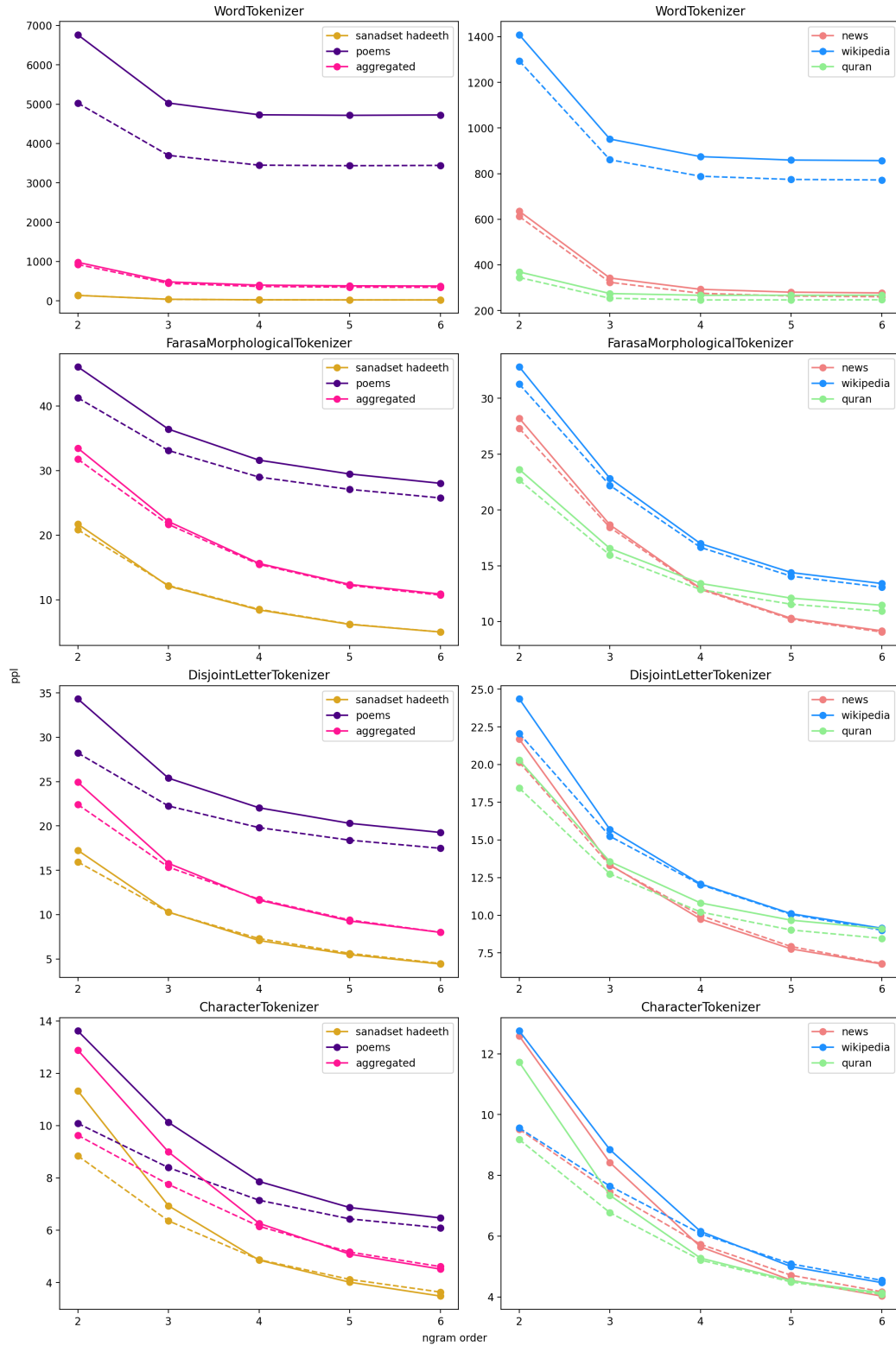


Figure 8: n-grams perplexities of various tokenizations applied to all proposed datasets

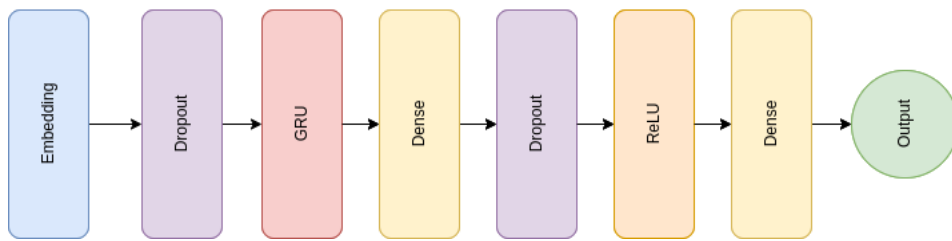


Figure 9: Neural language models architecture

dataset in all tokenizations although they dominate the aggregated dataset as they have more running text N than the non-rich datasets.

5.2 Neural language models

For the neural language model, we implemented a GRU-based model. The model comprises an embedding layer, a dropout layer (Gal and Ghahramani 2016), a layer of stacked GRU units (Chung et al. 2014), a dense layer utilizing ReLU activation, another dropout layer, an additional dense layer, and culminates in an output layer generating the model’s log probabilities. We used cross-entropy as our loss function and adam (Kingma and Ba 2014) as an optimizer. Figure 9 describes this architecture.

As we are experimenting with different datasets from different domains, with different tokenization schemes, the optimal set of hyperparameters may differ from one setup to another. For that reason, we implemented a tuning procedure that determines the best hyperparameters for a given dataset tokenized with a given tokenization scheme. It should be noted that this tuning was applied on the dotted experiment then the resulting best hyperparameters are applied to the dotless experiment. We used AHS scheduling technique (Li et al. 2018) that terminates bad trials early. We used ray framework (Liaw et al. 2018) which is a well-known framework for distributed computing. We optimized in a grid-search style on a subset of the training set the following hyperparameters:

- The number of GRUs chosen from (2, 3, 4).
- The GRU hidden size chosen from (256, 512).
- The embedding size chosen from (256, 512).
- The dropout chosen from (0.2, 0.333, 0.5)
- The initial learning rate chosen from (0.01, 0.001).

With word tokenization, or more generally any rich vocabulary dataset, the embedding and output layers parameters explode. It has been shown that both of these layers exhibit similar properties (Press and Wolf 2016), (Inan, Khosravi, and Socher 2016), (Nowlan and Hinton 1991), (Pappas, Werlen, and Henderson 2018). Therefore, they can share the same parameters resulting in a noticeable reduction in the model size. This technique is known as weight tying. We used it when applicable based on the best hyperparameters resulting from the implemented tuning procedure.

Dataset		Words	Farasa	Disjoint	Characters
Quran	PPL	179.02	22.89	15.76	5.28
	PPL'	195.69	18.54	14.90	4.66
Sanadset	PPL	33.11	5.16	4.87	2.79
	PPL'	35.85	5.31	5.11	2.74
Poems	PPL	1195.47	36.60	18.30	6.17
	PPL'	1022.35	31.4	16.8	5.6
News	PPL	166.90	7.72	7	3.17
	PPL'	179.09	8.6	7.29	3.18
Wikipedia	PPL	290.25	9.55	8.79	3.61
	PPL'	296.16	10.49	9.13	3.52

Table 9: Neural Language Models results of all datasets with different tokenizations

We reduced the learning rate on plateaus when little progress has been made. That is, we decrease the learning rate by a factor of 0.75 if no improvement was noticed after a full epoch. The training continued for 100 epochs at most unless no significant improvement with a margin greater than 5×10^{-3} is witnessed for consecutive 5 epochs. The best model that achieved the lowest validation loss is kept. We built the model with Pytorch (Paszke et al. 2019) exploiting the batteries-included features provided by Pytorch-Lightning package (Falcon et al. 2019). For datasets, we partitioned each dataset into 85%, 5%, and 10% splits for training, validation, and testing considering the batch size to be 64.

Vocabulary in the language follows Zipfian distribution (Piantadosi 2014), (Moreno-Sánchez, Font-Clos, and Corral 2016). This implies that a large subset of vocabulary has very little frequency while only a small portion of the vocabulary covers most of the running text. In our experiments, we selected vocabulary that covers 95% of the running text for word tokenization.

We fixed a sequence length for each dataset with a given tokenization. We discard the remaining text for samples larger than the specified sequence length and padded samples with smaller sequence lengths. This padding token is ignored in loss and perplexity calculations. We noticed that a tiny portion of samples of each sufficiently large dataset is significantly large. The sequence length of the maximum sample will result in lengthy training with little knowledge learned. Additionally, this issue becomes much more apparent with more fine-grained tokenization as GRU learning is less efficient for long sequences (Khandelwal et al. 2018), (Rumelhart, Hinton, and Williams 1985), (Werbos 1990). Hence, we considered different percentiles for different tokenization. We considered the length of the 0.99 percentile sample for word tokenization, 0.975 percentile sample for disjoint letters tokenization, and 0.95 percentile sample for character tokenization.

Table 9 presents the results of our neural language models applied on all of the proposed datasets with all proposed tokenizations. In that table, PPL is the perplexity of dotted text while PPL' is the perplexity of the undotted one.

Analyzing perplexity results from this table, it can be noted that rich datasets have very high perplexity. This can be seen in poems as compared to Sanadset and Wikipedia as compared to news datasets. Poems, as it is the richest dataset has the highest perplexity as compared to any other dataset across all tokenizations. This is

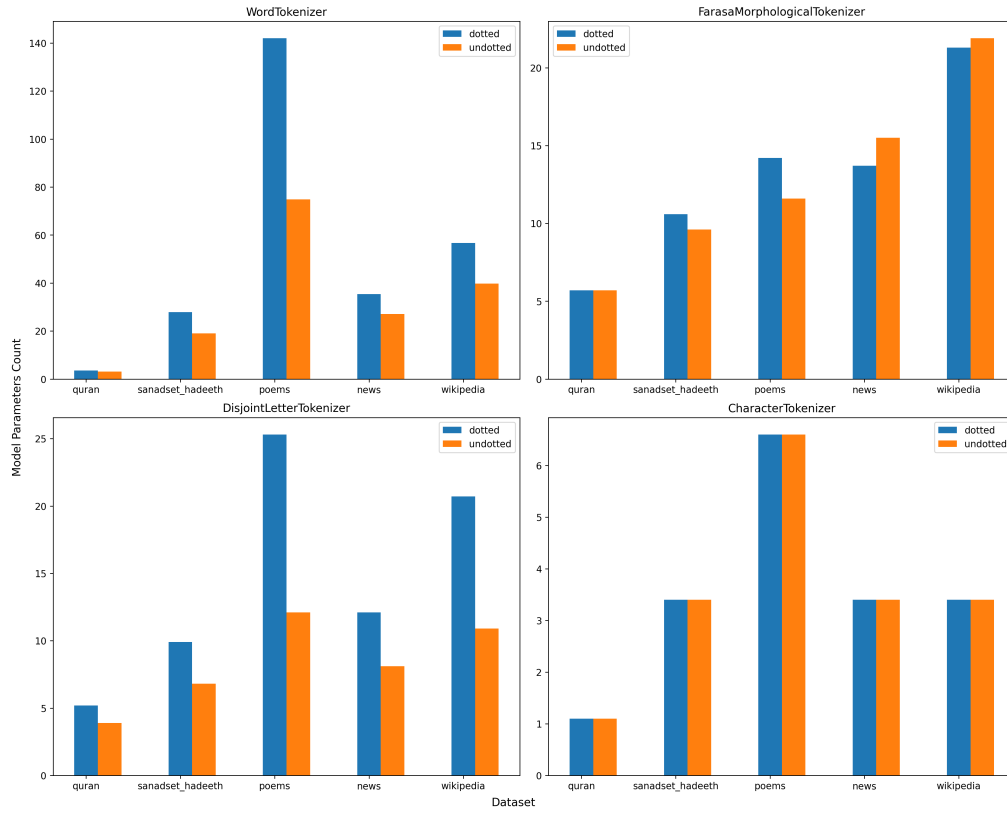


Figure 10: Neural Language Models parameters size comparison for dotted and dotless text

because rich datasets are rich in vocabulary and thus are less predictable making it harder for the learning algorithm to learn.

As the dotless text reduces the model vocabulary compared to the dotted, the embedding layer of the dotless text becomes smaller than the dotted embedding. This results in reducing the model parameters. Figure 10 plots the model parameters size of dotted text as compared to dotless on all our datasets across all tokenizations.

In this figure, it can be noted that the difference in model size is almost negligible for character tokenization. This is because the character's embedding size is generally small as compared to other layers due to the small vocabulary size. Furthermore, the difference between dotless and dotted vocabulary is insignificant.

An interesting observation arises in farasa tokenization, where, in certain datasets, the model size of the dotted representation is unexpectedly smaller than the dotless one. This trend emerges due to the ratio of dotless to dotted vocabulary decreasing as the vocabulary size increases, as illustrated in Figure 4. Such findings suggest the significance of carefully selecting the vocabulary size, particularly within the framework of this tokenization method, to effectively manage the dotless and dotted representations.

The figure shows a noticeable reduction in terms of model size, especially for rich datasets as compared to non-rich datasets. This distinction becomes more apparent in coarse-grained tokenizations. Furthermore, the model size difference is negligible for

small datasets as compared to mid-size and large-size datasets. This can be clearly seen in quran dataset.

6. Comparison with English

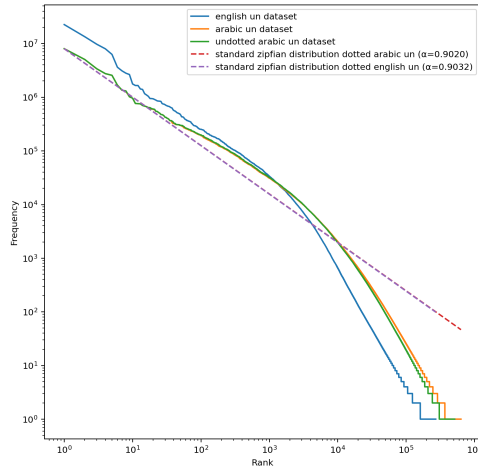
Arabic, known for its intricate morphology, exhibits a notably rich linguistic structure compared to widely spoken languages like English, resulting in a substantial expansion of the language's vocabulary. As demonstrated previously, the adoption of dotless text yields a reduction in vocabulary, consolidating multiple dotted vocabularies into a singular dotless vocabulary known as *rasm*. To explore and quantify this reduction in vocabularies, we conducted comparative analyses between dotless and dotted text within Arabic, compared against an English corpus. These comparisons were approached in two distinct manners. Firstly, utilizing a parallel corpus, we employed the UN corpus (Ziemski, Junczys-Dowmunt, and Pouliquen 2016) to facilitate this comparison. Secondly, we opted for two corpora drawn from similar domains, selecting the English Wikipedia corpus for comparison with the Arabic corpus. Notably, the English corpus substantially exceeded the Arabic corpus in terms of running text volume. To ensure a fair comparison, we curated a subset of English articles whose cumulative text size closely approximated that of the Arabic corpus. In both comparisons, the English text underwent processing by converting characters to lowercase and eliminating characters not belonging to the English character set, including numerals and punctuation. The statistics detailing these datasets are presented in Table 10.

Dataset	Char Entropy	V	N	Words Entropy
English UN	4.12	272,567	248,442,788	9.7
Arabic UN	4.19	637,778	208,325,620	12.21
Dotless Arabic UN	3.68	515,946		12.07
English Wikipedia	4.17	1,019,960	198,029,572	10.97
Arabic Wikipedia	4.27	1,811,244	177,422,512	13.2
Dotless Arabic Wikipedia	3.72	1,345,853		12.934

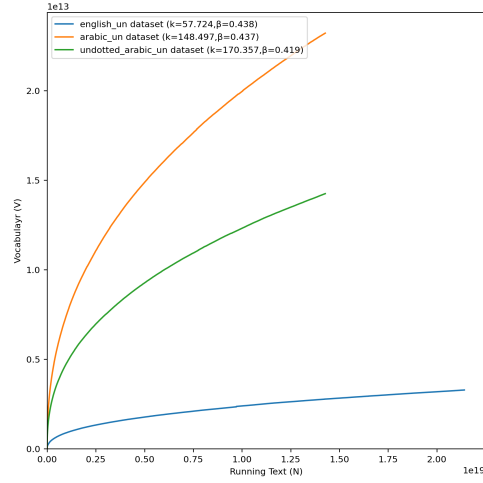
Table 10: Parallel Datasets Comparison considering dotted and dotless text

From this table, it can be seen that English has a lower character entropy than dotted Arabic, but dotless is the lowest. However, dotless Arabic has a higher entropy in words vocabulary very close to Wikipedia. English here is the lowest. This indicates that English generally is more predictable than dotted and dotless Arabic for word vocabulary. It is also interesting to note that although the UN corpus is parallel, English running text is 40M more, around 17%, in terms of running text, N.

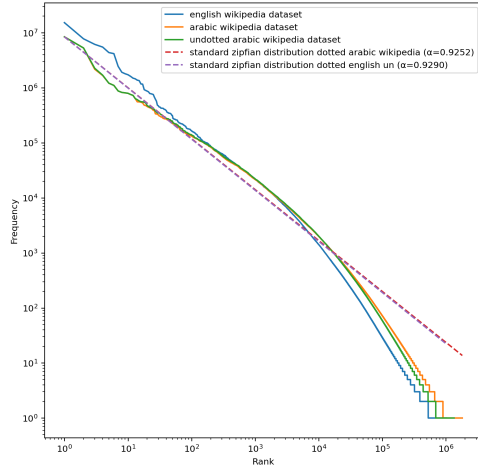
To further study vocabulary frequencies and growth, we plot Zipf's and Heap's law for both of the parallel datasets in figure 11. From this graph, we can notice that the Zipfian fit is overlapping for both datasets. It is interesting to generalize this across languages. It can be noted that in UN corpus, the English graph seems to be shifted upward for the most frequent vocabulary and downward for the least frequent vocabulary. That is, the graphs for both languages are identical if the shift effect is neglected. Following a similar pattern for Wikipedia dataset, it seems that the most frequent vocabulary for English is more frequent than Arabic. However, the least frequent English vocabularies are less frequent than Arabic.



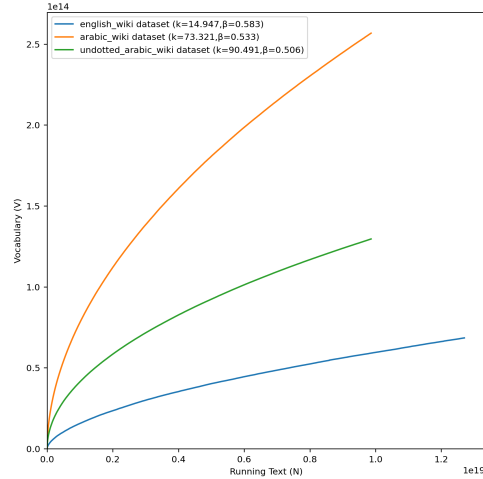
(a) UN parallel Zipf's law



(b) UN parallel Heap's law



(c) Wikipedia Zipf's law



(d) Wikipedia Heap's law

Figure 11: Zipf's and Heap's law for English-Arabic Parallel datasets

For the vocabulary growth described by Heap's law, it can be seen that Arabic is larger by order of magnitudes in terms of vocabulary although the English running text is larger for both datasets, UN and Wikipedia. It can also be noted from the graph that the Wikipedia dataset is richer than the UN dataset. Following our previous findings that the vocabulary growth for dotless text is slower for rich datasets, we can see this pattern clearly in the UN dataset as compared to Wikipedia. In fact, in Wikipedia, the dotless text growth is close to the average of both English and Arabic vocabulary growth.

7. Conclusion and future directions

In this research paper, we introduce an alternative representation of Arabic text for NLP, rooted from the idea of removing dots. Our investigation involves a comprehensive comparison between this new dotless representation and the conventional dotted Arabic text across five diverse datasets encompassing different domains and sizes, employing various tokenization schemes. As a case study within the realm of NLP tasks, we conduct a comparative analysis between dotted and dotless text specifically focused on language modeling. We performed the experimentation with two different language modeling techniques with various datasets and tokenizations. Our findings, briefly summarized below, have been extensively discussed in detail earlier in the paper.

Analyzing token statistics reveals that the ratio of dotless vocabulary to dotted varies depending on the utilized tokenization method. When considering the aggregated dataset as a representation of the language, the observed ratio of V' to V amounts to 68.1%, indicating a reduction of approximately 22% in vocabulary attributed to the removal of dots. Notably, this reduction tends to increase as the dataset richness grows. Moreover, a similar reduction pattern is observed across various levels, specifically for farasa and disjoint letters tokenization. Among these, the disjoint letters tokenization yields the most significant vocabulary reduction, reaching approximately 60% on the aggregated dataset.

All types of tokenizations presented in this work follow Zipf's law. The more coarse-grained the tokenization, the more Zipfian it is. Dotless vocabulary frequencies have similar behavior to dotted ones in Zipfian settings. This observation is dataset and tokenization-independent. Dotted and dotless vocabulary adheres to heap's law. However, dotless growth is less than dotted. This is more apparent for rich datasets. Further, disjoint growth is much slower confirming our previous findings of V'/V being the smallest.

The undotted text holds entropy close to the dotted text with a maximum difference of less than 10% in character tokenization. However, it reduces the character set by a margin of 61.3%. Based on the entropy analysis of dotted text as compared to dotted text, the dots reduced text redundancy by 5%. This implies that there are other sources of redundancy other than dots that are also present in the dotless text. The dotted text has a slightly higher entropy compared to the dotless text. This is dataset size and tokenization independent. Also, this applies to rich and non-rich datasets.

Although dotless characters' entropy is the lowest, English words' entropy is almost one-third lower than dotted and dotless Arabic words' entropy. For a parallel English-Arabic corpus, English's most frequent vocabularies are more than the Arabic most frequent ones and English's least frequent vocabularies are less than those of Arabic. This implies, the vocabulary difference is not uniform.

The difference between dotless and dotted n-grams increases as the tokenization becomes fine-grained. This supports our previous hypothesis that fine-grained tokenization covers more dotted text as compared to coarse-grained. However, this difference is generally small. Moreover, as the n-grams order grows, rich datasets' n-grams counts decrease faster than the non-rich datasets. This applies to dotted and dotless text.

Generally, for statistical language models, the dotless text perplexity becomes closer to dotted as the n-grams order grows. This observation becomes clearer as the tokenization becomes more fine-grained. Characters' perplexity for neural language models is very close for dotted and dotless text, closer than any other tokenization.

In neural language models, dotless text consistently reduces the model size due to the reduction in vocabulary across various tokenization methods, with one notable

exception found in farasa tokenization. In specific cases within farasa tokenization, the dotless vocabulary does not exhibit a reduction compared to the dotted version, due to the choice of vocabulary size. Moreover, the impact of dotless vocabulary is minimal in character tokenization, as the variance in vocabulary size between dotless and dotted representations is negligible in this tokenization scheme.

As a future work, it is interesting to study the performance of the dotless text on downstream tasks such as text classification, POS tagging and machine translation. It is also interesting to extend the use of this representation to other tasks like speech recognition and optical character recognition and compare the performance with standard text representation including dots.

Acknowledgment

The authors would like to thank Saudi Data and AI Authority (SDAIA) and King Fahd University of Petroleum and Minerals (KFUPM) for supporting this work through SDAIA-KFUPM Joint Research Center for Artificial Intelligence grant number JRC-AI-RFP-10.

References

- Abandah, Gheith A, Alex Graves, Balkees Al-Shagoor, Alaa Arabiyat, Fuad Jamour, and Majid Al-Tae. 2015. Automatic diacritization of arabic text using recurrent neural networks. *International Journal on Document Analysis and Recognition (IJ DAR)*, 18:183–197.
- Abdelali, Ahmed, Kareem Darwish, Nadir Durrani, and Hamdy Mubarak. 2016. Farasa: A fast and furious segmenter for arabic. In *Proceedings of the 2016 conference of the North American chapter of the association for computational linguistics: Demonstrations*, pages 11–16.
- Abushariah, Mohammad AM, Raja N Ainon, Rozati Zainuddin, Moustafa Elshafei, and Othman O Khalifa. 2010. Natural speaker-independent arabic speech recognition system based on hidden markov models using sphinx tools. In *International Conference on Computer and Communication Engineering (ICCCE'10)*, pages 1–6, IEEE.
- Ahmad, Irfan and Gernot A Fink. 2015. Multi-stage hmm based arabic text recognition with rescoring. In *2015 13th international conference on document analysis and recognition (ICDAR)*, pages 751–755, IEEE.
- Ahmad, Irfan and Gernot A Fink. 2019. Handwritten arabic text recognition using multi-stage sub-core-shape hmms. *International Journal on Document Analysis and Recognition (IJ DAR)*, 22(3):329–349.
- Al-Kadi, Ibrahim A. 1998. Study of information-theoretic properties of arabic based on word entropy and zipf's law. *Journal of King Saud University-Computer and Information Sciences*, 10:1–14.
- Alarifi, Abdulrahman, Mansour Alghamdi, Mohammad Zarour, Batoul Aloqail, Heelah Alraqibah, Kholood Alsadhan, and Lamia Alkwai. 2012. Estimating the size of arabic indexed web content. *Scientific Research and Essays*, 7(28):2472–2483.
- Alhathloul, Zainab and Irfan Ahmad. 2022. Automatic dottization of arabic text (rasms) using deep recurrent neural networks. *Pattern Recognition Letters*, 162:47–55.
- Alimardani, Mahsa and Mona Elswah. 2021. Digital orientalism:# savesheikhjarrah and arabic content moderation. *Alimardani, Mahsa and Elswah, Mona. Digital Orientalism:# SaveSheikhJarrah and Arabic Content Moderation (August 5, 2021). In POMEPS Studies*, 43.
- Almeman, Khalid and Mark Lee. 2013. Automatic building of arabic multi dialect text corpora by bootstrapping dialect words. In *2013 1st International Conference on Communications, Signal Processing, and their Applications (ICCSPA)*, pages 1–6, IEEE.
- Alotaiby, Fahad, Ibrahim Alkharashi, and Salah Foda. 2009. Processing large arabic text corpora: Preliminary analysis and results. In *Proceedings of the second international conference on Arabic language resources and tools*, pages 78–82.
- Alotaiby, Fahad, Salah Foda, and Ibrahim Alkharashi. 2014. Arabic vs. english: Comparative statistical study. *Arabian Journal for Science and Engineering*, 39:809–820.

- Aloufi, Khalid. 2019. Quran dataset.
- Alyafeai, Zaid and Maged Al-Shaibani. 2020. ARBML: Democratizing Arabic natural language processing tools. In *Proceedings of Second Workshop for NLP Open Source Software (NLP-OSS)*, pages 8–13, Association for Computational Linguistics, Online.
- Alyafeai, Zaid, Maged S Al-shaibani, Mustafa Ghaleb, and Irfan Ahmad. 2023. Evaluating various tokenizers for arabic text classification. *Neural Processing Letters*, 55(3):2911–2933.
- Azmi, Aqil M and Reham S Almajed. 2015. A survey of automatic arabic diacritization techniques. *Natural Language Engineering*, 21(3):477–495.
- Baevski, Alexei, Yuhao Zhou, Abdelrahman Mohamed, and Michael Auli. 2020. wav2vec 2.0: A framework for self-supervised learning of speech representations. *Advances in neural information processing systems*, 33:12449–12460.
- Bengio, Yoshua, Réjean Ducharme, and Pascal Vincent. 2000. A neural probabilistic language model. *Advances in neural information processing systems*, 13.
- Bentouati, Abdelkader. 2020. Arabic language controls before the descent of the holy quran, points - shape - sizes. *Bidayat*, 2(1):18–28.
- Blasius, Bernd and Ralf Tönjes. 2009. Zipf’s law in the popularity distribution of chess openings. *Physical Review Letters*, 103(21):218701.
- Boudad, Naaima, Rdouan Faizi, Rachid Oulad Haj Thami, and Raddouane Chiheb. 2018. Sentiment analysis in arabic: A review of the literature. *Ain Shams Engineering Journal*, 9(4):2479–2490.
- Brown, Tom, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared D Kaplan, Prafulla Dhariwal, Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, et al. 2020. Language models are few-shot learners. *Advances in neural information processing systems*, 33:1877–1901.
- Chernyavskiy, Anton, Dmitry Ilvovsky, and Preslav Nakov. 2021. Transformers: “the end of history” for natural language processing? In *Machine Learning and Knowledge Discovery in Databases. Research Track: European Conference, ECML PKDD 2021, Bilbao, Spain, September 13–17, 2021, Proceedings, Part III* 21, pages 677–693, Springer.
- Chung, Junyoung, Caglar Gulcehre, KyungHyun Cho, and Yoshua Bengio. 2014. Empirical evaluation of gated recurrent neural networks on sequence modeling. *arXiv preprint arXiv:1412.3555*.
- Darwish, Kareem, Hamdy Mubarak, and Ahmed Abdelali. 2017. Arabic diacritization: Stats, rules, and hacks. In *Proceedings of the third Arabic natural language processing workshop*, pages 9–17.
- Dehghani, Mostafa, Stephan Gouws, Oriol Vinyals, Jakob Uszkoreit, and ukasz Kaiser. 2018. Universal transformers. *arXiv preprint arXiv:1807.03819*.
- Devlin, Jacob, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. 2018. Bert: Pre-training of deep bidirectional transformers for language understanding. *arXiv preprint arXiv:1810.04805*.
- Diab, Mona, Mahmoud Ghoneim, and Nizar Habash. 2007. Arabic diacritization in the context of statistical machine translation. In *Proceedings of Machine Translation Summit XI: Papers*.
- El-Khair, Ibrahim Abu. 2016. 1.5 billion words arabic corpus. *arXiv preprint arXiv:1611.04033*.
- Falcon, William et al. 2019. Pytorch lightning. *GitHub*. Note: <https://github.com/PyTorchLightning/pytorch-lightning>, 3(6):11.
- Foundation, Wikimedia. Wikimedia downloads.
- Gal, Yarín and Zoubin Ghahramani. 2016. A theoretically grounded application of dropout in recurrent neural networks. *Advances in neural information processing systems*, 29.
- Guellil, Imane, Houda Saâdane, Faical Azouaou, Billel Gueni, and Damien Nouvel. 2021. Arabic natural language processing: An overview. *Journal of King Saud University-Computer and Information Sciences*, 33(5):497–507.
- Hamed, Osaama and Torsten Zesch. 2017. A survey and comparative study of arabic diacritization tools. *Journal for Language Technology and Computational Linguistics*, 32(1):27–47.
- Heafield, Kenneth. 2011. Kenlm: Faster and smaller language model queries. In *Proceedings of the sixth workshop on statistical machine translation*, pages 187–197.
- Heafield, Kenneth. 2013. *Efficient Language Modeling Algorithms with Applications to Statistical Machine Translation*. Ph.D. thesis, Ph. D. thesis, Carnegie Mellon University.
- Heafield, Kenneth, Ivan Pouzyrevsky, Jonathan H Clark, and Philipp Koehn. 2013. Scalable modified kneser-ney language model estimation. In *Proceedings*

- of the 51st Annual Meeting of the Association for Computational Linguistics (Volume 2: Short Papers), pages 690–696.
- Heaps, Harold Stanley. 1978. *Information retrieval, computational and theoretical aspects*. Academic Press.
- Hochreiter, Sepp and Jürgen Schmidhuber. 1997. Long short-term memory. *Neural computation*, 9(8):1735–1780.
- Inan, Hakan, Khashayar Khosravi, and Richard Socher. 2016. Tying word vectors and word classifiers: A loss framework for language modeling. *arXiv preprint arXiv:1611.01462*.
- Katz, Slava. 1987. Estimation of probabilities from sparse data for the language model component of a speech recognizer. *IEEE transactions on acoustics, speech, and signal processing*, 35(3):400–401.
- Khandelwal, Urvashi, He He, Peng Qi, and Dan Jurafsky. 2018. Sharp nearby, fuzzy far away: How neural language models use context. *arXiv preprint arXiv:1805.04623*.
- Khreisat, Laila. 2009. A machine learning approach for arabic text classification using n-gram frequency statistics. *Journal of Informetrics*, 3(1):72–77.
- Kingma, Diederik P and Jimmy Ba. 2014. Adam: A method for stochastic optimization. *arXiv preprint arXiv:1412.6980*.
- Kneser, Reinhard and Hermann Ney. 1995. Improved backing-off for m-gram language modeling. In *1995 international conference on acoustics, speech, and signal processing*, volume 1, pages 181–184, IEEE.
- Kudo, Taku and John Richardson. 2018. Sentencepiece: A simple and language independent subword tokenizer and detokenizer for neural text processing. *arXiv preprint arXiv:1808.06226*.
- Li, Liam, Kevin Jamieson, Afshin Rostamizadeh, Ekaterina Gonina, Moritz Hardt, Ben Recht, and Ameet Talwalkar. 2018. Massively parallel hyperparameter tuning.
- Li, Wentian. 1992. Random texts exhibit zipf’s-law-like word frequency distribution. *IEEE Transactions on information theory*, 38(6):1842–1845.
- Liaw, Richard, Eric Liang, Robert Nishihara, Philipp Moritz, Joseph E Gonzalez, and Ion Stoica. 2018. Tune: A research platform for distributed model selection and training. *arXiv preprint arXiv:1807.05118*.
- Linders, Guido M and Max M Louwerse. 2022. Zipf’s law revisited: Spoken dialog, linguistic units, parameters, and the principle of least effort. *Psychonomic Bulletin & Review*, pages 1–25.
- Magueresse, Alexandre, Vincent Carles, and Evan Heetderks. 2020. Low-resource languages: A review of past work and future challenges. *arXiv preprint arXiv:2006.07264*.
- Manaris, Bill Z, Luca Pellicoro, George J Pothering, and Harland Hodges. 2006. Investigating esperanto’s statistical proportions relative to other languages using neural networks and zipf’s law. In *Artificial Intelligence and Applications*, pages 102–108.
- Metwally, Aya S, Mohsen A Rashwan, and Amir F Atiya. 2016. A multi-layered approach for arabic text diacritization. In *2016 IEEE International Conference on Cloud Computing and Big Data Analysis (ICCCBDA)*, pages 389–393, IEEE.
- Mghari, Mohammed, Omar Bouras, and Abdelaziz El Hibaoui. 2022. Sanadset 650k: Data on hadith narrators. *Data in Brief*, 44.
- Moreno-Sánchez, Isabel, Francesc Font-Clos, and Álvaro Corral. 2016. Large-scale analysis of zipf’s law in english texts. *PloS one*, 11(1):e0147073.
- Mubarak, Hamdy, Ahmed Abdelali, Hassan Sajjad, Younes Samih, and Kareem Darwish. 2019. Highly effective arabic diacritization using sequence to sequence modeling. In *Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long and Short Papers)*, pages 2390–2395.
- Newman, Mark EJ. 2005. Power laws, pareto distributions and zipf’s law. *Contemporary physics*, 46(5):323–351.
- Ney, Hermann, Ute Essen, and Reinhard Kneser. 1994. On structuring probabilistic dependences in stochastic language modelling. *Computer Speech & Language*, 8(1):1–38.
- Nowlan, Steven and Geoffrey E Hinton. 1991. Adaptive soft weight tying using gaussian mixtures. *Advances in Neural Information Processing Systems*, 4.
- Pappas, Nikolaos, Lesly Miculicich Werlen, and James Henderson. 2018. Beyond weight tying: Learning joint input-output embeddings for neural machine translation. *arXiv preprint arXiv:1808.10681*.
- Paszke, Adam, Sam Gross, Francisco Massa, Adam Lerer, James Bradbury, Gregory Chanan, Trevor Killeen, Zeming Lin, Natalia Gimelshein, Luca Antiga, et al.

2019. Pytorch: An imperative style, high-performance deep learning library. *Advances in neural information processing systems*, 32.
- Piantadosi, Steven T. 2014. Zipf's word frequency law in natural language: A critical review and future directions. *Psychonomic bulletin & review*, 21:1112–1130.
- Press, Ofir and Lior Wolf. 2016. Using the output embedding to improve language models. *arXiv preprint arXiv:1608.05859*.
- Rom, Aviad and Kfir Bar. 2021. Supporting undotted arabic with pre-trained language models. *arXiv preprint arXiv:2111.09791*.
- Rumelhart, David E, Geoffrey E Hinton, and Ronald J Williams. 1985. Learning internal representations by error propagation. Technical report, California Univ San Diego La Jolla Inst for Cognitive Science.
- Sano, Yukie, Hideki Takayasu, and Misako Takayasu. 2012. Zipf's law and heap's law can predict the size of potential words. *Progress of Theoretical Physics Supplement*, 194:202–209.
- Selab, Essma and Ahmed Guessoum. 2015. Building talaa, a free general and categorized arabic corpus. In *ICAART (1)*, pages 284–291.
- Serrà, Joan, Álvaro Corral, Marián Boguñá, Martín Haro, and Josep Ll Arcos. 2012. Measuring the evolution of contemporary western popular music. *Scientific reports*, 2(1):1–6.
- Sicilia-Garcia, Elvira I, Ji Ming, Francis J Smith, et al. 2002. Extension of zipf's law to words and phrases. In *COLING 2002: The 19th International Conference on Computational Linguistics*.
- Sicilia-Garcia, Elvira I, Ji Ming, Francis J Smith, et al. 2003. Extension of zipf's law to word and character n-grams for english and chinese. In *International Journal of Computational Linguistics & Chinese Language Processing, Volume 8, Number 1, February 2003: Special Issue on Word Formation and Chinese Language Processing*, pages 77–102.
- Smith, Ray. 2011. Limits on the application of frequency-based language models to ocr. In *2011 International Conference on Document Analysis and Recognition*, pages 538–542, IEEE.
- Song, Xinying, Alex Salcianu, Yang Song, Dave Dopson, and Denny Zhou. 2020. Fast wordpiece tokenization. *arXiv preprint arXiv:2012.15524*.
- Stolcke, Andreas. 2002. Srilm-an extensible language modeling toolkit. In *Seventh international conference on spoken language processing*.
- Vaswani, Ashish, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, ukasz Kaiser, and Illia Polosukhin. 2017. Attention is all you need. *Advances in neural information processing systems*, 30.
- Welsh, Dominic. 1988. *Codes and cryptography*. Oxford University Press.
- Werbos, Paul J. 1990. Backpropagation through time: what it does and how to do it. *Proceedings of the IEEE*, 78(10):1550–1560.
- Yu, Shuiyuan, Chunshan Xu, and Haitao Liu. 2018. Zipf's law in 50 languages: its structural pattern, linguistic interpretation, and cognitive motivation. *arXiv preprint arXiv:1807.01855*.
- Zhang, Hui and David Chiang. 2014. Kneser-ney smoothing on expected counts. In *Proceedings of the 52nd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 765–774.
- Ziemski, Michał, Marcin Junczys-Dowmunt, and Bruno Pouliquen. 2016. The united nations parallel corpus v1. 0. In *Proceedings of the Tenth International Conference on Language Resources and Evaluation (LREC'16)*, pages 3530–3534.
- Zipf, George Kingsley and Human Behavior. 1949. *the Principle of Least Effort*. New York, Hafner Publishing Company.

