The Thesis committee for Juan Daniel Pinto
Certifies that this is the approved version of the following thesis:

# Creating a Conversational
# Hebrew Vocabulary List

**APPROVED BY**
**SUPERVISING COMMITTEE:**

———————————————————————

Esther Raizen, Supervisor

———————————————————————

Elaine Horwitz, Co-Supervisor

# Creating a Conversational Hebrew Vocabulary List

by

Juan D. Pinto

**Thesis**

Presented to the Faculty of the Graduate School
of the University of Texas at Austin
in Partial Fulfillment
of the Requirements
for the Degree of

**Master of Arts in Hebrew linguistics**

The University of Texas at Austin
May 2018

# Dedication

Dedicated to

# Acknowledgements

Interdum et malesuada fames ac ante ipsum primis in faucibus. Aliquam congue fermentum ante, semper porta nisl consectetur ut. Duis ornare sit amet dui ac faucibus. Phasellus ullamcorper leo vitae arcu ultricies cursus. Duis tristique lacus eget metus bibendum, at dapibus ante malesuada. In dictum nulla nec porta varius. Fusce et elit eget sapien fringilla maximus in sit amet dui.

Mauris eget blandit nisi, faucibus imperdiet odio. Suspendisse blandit dolor sed tellus venenatis, venenatis fringilla turpis pretium. Donec pharetra arcu vitae euismod tincidunt. Morbi ut turpis volutpat, ultrices felis non, finibus justo. Proin convallis accumsan sem ac vulputate. Sed rhoncus ipsum eu urna placerat, sed rhoncus erat facilisis. Praesent vitae vestibulum dui. Proin interdum tellus ac velit varius, sed finibus turpis placerat.

# Creating a Conversational
# Hebrew Vocabulary List

by

Juan Daniel Pinto

The University of Texas at Austin, 2018

SUPERVISORS: Esther Raizen, Elaine Horwitz

Indent and begin abstract here. It should be a concise statement of the nature and content of the ETD. The text must be either double-spaced or 1.5spaced. Abstracts should be limited to 350 words.

# Table of Contents

# List of Tables

# 1   Introduction

This thesis provides an in-depth look at the creation of the Conversational Hebrew Vocabulary List (hereafter CHVL)—a list of the most common words in spoken Modern Hebrew. Its two-fold aim is (1) to explore the theory behind the creation of the CHVL, along with implications for similar projects, and (2) to describe the methods and provide the tools to make the process as reproducible as possible.

The complete list itself, consisting of 5,000 items, is included as an electronic supplement and can be downloaded free of charge.[1] A partial list of the first 1,000 items can be found in *Appendix 1.*

A review of the literature will first highlight the gap that exists for less commonly taught languages (LCTLs). Because the overwhelming majority of the previous research in vocabulary frequency lists has focused on English (and a handful of other European languages), some important nuances are yet to be addressed. More often than not, the few non-English word lists that do exist, along with much of the research in vocabulary acquisition, have taken at face value some of the findings of this limited-scope research—often without questioning whether the same methodologies and conclusions should be applied to different languages.

The present paper is, therefore, an effort to partially fill that gap in order to help educators interested in creating and/or using word lists for their own classrooms, for wider dissemination, or simply for general research purposes. In doing so, it will provide an overview of some of the key decisions that must be taken into account for such a project.

The various uses of word frequency lists can be loosely classified into research applications and practical applications. Examples of research applications include traditional linguistic studies that look for common morphological patterns, corpus-linguistic studies seeking to understand language through "real world" texts, and psycholinguistic studies that explore connections between a speaker's mental lexicon and word frequency. Practical applications of word lists include curriculum and

---

[1]Supplements can be downloaded directly from the thesis archive of the University of Texas at Austin. A separate repository at GitHub also contains the complete CHVL at *https://github.com/ juandpinto/opus-lemmas.*

textbook planning for language teachers, vocabulary selection for graded readers and dictionaries, and even independent language study. Of course, some of the most influential studies straddle both sides of this divide and attempt to answer questions such as: How can vocabulary knowledge be appropriately tested and measured? What is the role of extensive reading (as opposed to intensive reading) in incidental vocabulary acquisition? What level of vocabulary do learners need in order to read extensively for pleasure? What level of vocabulary do learners need in order to succeed in an academic setting? What role does specialized vocabulary play in reaching understanding? These questions and their answers rely heavily on the creation and use of trustworthy word frequency lists. Yet due to the resources and effort required to create these lists, they are rarely found for less commonly taught languages.

The primary research question guiding this project is this: *What are the most-frequently used words in conversational Modern Hebrew?* The resulting study also addresses the following secondary research questions, which were necessary to address in order to answer the aforementioned question: *What effect does a corpus of unvocalized texts have on the identification of word families in the computerized creation of a vocabulary frequency list? What factors affect the way that boundaries are demarcated for various levels of word families in Modern Hebrew?* And finally: *What implications might these findings have for word list creation and use as it pertains to other less commonly taught languages?*

The literature review will serve as a basis for many of the important decisions taken during the creation of the CHVL. These decisions—surrounding both corpus and list creation—along with their reasoning, will be explained further in an analysis of the literature. For the sake of clarity, these decisions are listed here at the outset. They are as follows:

**Corpus design** - *Size:* - *Text types:* The corpus consists of a single text type: conversation. This is to best fit with the list's intended audience. In order to accomplish this, movie and television subtitles compose the core of the corpus. **List design:** - *Use:* The primary intended audience for the CHVL is composed of beginning-to-low-intermediate learners of Hebrew as a foreign language. It is designed for both receptive and productive language use. - *Word family levels:* The word family level that is best suited for the CHVL's intended audience is the lemma. - *Criteria:* The

CHVL was created using exclusively objective criteria, meaning that it is the product of calculations, and it was not manually tweaked in any way. The empirical criteria used were frequency and range.

Following the review of literature and explanation of theory, the process of the CHVL's creation will be explained in detail, along with findings from the project. Possible implications for other less commonly taught languages will then be discussed. Finally, the CHVL and any code used will be provided in the appendix.

# 2 Background: Review of the literature

The theoretical foundation for the creation and use of word frequency lists rests on the observation, made popular by the linguist George Kingsley Zipf in the 1930s and 40s, that if one were to create a frequency list of words in a large enough text, the first word would occur roughly twice as often as the second word, three times as often as the third word, and so on (1935, 1949).

This exponential distribution is significant because it means that a small number of words make up the bulk of a text, whereas the majority of the words occur very few times(Sorell, 2012). Paul Nation, one of the most influential scholars in the field of vocabulary acquisition, has pointed out that Zipf's Law—as it is has come to be known—can serve as motivation to language learners and teachers, since learning the most common vocabulary in a language covers so much of the communication that naturally occurs (2013, p. 34).

This observation guides the entire endeavor of word list creation and use. Though the CHVL is not sorted using raw frequency alone[2], the effect of Zipf's law can be easily seen in the listed frequencies that accompany each item on the list.

One level above this theoretical basis lie the theoretical considerations of the process that serve as the structure upon which the CHVL is built. These include corpus size and text type, general vs. specialized lists, word family levels, and objective criteria. Each of these issues will be treated separately throughout this literature review.

## 2.1 Corpus Design

Before designing a word list, a careful, clear plan must be made for the design of the corpus from which the list is extracted. The corpus must be representative of the language context that the word list wishes to analyze. Of course, it is impossible to capture all of the communications that take place in a particular language. For this simple reason, researchers must make do with an approximation of the whole: a bounded corpus of language.

---

[2]The sorting method is explained in the sections *Objective Criteria*, *Dispersion*, and *Sort and Export*.

Though the focus of this literature review is the creation of word frequency lists, the truth is that relatively few corpora have been created for this specific purpose. Most corpora have aimed at being general collections that cover the language (usually English) as a whole in an attempt to serve different theoretical and applied uses. Yet despite this broad objective, the creation of corpora has historically revolved around two big questions: (1) how large should the corpus be, and (2) what kinds of texts should it include. These questions are important not only for corpus creation, but also for corpus selection. Both of these points will be addressed here, with the recurring emphasis being corpus use for word list creation.

### 2.1.1 Corpus Size

Conventional wisdom in corpus creation states that more is better. If a word list is to accurately reflect the frequencies of words in the language as a whole, then a corpus must contain enough text to approximate the overall use of discourse. This line of thinking is equivalent to the maxim in quantitative research that a sample should be as representative of the target population as possible. And in order to maximize the statistical probability of this representation, the sample must be of an appropriate size for the study.

True, larger sample sizes often increase this probability, but they also tend to be more resource-intensive for the researcher. The same is true of corpus size. When creating a vocabulary list, then, what is an "ideal" corpus size?

Corpora composed of millions of tokens are easy to access today. This is especially true of corpora of written material—corpora of spoken language are still comparatively small. And thanks to advances in computing power, it is finally becoming plausible for more researchers without access to extensive resources to use these mega-corpora for the purpose of word list creation.

The first project to create a one-million-token corpus was a joint effort by Henry Kučera and W. Nelson Francis of Brown University to compile a corpus of American English texts printed in 1961 (Kučera & Francis, 1967), known today simply as the *Brown Corpus*. They strived to create a corpus with equal amounts of texts from different sources by randomly selecting 500 passages of 2,000 words each from

different published materials found at the Brown University Library and the Providence Athenaeum. This mixed design would be used as a model by many of the corpora created during the next few decades: . These began to be compiled at increasingly faster rates. Many of these corpora were created—in part—to serve as parallel corpora of different varieties of English.

As an example of how quickly corpora have grown in recent decades, consider the history of COBUILD. What began in 1980 as a collaboration between Collins Publishing and a group of researchers led by John Sinclair—the Collins Birmingham University International Language Database (COBUILD)—led to the creation of the *Collins Corpus* of 7-million-tokens by 1982. It continued expanding until transforming into the *Bank of English* in the 1990s, which reached 320 million words in 1997. In 2005, as part of the Collins World Web, which also comprises French, German, and Spanish corpora, it reached 2.5 billion words (*Collins Cobuild English grammar*, 2005). The Collins Corpus now contains over 4.5 billion words ("The history of Collins COBUILD," n.d.).

Today, with the use of web-crawling applications that scour the internet and collect text at unprecedented speed, we can now use the *enTenTen12* corpus of 12 billion English tokens, which was collected in 12 days (Jakubíček, Kilgarriff, Kovář, Rychlý, & Suchomel, 2013)!

Clearly, then, the sky's the limit when it comes to ever-growing corpora of language. But when it comes to word list creation, is there a corpus size that can be considered sufficient?

Studies have approached this specific problem of corpus size for word list creation by creating multiple frequency lists—each from a different size of corpus—and then comparing the efficacy of these lists themselves. But what makes a word list effective? Different researchers have tackled this question differently.

One way to judge the effectiveness of a word list is to find how closely it correlates with reaction times in a lexical decision task—a widely-used procedure in psychological and psycholinguistic research. In a lexical decision task, participants are presented with a series of words and non-words, one after the other, and they are asked to judge which is which as quickly as possible. The reaction times are then

analyzed for each word. The basic idea behind this experiment is that the average time it takes participants to react to a word reflects something about the way the word is accessed in participants' mental lexicons. Given enough data, it is possible to make certain inferences about the arrangement of this internal lexicon, which has led to various psycholinguistic theories over the years. But what does this have to do with words on a frequency list? There is well-attested evidence to suggest that there exists an inverse correlation between word frequency and reaction time on a lexical decision task (Whitney, 1998; Balota and Chumbley, 1984). In other words, more common words are accessed and recognized more quickly than less common words. Therefore—the thinking goes—an effective word frequency list should correspond to and reflect this reality.

This was precisely the approach taken by Brysbaert & New (2009), who compared respond times collected as part of the massive Elexicon Project (Balota, et al., 2007) to words on a series of frequency lists made from increasingly larger corpora. The corpora used were all subcorpora extracted from the British National Corpus (BNC). With each subsequent increase in token count, the word list correlated more and more closely with the response times from lexical decision tasks. This observation validates the line of thinking described at the beginning of this section regarding the need for large corpora. Brysbaert and New hoped to find an "ideal" corpus size after which the increase in effectiveness would no longer be significant enough to justify the additional cost of resources. After conducting several regression analyses on the two sets of data, they found that the variance in the response times that could be accounted for by corpus size reached a plateau at about 16 million words. In other words, for corpora with less than 16 million words, the size of the corpus had a significant effect on the correlation between word frequencies and average response times for those words on lexical decision tasks. For corpora with more than 16 million words, the effect of increasing corpus size became considerably more subtle. In the end, they concluded that in order to construct an effective word list for *high-frequency* words, a corpus of about 1 million tokens is needed. However, in order to reach the same effectiveness for *low-frequency* words, a corpus size of at least 16 million words is preferable.

A different, more straightforward methodology is to directly compare word lists made from corpora of different sizes. Rather than judging the "effectiveness" of a list, this

approach measures similarities shared between different lists. Hypothetically, doing this at increasing corpus sizes should allow one to find a size after which the variance between lists only minimally decreases. As with the previous approach, the goal here is to find a point at which the benefits of increasing size no longer outweigh the additional needed resources.

Essentially, then, all corpora of sufficient size should result in nearly the same word frequency list—a theory based on a strict interpretation of Zipf's law applied to all natural language. If the appropriate criteria can be found—Sorell (2013) suggests—then this would, at last, provide a solution to Nation's (2001, 2013) observation that, problematically, word lists tend to disagree rather drastically on both the words included and their respective ranking.

Inspired by the computational linguistic measure of *rank distance* (Popescu and Dinu, 2008)—a method for comparing stylistic differences between texts—Sorell (2013) developed a variant of this methodology. First, he used different corpora of the same size to create multiple word lists, one for each corpus, ranked entirely by frequency. He then identified the percentage of words that are *not* shared between each set of two lists. Finally, he averaged these percentages to find the level of variability created at that specific corpus size. The levels of variability he found were remarkably close to each other—despite using a wide variety of entirely different corpora (with no overlap on texts within each one). He then increased the size of each corpus and repeated the process.

In order to calculate this level of variability, Sorell used a modified version of a complex formula that he borrowed from the natural sciences, and called his resulting calculation the *Dice distance.* Though this Sørensen–Dice coefficient that he altered (also known by other names) is widely used in botany and other fields to measure similarity in areas and samples of different sizes, the frequency lists measured by Sorell were all purposefully of the same size. What this means is that—apparently without realizing it—his *Dice distance* was ultimately just a simple percentage: *number of different words between frequency lists / size of frequency list (total words).* Regardless of the round-about way he used to calculate it, his resulting measure for each corpus size—the level of variability—can be accurately described as the average proportion of difference for word lists at that particular corpus size.

Sorell found that a stable list (about 2% variation) of the most frequent 1,000 words, or a reasonably stable list (less than 5% variation) of the most frequent 3,000, words can be created using a corpus of 50 million tokens. In other words, 1,000-type word lists created from different 50-million-token corpora will likely only differ by 20 words. At the 3,000-type level using the same sizes of corpora, the lists will likely vary by less than 150 words. This is a remarkable level of similarity. Expanding the list to 9,000 types will still only have about 4–7% variation, or 360–630 words. Even corpora of 20 million tokens can be considered sufficient in many cases, since they will result in 3,000-type word list with roughly 5% variation, and 9,000-type word list with less than 10% variation.

Taking a similar approach, though with significant variations, Brezina and Gablasova (2015) compared four corpora of various sizes: The Lancaster-Oslo-Bergen Corpus (LOB), The BE06 Corpus of British English (BE06), The British National Corpus (BNC), and EnTenTen16. These corpora had respective token sizes of 1 million, 1 million, 100 million, and 12 billion. The word list created from each corpus was, in this case, a combination of frequency and dispersion—a measure that will be discussed in more detail later in this paper. The resulting word lists were then compared, and the percentage of shared vocabulary words calculated. Additionally, the researchers also calculated the correlation between the ranking for each word that was shared between word lists. Contrary to Sorell, Brezina and Gablasova considered this final comparison an important part of understanding the effect of corpus size.

The aim of this study was not to find a corpus size after which the difference was negligible, but rather to find if there was a significant difference between word lists made from corpora of different sizes. The study found a 78%–84% overlap between each of the 3,000–lemma word lists. 71% of the words were shared among all four of the lists. Based on this number, Brezina and Gablasova concluded that regardless of corpus size—at least for anything larger than one million tokens—"similar results" are obtained.

This conclusion differs significantly from Sorell's, who concluded that a corpus of at least 20 million tokens (though 50 million is preferable) is needed for a stable word list with low variability. These disagreements are primarily the result of a difference in what should be considered "stable." At 71% vocabulary overlap—which is sufficient

for Brezina and Gablasova—870 words were only found in one of the four lists. This is drastically higher than Sorell's threshold, which at the 3,000-word level varies in roughly 150 words. Note that Nation and Hwang (1995) found a level of overlap similar to Brezina and Gablasova when comparing the GSL, the LOB, and the Brown corpora—a percentage of overlap that they deemed to be not particularly high. As Nation later put it, "Brezina and Gablasova are a bit too tolerant in accepting that 71% or even 78%-84% overlap is good enough. If roughly one out of every four or five words is different from one list to another, that is a lot of difference" (2016, p. 100).

Another difference to mention between these two studies is the unit of counting used. Sorell made lists based on *types*, whereas Brezina and Gablasova preferred the use of *lemmas*. I will explain this important distinction in a later section of this review ("Identifying Words"). For now, it is sufficient to say that the effect of these different measures in comparing word lists created from corpora of different sizes has (to my knowledge) not been studies. This is one area that could benefit from further research.

Lastly, the corpora used by Brezina and Gablasova were all-inclusive: each built on its own philosophy on the way that different types of texts should be balanced in a corpus, but all seeking to be representative of English as a whole. This is also true of the corpora used by Brysbaert and New in their study using response times from a lexical decision task. Contrast this with Sorell's word lists, which were systematically created from corpora that consisted of only one specific text type. Surely, this is a factor to consider in corpus design.

Therefore, having a sufficiently large corpus is important, as demonstrated in this section. But is it enough? How much do the types of texts included in a corpus factor into its effectiveness for word list creation?

### 2.1.2  Text Types

There's been a lot of debate about the "best" way to balance a corpus' text types. This is a major aspect of corpus design, and one worth delving into. At the end of the day, much of it comes down to the purpose of the corpus. When used for the

creation of word lists, one must also consider the intended purpose of the word list itself. Is it for general use or for one of many possible specialized uses? More on this in the next section.

In order to design a corpus with different amounts of text types (i.e. narrative, conversational, academic), clear definitions for these text types are necessary. But is there a better way than the use of subjective genres to classify texts?

Or is there a better methodology than simply mixing a bunch of different texts together, with the hope that the resulting word list covers the language as a whole? This is the most common way of creating frequency lists, but it tends to result in a mix of words that have little relevance to any one purpose. Esoteric, academic words in a beginners' vocabulary list? Science fiction terms in a vocabulary list for business managers? It's obvious that a list is only as good as the corpus from which it's made, which is why a clear delineation of different text types and their qualities is critical.

When speaking of corpus balance, I refer to the proportion of different text types that make up a corpus. Published corpora have taken different approaches in this regard, and published word lists have made use of a variety of strategies for balancing the corpora from which they are made. Coxhead's *Academic Word List* (2000) was created from a carefully-designed corpus that used equally-sized sub-corpora of texts from different disciplines. This suited the purpose of her word list well, since it was intended to serve students from a variety of disciplines.

The importance of identifying a taxonomy of text types based on objective criteria: are there distinguishable linguistic differences between an informal correspondence and a narrative work of fiction? What about between a romance and a fantasy novel?

Biber's early work (1988) conducted an analysis of a wide variety of texts using large corpora to tag syntactic markers and other linguistic attributes that could potentially be used to define different types of texts. In this study, he found a series of five categories (each consisting of two opposite ends of a spectrum) in which texts varied: involved vs. informational, narrative, situated vs. elaborated, persuasive, and abstract. He then conducted a very large study, which he published as a book, (1995) that found eight distinct, recurring patterns of different combinations of these

categories. These groupings serve as a linguistically-based taxonomy that divides texts along objective lines, rather than subjective, culturally-defined genres.

Similar but independent studies were conducted for Somali, Korean, Nukulaelae Tuvuluan, Taiwanese, and Spanish (Biber, 1995; Jang, 1998). For each language, a unique set of text types were identified. However, the texts were found to align along similar distinguishing linguistic dimensions as the English texts.

Sorell (2013) sought to simplify Biber's eight text types into categories suitable for corpora study. He did this by noticing the closely similar ways that some of the text types lined up along Biber's five linguistic categories, also incorporating some extra-linguistic features, such as shared contexts (e.g. predominantly spoken types). He also dropped Biber's two smallest text types, deeming them impractical for corpus study and difficult to isolate. In doing this, he came up with four simplified text types: interactive (conversation), general reported exposition (general writing), imaginative narrative (narrative writing), and academic. Regarding this last type, Biber's study found a sosignificant difference between academic writing in the natural sciences ("scientific exposition") and the humanities ("learned exposition")—he found that natural science uses more concrete language, whereas the humanities tend to use more abstract language. However, Sorell sought to unify these for the sake of simplicity, simply leaving their distinction to "a future study" (p. 68). Sorell acknowledged that his wasn't the first attempt at simplification of Biber's text types, a surprisingly similar effort having been made in the *Longman Grammar of Spoken and Written English* (Biber, Johansson, Leech, Conrad, & Finegan, 1999: 16) and the *Longman Student Grammar of Spoken and Written English* (Biber, Conrad, & Leech, 2002: 23).

Sorell found that each of his four simplified text types yielded a vocabulary frequency list that was as unique as the linguistic criteria that Biber had used. He also measured how different they were from each other, and found all four to be equidistant from the next in this order: conversation, narrative, general writing, and academic writing (See section on corpus size for an explanation of this measurement). Sorell, therefore, claims that his own study of vocabulary frequency using his simplified text types as a base has "validated Biber's studies by adding a vocabulary dimension to the description of each of the key text types" (201).

Despite the importance of spoken language—or the conversation text type—for language learners and linguistic studies, the number of conversation corpora that exist, as well as their size, is very limited. This is clearly because of the difficulty of gathering large amounts of spoken data that then needs to be transcribed by hand in order to be analyzed. It is true that speech recognition software has come a long way in recent years, but its rate of error remains too high for research purposes. It has been estimated that it takes 40 hours to professionally transcribe one hour of audio recording, making the task too costly. For this reason, some researchers have begun looking at alternative sources for a conversation corpus, including the internet and movie subtitles.

New, et al. (2007) created a 50-million-token corpus of French subtitles. They divided this into four subcorpora, one for each of the type of media from which the subtitles were extracted: French films, English movies, English television series, and non-English-language European films. The reason for using French subtitles from English media is the sheer dominance of English in the film industry. In order to counter-balance the much larger sizes of the two subcorpora extracted from English media, the researchers measured word frequencies for each subcorpora separately, then averaged them to arrive at the final frequency used for their ranked word list.

In order to test the validity of their new approach, New, et al. used two different methods. First, they compared their subtitle word list with word lists created from more traditional corpora. Second, they used lexical decision times—similar to Brysbaert and New (2009) above—to test the rankings of words on their list.

The first test found a .73 correlation with a classical French spoken corpus, the "Corpus de Référence du Français Parlé" (CRFP; Equipe DELIC, 2004). However, when looking at the specific words and semantic categories that differ the most, it's clear that most major differences are caused by the monologue-nature of the CRFP. This corpus was created from a large number of interviews (each asking the same questions to the interviewee), whereas movie subtitles tend to be composed primarily of people interacting in conversations. This results in more colloquial expressions having higher frequencies in the subtitle corpus. The nature of movies themselves also played a role, resulting in an overrepresentation of words related to action movies and police matters—words like *tuer* [to kill], *prison* [jail], and *armes*

[weapons] (p. 665).

For the second test of the subtitle word list, the researchers used the lexical decision times from two previous experiments. They found that the subtitle list's ability to predict lexical decision times was at least equally as accurate as the CRFP frequencies or those from a traditional corpus of written French. In many cases, it actually fared much better, surprising even the researchers themselves. However, this latter test was based on the rather small sample sizes of the two previous experiments (234 and 240 words), limiting the reliability of this test.

Picking up on these findings, and expanding the lexical decision task to a much larger sample size, Brysbaert and New (2009) compiled a corpus of English subtitles (SUBTLEX$_{US}$) and evaluated it as part of their study. This corpus is composed of subtitles from a wide variety of American films since 1900, though a majority are from 1990, as well as a large number of American television series. They found that the subtitle frequencies were especially good at predicting the lexical decision times of short words, often surpassing the accuracy of rankings based on the many written corpora they tested. It had more difficulty explaining the response times of longer words, which are more rarely found in film than in literature. Overall, their own conclusion confirmed that of the New, et al. (2007) study, that word frequencies derived from subtitle corpora seem to have a clear advantage over other types of corpora.

Though these two studies arrive at the same conclusion regarding the use of subtitles, more research is needed in this area. If, indeed, subtitles can be considered as appropriate sources for corpora of the conversation text type, their availability will open many possibilities previously made nearly impossible by the difficulty of the collection medium.

## 2.2   LIST DESIGN

Perhaps even more complex than appropriately designing the corpus from which to extract vocabulary for a word list, researchers have found a wide range of variables that play a role in the design of the list itself. Questions addressed in the literature deal with the difference between a general service list and a specialized list, differences

in the way that a "word" is defined and measured, different ranking criteria used, and the influence of subjective criteria on list creation, among other issues.

## 2.2.1   General Use vs. Specialized Use

Nation (2016) emphasized the importance of identifying the purpose of a word list before beginning the creation process. He believes that the main purpose of most general-use lists is to select vocabulary that language learners should learn during their first years of study. Though this may be the stated goal of some general-use lists, it is clear that they in fact serve a wide variety of purposes. He rightfully suggests, however, that the goal of serving language learners is far too broad to be very helpful. Language learners come to the task at different ages, with different language needs, and with different reasons for learning the language. A word list that is useful for adult learners intent on attending university will likely not be helpful for young leaners whose language focuses on animals, colors, and other age-appropriate material. And yet general-use lists are far more common than specialized-use lists. This is largely due to attempt at finding the language's core vocabulary.

The majority of word lists in use attempt to describe the vocabulary of the language as a whole. They are designed to be broad and all-encompassing so that they can serve any number of uses and scenarios. Essentially, they are lists that are created for general use. This broad nature of general use lists is reflected in the name of the most widely-used word list, West's *General Service List* (1953). Others include Nation's BNC/COCA lists, Browne's *New General Service List* (2014), Brezina and Gablasova's *New General Service List* (2015), and Dang and Webb's *Essential Word List* (Nation, 2016).

Another way of understanding general-use lists is that their objective is to find what is often termed the *core* vocabulary. Though not always explicitly stated, the philosophy behind this approach is that the language being used—usually English—has at its center a self-contained lexicon of essential, primary, basic, fundamental vocabulary that then runs through the entire language. There are layers of frequency and increasing complexity beyond this, with regions of specialized language demarcated for specific purposes such as fields of study or external dialects. Still, this core vo-

15

cabulary is at the center of it all, and the purpose of a word list is to identify what words fall within its boundaries. Sorell (2013) evaluated a number of definitions of core vocabulary found in the literature. He suggests that general use lists, such as West's GSL, serve as intuitively-selected lists of core written communication, whereas survival vocabulary lists—often found in travel guides or similar materials—are core vocabulary lists of oral communication.

Relatively fewer researchers have created word lists aimed at a more specific purpose or target audience. Specialized-use lists can be designed to only include words that belong to a specific domain, such as a discipline or trade. They can also encompass vocabulary found in a broad range of disciplines, but which are common in a specific context, such as academic texts. In this case, they usually serve as supplements to aid language learners who are already familiar with the core vocabulary of the language.

Perhaps the most well-known example of a specialized-use list is Coxhead's Academic Word List (2000), which replaced the University Word List (Xue & Nation, 1984) as the go-to vocabulary list for aspiring students intent on attending an English-speaking university or those entering the academic world. This is considered a *general* academic word list, since it is for academic use in general, and not for a specific discipline.

More specialized lists include those designed for business English courses, or medical English courses. This is sometimes designated *technical vocabulary*. Nation (2016) explains that technical vocabulary is most often taught after students have mastered general-use vocabulary, and after they have some familiarity with academic vocabulary. Chung and Nation (2003) looked into the nature of a technical vocabulary. By studying specialized words in the fields of anatomy and applied linguistics, they found that a large number of technical words are also found in the language's core vocabulary, or have a general academic use as well. However, when used in a technical text, these words take on a specialized definition that is particular to that domain. This means that much vocabulary is shared across layers of vocabulary, though they may vary semantically, based on context.

## 2.2.2 Identifying Words (Word Family Levels)

One of the most essential questions that needs to be answered when designing a word list is how one is defining a *word.* Though this may seem like a straight-forward decision, it requires thorough planning and a solid understanding of the theory behind the decision. Should *jump* and *jumped* be counted as two different words or just one? What about irregular inflections such as *go* and *went*? In an article aimed at raising awareness of what he calls the "*Word* dilemma," Gardner (2007) points out that the validity of much vocabulary research hinges "on the various ways that researchers have operationalized the construct of *Word* for counting and analysis purposes" (2007, p. 242).

The literature has generally come to accept some key terms that are helpful when speaking of the way words are counted. Beginning with the most basic measurement and progressing to the most complex, we can choose to count tokens, types, lemmas, or word families.

Measuring *tokens* means simply measuring the total number of words. The sentence "I like small dogs, big dogs, and every other kind of dog" contains twelve tokens—twelve words in total. Counting *types* refers to the number of separate and distinct words. That is, *dog* and *dog* are the same type, but *dogs* is a different type—even a single difference makes them different types. The sentence above is composed of eleven types. A level above this, the *lemma* includes the stem of the word and its inflected forms, but not any derived forms of the word (derived forms are usually considered a different part of speech). So *do*, *does*, and *did* are all the same lemma, but *doable* is not. This is because *doable* has the derivational affix *-able*, which turns it into an adjective. Francis and Kučera define lemma as "a set of lexical forms having the same stem and belonging to the same major word class, differing only in inflection and/or spelling" (1982, p. 1).

Finally, the term *word family* is used to describe an even more inclusive level than the lemma. However, its precise definition has often varied among researchers. Bauer and Nation (1993) sought to rectify this problem through an in-depth classification of English affixes. Borrowing from Thorndike's (1941) study of English suffixes, their grouping was based on a series of eight criteria: frequency, productivity, predictabil-

ity, regularity of the written form of the base, regularity of the spoken form of the base, regularity of the spelling of the affix, regularity of the spoken form of the affix, and regularity of function. (pp. 255–56) They identified seven "levels" of word families, with each successive one including a larger number of affixes, and therefore a larger number of types per word family. One very useful aspect of their particular system is that it places all the previous levels (type, lemma, etc.) within the same framework. Under their schema, a level 1 word family is the same as a type, a level 2 word family is a lemma (including all regular inflected affixes), and level 7 (the highest level) consists of classical roots and affixes beyond what most speakers any longer consider separate affixes.

Nation himself suggests that for the purposes of language learning, these specific family word levels can be used simply "as a starting point as an initial framework of reference" (2016, p. 36). That is, they are one interpretation of how to systematically count words for a frequency list. These levels are based on criteria that reflect the needs of language learners, rather than on any psycholinguistic theory of how speakers' mental lexicon is arranged. Still, the idea of word families aligns closely with theoretical models that dictate morphological decomposition as a constant. These theories propose that words are often deconstructed into independent morphemes in receptive tasks and recognized that way, for example by deconstructing *jumping* into *jump* and *-ing*. At the other end of the spectrum stand theories that would place *jump* and *jumping* as separate lexical entries (Brysbaert and New, 2009, 982–83).

Either way, there is strong evidence to suggest that inflected/derived forms and their base forms do affect each other in some way, suggesting that word families are a measure of a real representation in speakers' mental lexicon. In one such study, Nagy et al. (1989) explored the effect of both inflectional and derivational family frequency during a lexical decision task. They found that both types of morphological relationships lowered word recognition times, leading to the conclusion that inflections and derivational relationships are both represented in the mental lexicon, either through the grouping of related words under the same entry, or through linked entries. However, all the participants were native English speakers, so to what extent do L2 learners' lexicons reflect the same level of linking?

More recent studies have found that L2 learners' morphological knowledge and word-

building ability are not nearly as developed. Ward and Chuenjundaeng (2009) conducted a study that tested the receptive ability of Thai engineering and doctoral students learning English. They were tested for their knowledge of a series of base words, together with various derived forms of the same words. They found a surprising lack of familiarity with the derived words, even when participants knew the base forms from which they were derived. Similarly, but from a productive and not receptive standpoint, Schmitt and Zimmerman (2002) found that learners could produce only a limited number of derived forms when presented with a word family headword. These results challenge the common assumption that "once the base word or even a derived word is known, the recognition of other members of the family requires little or no extra effort" (Bauer and Nation, 1993, p. 253).

There is evidence (Mochizuki and Aizawa, 2000; Schmitt and Meara, 1997) to suggest a positive correlation between vocabulary size and morphological knowledge. If this is the case, then using higher-level word families in Bauer and Nation's framework for word list creation (as is the case in ), may not be appropriate for learners with limited knowledge of vocabulary—the very learners that many of these lists target.

Similarly, a study by Jeon (2011) found that L2 learners' morphological knowledge leads to greater reading comprehension. Since many word lists are designed to increase reading comprehension in learners, it follows that they will likely be used by students without strong word-building abilities.

Clearly, then, when it comes to creating a word list, the unit of counting needs to fit the purpose and target audience of that list. Brezina and Gablasova (2015) contend that Bauer and Nation's (1993) higher word family levels ignore the lack of transparency that exists between many of the entries that would be placed under the same word family. Especially when creating a word list for beginners, Brezina and Gablasova point out that the morphological knowledge of language learners is often not developed enough. Because their New General Service List was created for beginners, and since it is intended to aid vocabulary acquisition for both receptive and productive purposes, Brezina and Gablasova chose the lemma as their unit of measure.

Seeking to quantify the effect of choosing to measure word families as opposed to word types, Sorell (2013) compared the text coverage of frequency lists made from

the same four corpora. Each corpus corresponded to one of Sorell's text types (see above). Sorell's definition of "word families" was a slightly modified version of Bauer and Nation's (1993) sixth level of affix inclusion. He found, as would be expected, that the most frequent word families have a much larger text coverage than the most frequent types. This is especially true when measuring type coverage—the most frequent word families accounted for roughly 4–6 times as many types in each corpus. However, when measuring overall token coverage, the top word families only covered about 3–10% more than the same number of most frequent types. Sorell also found that the most frequent 1,000 word families consisted of 6,557 word types in the general writing corpus. The number was similar in the other text types, though somewhat lower.

### 2.2.3   Objective vs. Subjective Design

(Nation 2016:133) >There are two major approaches to making corpus-based word lists. One is to stick strictly to criteria based on range, frequency and dispersion (Brezina & Gablasova, 2015; Dang & Webb, Chapter 15 this volume; Leech, Rayson & Wilson, 2001). The other is to use a similar statistical approach but to adjust the results using other criteria such as ensuring that lexical sets such as numbers, days of the week, months.

Brezina and Gablasova (2015), p. 3: > Seen from the perspective of current corpus linguistic research (cf. McEnery and Hardie 2011), one of the main problems of West's GSL lies in the fact that its compilation involved a number of competing principles that brought a large element of subjectivity into the final product. When reviewing the compilation principles of the GSL, we can see that in addition to the quantitative measure of word frequency, West also used a number of 'qualitative' criteria for the selection of individual lexical items. These include (i) the ease of learning, (ii) necessity, (iii) cover, and (iv) stylistic and emotional neutrality (West 1953: ix–x). Let us now briefly discuss these principles.

### 2.2.4 Objective Criteria (Frequency, Range, Dispersion)

Nation (2016), p. 103: > Dividing a corpus into sub-corpora allows the creation of range and dispersion figures. In some ways range figures are more important than frequency figures, because a range figure shows how widely used a word is, and this indicates its "general service". Brysbaert and New (2009) found that a range measure was a good predictor of lexical decision times. Carroll, Davies and Richman (1971) found in their study that frequency and their measure of dispersion correlated at .8538 (page xxix), showing that the more widely used a word is, the more likely it is to be frequent. Some words however are frequent in just one or two texts or sub-corpora and may not even occur in others. The use of a range or dispersion figure or both can indicate such words.

Brysbert and New (2009), pp. 984–5: > Another variable that has been proposed as an alternative to WF frequency is the contextual diversity (CD) of a word (Adelman, Brown, & Quesada, 2006). This variable refers to the number of passages (documents) in a corpus containing the word. So, rather than calculating how often a word appeared in the BNC, Adelman et al. measured how many of the 3,144 text samples in the corpus contained the word. They found that the CD measure explained 1%–3% more of the variance in the Elexicon data.

Brezina and Gablasova (105), p. 8: > ARF is a measure that takes into account both the absolute frequency of a lexical item and its distribution in the corpus (Savicky´and Hlava´c ˘ ova´2002; Hlava´c ˘ ova´2006). Thus if a word occurs with a relatively high absolute frequency only in a small number of texts, the ARF will be small (cf. Cerma´k and Kr ˘ en 2005; Kilgarriff 2009). All four wordlists were then sorted according to the ARF that ensured that only words that are frequent in a large variety of texts appeared in the top positions in the wordlists.

Sorell (2013), p. 89: Dispersion.

### 2.3 Modern Non-English Word Lists

Gardner, D. (2007), p. 242: > Hazenberg and Hulstijn 1996—Dutch language;

# 3 Methods: Creating the Conversational Hebrew Vocabulary List (CHVL)

As we have seen, the brunt of the work in high-quality vocabulary frequency list creation has focused on *English* frequency lists. Outside of the English-speaking world, and especially when dealing with less commonly taught languages, it's difficult to find well-researched word lists, if they exist at all. Why have not more educators—those who may benefit from these lists the most—decided to undertake such a task?

This need not be a project that one starts from scratch every time. Many tools already exist to make the process smoother. Still, with the rapid pace at which technology changes, these tools tend to quickly become obsolete. They are also usually restrictive to the specific preferences of their creators.

Rather than using these tools, I chose to create a series of simple scripts to create the Conversational Hebrew Vocabulary List.

The two most widely-used languages for the type of data analysis involved in a word list creation are Python and R. I chose to use Python for this project. Python was designed specifically to be a very readable programming language. That is, it is easy to read and understand the purpose and flow of the code. This was one of my primary reasons for choosing to use it, since it increases the ease with which this project can be reproduced by other researchers and educators to create their own word lists. R, on the other hand, requires a deeper familiarity with the syntax and conventions of the language in order to understand.

The second characteristic that makes Python ideal for an open-source project of this nature is its mild learning curve. Though considerable effort must be made to learn any programming language, Python is widely considered good for beginners because of its simplicity. With only a rudimentary knowledge of Python, even educators or enthusiasts without a coding background will be able to modify the scripts used here to suit their own needs. To this end, I will also carefully explain what, exactly, the code does.

Though all of the code is included in this thesis (*Appendix 2*), it can also be found

in an online repository at https://github.com/juandpinto/opus-lemmas. The repository can easily be cloned, or individual files can be downloaded, for modification and use. The repository uses the version control system *Git*. This means that anyone can easily look through the history of each file to see specific changes that have been made over time.

Suggestions for improvements can also be submitted through the GitHub interface, allowing for a system of cooperation and incremental innovation among researchers. The exported Conversational Hebrew Vocabulary List, in its entirety, can also be found in the repository.

This thesis, then, beyond explaining the theory behind the creation of the CHVL, aims to make the process as reproducible as possible. This section contributes to that aim by carefully documenting each step of the process.

## 3.1 THE CORPUS

Before coding or analyzing anything, it's important to find an appropriate corpus to use and to become familiar with its structure. A useful place to begin is OPUS[3], which is part of the Nordic Language Processing Laboratory (NLPL), and hosted by the CSC IT center in Finland. OPUS is a database of many open, parallel corpora. These include corpora of movie and television subtitles, TED talks, web-crawled data, newspapers, and of course, books. The corpora are all free and open to the public.

The CHVL was created using one of OPUS's corpora, the OpenSubtitles2018[4] corpus. The corpus can be downloaded in a variety of formats, and can be downloaded either as *parallel* corpora, or as a monolingual corpus. A parallel corpus consists of two languages interwoven together. For example, a line from the English subtitles of a movie will be paired with the same line from the French subtitles of the same movie. In theory, this means that each line of the corpus should have the same meaning in two different languages. The creation of parallel corpora has made possible many interesting and useful tools for linguistics, translators, and language learners. These

---

[3]http://opus.nlpl.eu
[4]http://opus.nlpl.eu/OpenSubtitles2018.php

include the open-source CASMACAT[5] project and the ReversoContext[6] tool.

For the purpose of creating a word list, a monolingual corpus is best. Note that parallel corpora will often be composed of less tokens than monolingual ones. This is because parallel corpora will only include movies for which the subtitles exist in both selected languages.

Though it's possible to download plain text files, the most useful format available for download is XML. Indeed, the most common file format used for large corpora is XML. The XML structure allows for nested key-value pairs, which are especially useful for parsed corpora that contain extensive metadata. XML is comparable to JSON, which we will use later to extract specific movie metadata directly from a database.

Another factor to consider is whether to download an untokenized, tokenized, or parsed corpus. An untokenized corpus contains simply the raw lines of text as found in the original subtitle files (divided into lines as they would appear while watching the movie, and labeled with the appropriate time for them to be shown):

```
<s id="49">
  <time id="T39S" value="00:03:22,280" />
שרלוק ,אומר אתה מה?
  <time id="T39E" value="00:03:24,120" />
</s>
```

A tokenized corpus has further been split into individual words and punctuation, such that each word is tagged on its own:

```
<s id="49">
  <time id="T39S" value="00:03:22,280" />
  <w id="49.1">מה</w>
  <w id="49.2">אתה</w>
  <w id="49.3">אומר</w>
```

---

[5]http://www.casmacat.eu
[6]http://context.reverso.net/translation/

```
  <w id="49.4">,</w>
  <w id="49.5">שרלוק</w>
  <w id="49.6">?</w>
  <time id="T39E" value="00:03:24,120" />
</s>
```

A parsed corpus contains much more information for each token. The data included depends on the features of the language and on the parsing script used, but it can include things such as part of speech, syntactic role, lemma, and even specific features like gender, person, and number. Here is an example:

```
<s id="49">
  <time value="00:03:22,280" id="T39S" />
  <w xpos="ADV" head="49.3" feats="PronType=Int" upos="ADV"
  ↪  lemma="מה"
     id="49.1" deprel="obj">מה</w>
  <w xpos="PRON" head="49.3" feats="Gender=Masc|Number=Sing|Person=2|
     PronType=Prs" upos="PRON" lemma="הוא" id="49.2"
↪ deprel="nsubj">אתה</w>
  <w xpos="VERB" head="0"
  ↪  feats="Gender=Masc|HebBinyan=PAAL|Number=Sing|
     Person=1,2,3|VerbForm=Part|Voice=Act" upos="VERB"
↪ misc="SpaceAfter=No"
     lemma="אמר" id="49.3" deprel="root">אומר</w>
  <w xpos="PUNCT" head="49.3" upos="PUNCT" lemma="," id="49.4"
     deprel="punct">,</w>
  <w xpos="NOUN" head="49.3" feats="Gender=Masc|Number=Sing"
  ↪  upos="NOUN"
     misc="SpaceAfter=No" lemma="שרלוק" id="49.5"
↪ deprel="obj">שרלוק</w>
  <w xpos="PUNCT" head="49.3" upos="PUNCT" misc="SpaceAfter=No"
  ↪  lemma="?"
     id="49.6" deprel="punct">?</w>
```

```
  <time value="00:03:24,120" id="T39E" />
</s>
```

All of the data used to create the CHVL came from a monolingual parsed corpus of Hebrew. The parsing was all done automatically using .

## 3.2  CLEANSING THE CORPUS

Unlike many corpora, the OpenSubtitles2018 corpus as presented in its downloadable form has already undergone significant preprocessing by the OPUS team.(Lison & Tiedemann, 2016) This is good news, since data cleansing is often the most laborious part of the process. However, there is one issue that must be addressed before the corpus can be used to create a word list: deduplication.

The files inside the downloaded folder are organized as follows:

```
Zipped folder in GZ format
    Folder for year X
        Folder for movie A
            Zipped XML in GZ format
            Zipped XML in GZ format
            Zipped XML in GZ format
        Folder for movie B
            Zipped XML in GZ format
            Zipped XML in GZ format
    Folder for year Y
        Folder for movie C
            Zipped XML in GZ format
        Folder for movie D
            Zipped XML in GZ format
            Zipped XML in GZ format
            Zipped XML in GZ format
        Folder for movie E
```

```
        Zipped XML in GZ format
        Zipped XML in GZ format
    Folder for year Z
        Folder for movie F
            Zipped XML in GZ format
            Zipped XML in GZ format
```

This organization is straight-forward, except for the fact that there are multiple XML files for each movie. The subtitle files that OPUS has collected, parsed, organized, and made available for mass download were all obtained from the Open Subtitles[7] project (hence the name of the corpus). Because this is a database where users can upload the subtitle files they extract from their own movie collection, there are often multiple uploads for the same movie. For our purposes, this results in movies that can have anywhere from a single subtitle file to dozens of them. Unfortunately, though the tokens in the files themselves are usually the same (with only minor variations in the XML metadata), this is not always true. Some few variations seem to be different and independent translations.

Part of cleansing the corpus, then, entails getting rid of these duplicates. As a means of simplifying the entire process, I chose simply to use the first file in each movie folder. I've included the short Python script for this in its entirety in *Appendix 2.3*. However, I will here explain what it does in detail so that it can be easily modified to fit different circumstances.

The script first makes a copy of the entire folder structure in the original downloaded (and unzipped!) corpus into a new directory. It then finds the first XML file in each movie folder and copies it into the appropriate place in the new folder structure. This means that it doesn't delete or otherwise change the files in the original corpus in any way.

The first block of code imports necessary modules that are used later in the script (`shutil` and `os`). Lines 7 and 8 define where the original corpus is (`source`), and where the new one will be placed (`destination`).

---

```
4  import shutil
5  import os
6
7  source = '../OpenSubtitles2018_parsed'
8  destination = './OpenSubtitles2018_parsed_single'
```

Next, a single line of code copies all directories and subdirectories into their new location.

```
11  shutil.copytree(source, destination,
↳     ignore=shutil.ignore_patterns('*.*'))
```

Lastly, we create a variable that holds all the XML files located in each movie folder, trim the list to just one, and copy that one into its new location. This process is carried out for one movie folder at a time. The originals are left untouched.

```
14  for dirName, subdirList, fileList in os.walk(source):
15      for fname in fileList:
16          if fname == '.DS_Store':
17              fileList.remove(fname)
18      if len(fileList) > 0:
19          del fileList[1:]
20          src = dirName + '/' + fileList[0]
21          dst = destination + dirName[27:] + '/'
22          shutil.copy2(src, dst)
```

With a newly organized version of the corpus, it's now possible to begin the process of reading and processing data. At this stage, I took some time to gather metadata for all the movies in the corpus in order to identify movies that were originally filmed with Hebrew as their primary language (as opposed to translated subtitles). Because I ultimately decided against this approach for the creation of the CHVL, I will skip that step here. However, a description of the entire process will be discussed later under *Using original-language movies exclusively*.

## 3.3 Reading data

Before calculating any measures such as frequency, individual lemmas must be extracted from the XML files in the downloaded corpus. There are two ways to go about this. Because XML consists of nested tags and key-value pairs, a dedicated XML parsing tool can be used to extract specific information. In this case we would be creating a list of all *values* in the `'lemma'` *key* within each `<w>` *tag.* The value that corresponds to the `'lemma'` tag below for the word אומר is אמר.

```
<w xpos="VERB" head="0"
↪  feats="Gender=Masc|HebBinyan=PAAL|Number=Sing|
   Person=1,2,3|VerbForm=Part|Voice=Act" upos="VERB"
↪  misc="SpaceAfter=No"
   lemma="אמר" id="49.3" deprel="root">אומר</w>
```

A different approach is to use *regular expressions* to search for a specific string of characters and extract every instance of that string. This is a more brute-force approach, since it ignores the structure of the XML file and treats it all simply as raw text. To find a lemma, a very simple regular expression is sufficient: `lemma="[א-ת]+"`. This will search for any instance of the characters `lemma="`, followed by a combination of any number of Hebrew letters (at least one), followed by the character `"`.

Despite the existence of various Python modules for parsing XML files, I found a simple search using regular expressions to be more efficient for various reasons. First, not all elements in the parsed corpus contain *lemma* attributes. Second, punctuation and non-Hebrew words are often lemmaticized. This means that even after extracting all the *lemma* values in a file, I would still need to use regular expressions to search through the results and delete any that contain non-Hebrew characters. I chose instead to skip the XML parsing step altogether.

I will now explain the code in the script used to create the CHVL. As with the other code, the entire script in its entirety can be found in *Appendix 2.1*.

After importing necessary packages and initializing variables, two functions near the beginning of the script serve to open a file and extract a list of lemmas from it.

```python
# Open XML file and read it.
def open_and_read(file_loc):
    with gzip.open(file_loc, 'rt', encoding='utf-8') as f:
        read_data = f.read()
    return read_data
```

```python
# Search for lemmas and add counts to "lemma_by_file_dict{}".
def find_and_count(doc):
    file = str(f)[40:-3]
    match_pattern = re.findall(r'lemma="[א-ת]+"', doc)
    for word in match_pattern:
        if word[7:-1] in lemma_by_file_dict:
            count = lemma_by_file_dict[word[7:-1]].get(file, 0)
            lemma_by_file_dict[word[7:-1]][file] = count + 1
        else:
            lemma_by_file_dict[word[7:-1]] = {}
            lemma_by_file_dict[word[7:-1]][file] = 1
```

We then run both of these functions for each XML file in the corpus directory (defined earlier in `corpus_path`).

```python
for dirName, subdirList, fileList in os.walk(corpus_path):
    if len(fileList) > 0:
        f = dirName + '/' + fileList[0]
        find_and_count(open_and_read(f))
```

The `find_and_count()` function finds each instance of the string described above using a regular expression, then adds the Hebrew part of the string—the lemma itself—to a dictionary. The dictionary is named `lemma_by_file_dict`, and its structure looks like this:

```
'lemma': {'path of file': 'frequency of lemma in file'}
```

A dictionary is at its core a list of key:value pairs. Much like an actual dictionary consists of words and their definitions, this dictionary's keys are made up of all the individual lemmas found by our search. For each lemma, the value is another dictionary—making it a nested dictionary, or a dictionary within a dictionary. The keys for each inner dictionary are the paths of all the XML files (movies) that the lemma appears in, and the value of each is an integer that represents how many times that lemma appears in that file (frequency).

After the script reads each file, it returns a complete dictionary. Here is a sample:

```
'ב': {
        '/he/0/5753574/6853341.xml': 168,
        '/he/0/3607000/5764778.xml': 94},
'פרק': {
        '/he/0/5753574/6853341.xml': 3},
'קודם': {
        '/he/0/5753574/6853341.xml': 6,
        '/he/0/3607000/5764778.xml': 2,
        '/he/0/1278351/3777598.xml': 1}
```

Throughout the rest of the script, this nested dictionary serves as the basis for all of the calculations needed.

## 3.4  Calculations

For each lemma, the CHVL includes three measures: frequency, range, and $U_{DP}$ (dispersion). It uses dispersion as its sorting value. Let's look at how each of these is calculated. Range will be addressed in the export section, since the script calculates it on the spot as the list is created.

### 3.4.1 Frequency

Since we've already calculated the frequency of each lemma for each individual file, calculating total frequency per lemma is straight forward. The script simply creates a new dictionary, `lemma_totals_dict`, and adds to it every lemma in the corpus as its keys, with the corresponding value being a sum of the frequencies in all files for that lemma. In other words, {'lemma1':'frequency1', 'lemma2':'frequency2', . . . }

```
116  for lemma in lemma_by_file_dict:
117      lemma_totals_dict[lemma] =
    ↪    sum(lemma_by_file_dict[lemma].values())
```

This returns Using the short example given above, this would result in the following dictionary:

'ב':262,
'פרק':3,
'קודם':9

### 3.4.2 U$_{DP}$ (dispersion)

Dispersion is more complicated. In theory, it should provide a single quantifiable measure that incorporates both frequency and range, and which can then be used to sort the word list. There is no agreed-upon, single way to calculate dispersion, and different researchers will use the words in slightly different contexts. The model of dispersion I have chosen to follow for this project is Gries' dispersion coefficient, or U$_{DP}$, () calculated from Gries' DP. ()

In order to calculate Gries' DP for lemma$_x$, we must first make two calculations for each file in the corpus (file$_i$): the lemma's *expected frequency* if it were perfectly distributed, and its *observed frequency*—or its actual frequency.

$$\textbf{expected frequency} \; = \; \frac{tokens \; in \; file_i}{tokens \; in \; corpus}$$

32

$$\textbf{observed frequency} \ = \ \frac{frequency \ of \ lemma_x \ in \ file_i}{frequency \ of \ lemma_x \ in \ corpus}$$

We must then subtract the lemma's observed frequency from its expected frequency, which will return a value between -1 and 1. We can normalize this result by finding the absolute value. Now the closer the result is to 0, the closer that lemma's frequency is in that particular file to what we would expect if it were perfectly distributed throughout the corpus. A higher number (closer to 1), would indicate a heavier load in that file that we would expect.

By performing this calculation for every file in the corpus, adding them all together, and dividing the result by two (since we're using the absolute value and are therefore adding values originally in both directions), we now have Gries' DP. Where $\texttt{n}$ is the number of files:

$$\textbf{DP} \ = \ 0.5 \sum_{i=1}^{n} | \ \text{expected frequency} \ - \ \text{observed frequency} \ |$$

A DP of 0 represents a perfectly even dispersion, and a DP close to 1 means a more uneven distribution, where fewer files contain a larger load of the lemma's overall frequency. A DP of 1 is not actually possible.

Gries' usage coefficient, or $U_{DP}$, is an attempt to make this number more useful. DP is first subtracted from 1 and the result is multiplied by the lemma's total frequency. The full equation for $U_{DP}$ is as follows:

$$\left(1 - 0.5 \sum_{i=1}^{n} \left| \frac{file_i \ tokens}{total \ tokens} \ - \ \frac{frequency_x \ in \ file_i}{total \ frequency_x} \right| \right) \times total \ frequency_x$$

In order to calculate this, the script must first find the number of tokens in each file. Like before, this is done by creating a dictionary, $\texttt{token\_count\_dict}$, which contains the key:value pairs of file:tokens. Since we already have a dictionary with the number of times that each lemma appears in each file, $\texttt{lemma\_by\_file\_dict}$, we don't need to open and read the files again. Instead, we can add the values in this

dictionary and rearrange them into what we want.

```python
for lemma in lemma_by_file_dict:
    for file in lemma_by_file_dict[lemma]:
        token_count_dict[file] = token_count_dict.get(
            file, 0) + lemma_by_file_dict[lemma][file]
```

We also need to know the total number of tokens in the entire corpus. This is a simple matter of adding all the values in the `token_count_dict` dictionary. The final count is saved into an integer variable, `total_tokens_int`.

```python
for file in token_count_dict:
    total_tokens_int = total_tokens_int + token_count_dict.get(file,
    ↪  0)
```

Finally, the script uses all these measures to calculate DP and then $U_{DP}$ for each lemma, and places them into their respective dictionaries, `lemma_DPs_dict` and `lemma_UDPs_dict`.

```python
# Calculate DPs
for lemma in lemma_by_file_dict.keys():
    for file in lemma_by_file_dict[lemma].keys():
        lemma_DPs_dict[lemma] = lemma_DPs_dict[lemma] + abs(
            (token_count_dict[file] /
             total_tokens_int) -
            (lemma_by_file_dict[lemma][file] /
             lemma_totals_dict[lemma]))
lemma_DPs_dict = {lemma: DP/2 for (lemma, DP) in
↪  lemma_DPs_dict.items()}

# Calculate UDPs
lemma_UDPs_dict = {lemma: 1-DP for (lemma, DP) in
↪  lemma_DPs_dict.items()}
```

With these values all calculated for each lemma, the only thing left is to sort and create the final list.

## 3.5  SORT AND EXPORT

In order to ensure that the words on the list do not have an abnormally high frequency in some subcorpora (movies) and are nearly absent in others, some have suggested setting a minimum range or dispersion. All words that fall below this threshold are discarded, and the remaining words can then be sorted by frequency.

Though this is a more systematic approach than that used to create many early frequency lists, it still depends on a subjective decision and the whim of the researcher.

Rather than setting an arbitrary bar, the CHVL is sorted entirely by Gries' usage coefficient of dispersion ($U_{DP}$). This *modus operandi* ensures that the order of words itself—not just which words make it onto the list and which don't—is decided by a combination of both relevant measures: frequency and dispersion. This approach also has the added benefit of being entirely objective.

Since we've already calculated the $U_{DP}$ for each lemma, sorting the list is simple.

```
148  UDP_sorted_list = [(k, lemma_UDPs_dict[k]) for k in sorted(
149      lemma_UDPs_dict, key=lemma_UDPs_dict.__getitem__,
150      reverse=True)]
```

A final table is then created (using a list of tuples, `table_list`), with each line consisting of a lemma, its overall frequency, its range, and its $U_{DP}$. This table is already sorted by $U_{DP}$ as it's being created.

Because the script has not calculated range by this point, it must do so on the spot as it's entering each lemma into the table. It does this with a simple dictionary comprehension that quickly counts the number of files included in the `lemma_by_file_dict`. Here is the resulting code:

35

```
153  for k, v in UDP_sorted_list[:list_size_int]:
154      table_list.append((k, lemma_totals_dict[k], sum(
155          1 for count in lemma_by_file_dict[k].values() if count > 0),
156          v))
```

Lastly, now that everything is organized into a table, the script opens (or creates, if it doesn't yet exist) a CSV file, writes a header line into it (`LEMMA, FREQUENCY, RANGE, UDP`), and exports the entire table into the file. It then closes it to clear the computer's memory cache.

```
199  result = open('./export/WordList.csv', 'w')
200  result.write('LEMMA, FREQUENCY, RANGE, UDP\n')
201  for i in range(list_size_int):
202      result.write(str(table_list[i][0]) + ', ' +
203                   str(table_list[i][1]) + ', ' +
204                   str(table_list[i][2]) + ', ' +
205                   str(table_list[i][3]) + '\n')
206  result.close()
```

The list is now complete. The next section will explore the list itself more in-depth.

# 4 The CHVL: A vocabulary list of conversational Modern Hebrew

The Conversational Hebrew Vocabulary List in its entirety can be found as an electronic supplement to this thesis (in CSV format) or at the following GitHub repository: *https://github.com/juandpinto/opus-lemmas*. It contains the most common 5,000 lemmas of conversation Modern Hebrew, as found in the OpenSubtitles2018 corpus. A sample of the first 1,000 lemmas is included in *Appendix 1*.

For discussion purposes, a small sample of the first 20 items is here presented.

Table 1: Sample of the first 20 items on the CHVL.

| RANK | LEMMA | FREQUENCY | RANGE | $U_{DP}$ |
|------|-------|-----------|-------|------|
| 1 | הוא | 23446109 | 43455 | 0.9480170255915042 |
| 2 | ל | 5638813 | 43448 | 0.9420130372643667 |
| 3 | ה | 9850733 | 43458 | 0.929266134661147 |
| 4 | ב | 4812778 | 43450 | 0.9292364864789281 |
| 5 | את | 6846782 | 43426 | 0.9285176069174289 |
| 6 | לא | 5272808 | 43433 | 0.9145688112131216 |
| 7 | ש | 3880654 | 43439 | 0.9088900047303463 |
| 8 | של | 3892328 | 43445 | 0.9067041511201389 |
| 9 | על | 1766990 | 43430 | 0.9042865019832009 |
| 10 | זה | 5118759 | 43441 | 0.9015544612816044 |
| 11 | מה | 2362419 | 43403 | 0.8922532708182579 |
| 12 | היה | 2579370 | 43420 | 0.8909904417204713 |
| 13 | מ | 1061614 | 43411 | 0.88900672760779 |
| 14 | כול | 1325676 | 43414 | 0.8860074112131449 |
| 15 | ו | 1906717 | 43429 | 0.8852706380348441 |
| 16 | יש | 1069358 | 43376 | 0.8770543442171884 |
| 17 | עם | 839575 | 43331 | 0.8668140051895192 |
| 18 | אם | 861163 | 43321 | 0.8654587702150129 |
| 19 | ידע | 1202416 | 43323 | 0.8586088803742931 |
| 20 | אבל | 921757 | 42963 | 0.8519038846130076 |

Besides each lemma and its respective rank on the list, the CHVL includes three pieces of information: frequency, range, and $U_{DP}$. Frequency in this case is not raw frequency—the total number of times the lemma appears in the corpus—but rather how many times the lemma appears for every million tokens in the corpus. Using frequency per million makes the number more meaningful since—in theory—it reflects the per-million count of all spoken Hebrew, not just the OpenSubtitles2018 corpus. The range is the number of sub-corpora—or in this case, movies—the lemma appears in.

The most important piece of information the list provides, however, is the $U_{DP}$, which refers to Griers' usage coefficient for dispersion. This is discussed more in-depth in the methods section above. $U_{DP}$ is also used as the sorting measure for the CHVL.

The percentage of the corpus that is covered by the first $n$ items on the list is referred to as coverage. This is a simple matter of finding the total number of tokens in the corpus, and dividing from it the sum of all the *raw* frequencies from the first $n$ items.

For example, the sum of the frequencies of the first 20 lemmas in *Table 1* (84,656,819) divided by the total size of the corpus (193,755,220) is 0.436926649. In theory, this means that by knowing just the first 20 lemmas on the CHVL one would be able to understand 43.7% of the words in the entire OpenSubtitles2018 corpus! That is a clear example of the power of Zipf's Law (see *Introduction* for more on Zipf's Law).

Table 2 presents a listing of some important coverages provided by different amounts of lemmas on the CHVL.

Table 2: Breakdown of coverage percentages.

| $n$ Lemmas | Frequency Sum | $\div$ Corpus Size | $=$ Coverage |
|---|---|---|---|
| 374 | 135,767,644 | 193,755,220 | 70% |
| 939 | 155,016,588 | 193,755,220 | 80% |
| 4,246 | 174,380,519 | 193,755,220 | 90% |
| 13,758 | 184,067,666 | 193,755,220 | 95% |

The entire CHVL consists of 5,000 lemmas. This number was chosen in order for it to include the required items for 90% coverage, while also making it an even factor of

1,000. In its entirety, the CHVL covers 90.8% of the corpus from which it is created.

## 4.1 Challenges and future direction

Throughout the course of this project, I have encountered several issues that are worth discussing. Some of these are questions that require further study in order to address adequately. Others are technical issues related to the complex task of pre-processing and parsing the corpus—something not directly dealt with in this thesis. Others yet are simple suggestions that I simply did not have time to implement given this project's time constraints. And finally, there are limitations that are the inevitable result of the tools at hand.

I have divided all of these issues into two categories: methodological challenges of a bigger nature, and functional challenges of a more limited scope.

### 4.1.1 Methodological challenges

One of the more obvious issues of this project is the use of a corpus of movie subtitles as substitute for a corpus of true conversational language. This issue in a way forms the backbone of the CHVL, and it is at the heart of what this project is all about. Though I discuss several points related to this in the *Background* section of this thesis, I will here discuss some of its implications for future work.

**4.1.1.1 Ideal vs. practical corpora** The use of a subtitle corpus has both positive and negative aspects. As mentioned earlier, the early research that has been done on the topic indicates that movie subtitles share many features with spontaneous, spoken language. This includes a high level of correlation between the two , as well as a strong ability to predict the outcomes of a lexical decision task .

One especially positive aspect of subtitle corpora is their accessibility. Thanks to the efforts of organizations such as *http://opensubtitles.com* and OPUS[8], very large

---

[8]http://opus.nlpl.eu

corpora are available to the public for free. And they already come pre-processed, as an additional incentive for the time-constrained researcher.

This free and open nature makes subtitle corpora excellent tools for research in languages that don't yet have large, high-quality corpora of spoken language. Though advances in technology are rapidly making this type of data-collection more accessible, the costs remain too high for many less-commonly taught languages as of now. This is largely due to the arduous process of transcribing audio recordings.(Izre'el, 2004)

An ideal corpus for this sort of task would consist of many millions of tokens of recorded, transcribed, and parsed, spontaneous spoken language. Several attempts have been made to create a corpus of this nature in Hebrew.

The most prominent of these is the Corpus of Spoken Israeli Hebrew (CoSIH)[9], created at Tel Aviv University between 2000 and 2002.(Izre'el, Hary, & Rahav, 2001) Designed and initiated by a team of distinguished scholars, it unfortunately ran out of funding long before its goals were met. The CoSIH website (*http://cosih.com/*) makes available to the public a total of 13.5 hours of recorded Hebrew, with just over five hours of it having been transcribed.

Though a few publications have used data from CoSIH, these have been primarily methodological studies for the design of the project itself.(Amir, Silber-Varod, & Izre'el, 2004; Izre'el et al., 2005; Mettouchi, Lacheret-Dujour, Silber-Varod, & Izre'el, 2007) At least one dissertation, by Nurit Dekel, uses data exclusively from CoSIH. Her entire corpus consists of 44,000 tokens. (2010, p. 7)

Other corpora of spoken Hebrew include the Haifa Corpus of Spoken Hebrew (Yael, 2014) and the Hebrew CHILDES corpus (Albert, MacWhinney, Nir, & Wintner, 2013; Gretz, Itai, MacWhinney, Nir, & Wintner, 2015). The first consists of 17.5 hours of audio recordings, along with a limited selection of transcribed text. The latter is a collection of recordings of interactions between adults and children, comprising a total of 417,938 transcribed tokens. The CHILDES corpus is unique in that the transcriptions are provided using a Latin-based phonemic transliteration. This was done in order to avoid many of the textual ambiguities of using the Hebrew script,

---

[9]http://cosih.com/

which are addressed below under *Functional challenges.*

Though ideal in some ways, these corpora remain far too small to be effectively used for the creation of frequency lists. Even combined into a single corpus (which would introduce a series of new issues to solve), the total size would not be bigger than two million tokens. As discussed earlier in this thesis, Sorell provides evidence to suggest that a corpus of 20–50 million tokens is the minimum for a stable word list.(2013)

Are movie and television subtitles an suitable substitute for spontaneous, spoken language? Early studies suggest it is at least adequate, but much more research is needed to answer this question definitively. For now, it remains as one practical option.

#### 4.1.1.2 Using original-language movies exclusively

One of the potential downsides of using the OpenSubtitles2018 corpus not yet discussed is that it includes all subtitles of a specific language, even *translated* subtitles from movies filmed in other languages. The question is, does a translated script represent true conversational language as faithfully as an original script?

This is a question that requires more research in order to answer satisfactorily. Though translated subtitles don't need to try to approximate the utterance length and visual cues that a dubbed script does, its quality still largely depends on the skills of a translator. Most importantly, a translation may not accurately reflect the register of the original, no longer serving as a representation of conversational language. Again, these are important points to consider.

One solution is to simply use movies that were originally filmed in the target language of the corpus. In theory, each XML file in a monolingual OpenSubtitles2018 file should contain a tag that identifies the original language of the movie. In practice, I found that the overwhelming majority of the files contained an empty `<lang>` tag instead. Luckily, there is a way to obtain the desired metadata for each movie in the corpus.

This can be done with a script that uses an application programming interface (API) to fetch specific information from an online movie database. The name of each movie folder in the corpus, which is simply a series of numbers, corresponds to that movies

IMDb ID, which is a unique ID registered with the Internet Movie Database[10]. This makes the process relatively easy, as we simply need to query the database using this ID to receive all of the movie's metadata.

Though IMDb does provide their own API, I decided instead to use an API created for the Open Movie Database (OMDb)[11]. This API can be used free-of-charge, but it has a 1,000 movie limit per day. Since the OpenSubtitles2018 Hebrew corpus contains nearly 50,000 movies, I decided instead to pay for a daily limit of 100,000 movies. This only requires a $1.00 donation for each month that one is registered to use the OMDb API.

Once an API key is obtained, a script can be written to obtain the information desired for every movie all at once. In this case, we want to know the original language(s) for each movie.

This script in its entirety is found in Appendix 2.2[12]. It uses an imported Python wrapper for the API, written by Derrick Gilland[13], which can be found at https://github.com/dgilland/omdb.py. This package can be installed through PIP by entering `pip install omdb` into the command line.

For practical purposes, the script requires one to enter a specific year (or, more accurately, corpus folder name). If desired, an asterisk can act as wildcard: `python OMDb-fetch.py 1988` will fetch data for movies from 1988, while `python OMDb-fetch.py 198*` will do it for all movies in the 1980s. In order to fetch data for all movies in the database at once, use `python OMDb-fetch.py *`. I don't recommend this, however, since it may overload the server and cause the script to time out.

The script begins by creating a list of all movie directory paths for the desired year.

```
15  for name in glob.glob(
16          './OpenSubtitles2018_parsed_single/parsed/he/' + year +
        ↪   '/*/'):
```

---

[10]http://www.imdb.com/
[11]http://www.omdbapi.com/
[12]14_appendix_2.md
[13]https://github.com/dgilland

```
17          IDs.append(name)
```

Each item in the list is then trimmed to include only the name of the movie folder, which is *almost* equivalent to the IMDb ID.

```
20    IDs = [os.path.basename(os.path.dirname(str(i))) for i in IDs]
```

In order to make the IDs match those in the database, additional zeros must be added to the beginning until they are seven digits long.

```
23    for i in IDs:
24        while len(i) < 7:
25            IDs[IDs.index(i)] = '0' + i
26            i = '0' + i
```

The list is then sorted numerically in order to more easily interpret the results: `IDs.sort()`.

The API key is set in line 32, but be sure to replace `906517b3` with your own key, which can be obtained at http://www.omdbapi.com/.

```
32    omdb.set_default('apikey', '906517b3')
```

The script then prints a table header, fetches the title, year, and language(s) for each movie, and prints the results directly into the computer terminal.

```
35    print('# ' + year + '\n' +
36          'IMDb ID\tTitle\tYear\tLanguage(s)')
```

```
39      for i in IDs:
40          doc = omdb.imdbid('tt' + i)
41          print('tt' + i + '\t' +
42                  doc['title'] + '\t' +
43                  doc['year'] + '\t' +
44                  doc['language'])
```

### 4.1.2 Functional challenges

A quick scan of the CHVL reveals some notable items. Some of these are mere quirks of the automatic parser, while others are the result of ambiguities.

For example, the very first lemma on the list is a bit unexpected. "הוא" is certainly not the most common lemma in Modern Hebrew. A quick look at some of the files in the corpus, however, reveals that all pronouns are grouped under this lemma. That is, אתה (you), היא (she), and אנחנו (we), just to name a few, are parsed as belonging to the lemma "הוא." Considering how common pronouns are in the majority of spoken dialogue (in many languages), its place at the top of the list ceases to be a surprise.

Another thing to note is that verbs are all listed in their traditional third-masculine-singular past conjugation. The first verb on the list is "היה"—a lemma referring to all forms of the verb להיות, including the infinitive. The same is true of "ידע" (item 19) and "דיבר" (item 60).

Many of the most common lemmas on the CHVL are prepositions. Note that even inseparable prepositions, such as -ה and -ב are considered independent lemmas by the parser, and are listed respectively as the lemmas "ה" and "ב".

Other issues, however, are more difficult to explain.

**4.1.2.1 Textual ambiguity of Hebrew orthography** The flexible spelling conventions of Hebrew are at the root of many of the problems with the CHVL. For example, דִּבֵּר *he spoke* can be written as either דיבר ("full spelling") or דבר ("defective spelling"). There is also a noun, דָּבָר *thing*, that looks identical to the verb's defective

spelling (דבר). Though the difference is usually clear from context, the automatic parser has some difficulty with this orthographic ambiguity.

The lemma "דבר" (item 27) includes instances of both the verb and the noun, which are completely unrelated. A simple search through the corpus reveals multiple examples of the noun דבר tagged with `lemma="דבר"`:

```
<w xpos="NOUN" head="579.3" feats="Gender=Masc|Number=Sing"
↪  upos="NOUN" lemma="דבר" id="579.2" deprel="nsubj">דבר</w>

<w xpos="NOUN" head="200.11" feats="Gender=Masc|Number=Plur"
↪  upos="NOUN" lemma="דבר" id="200.12" deprel="obj">דברים</w>
```

We also find plenty of examples of the verb with the same lemma tag:

```
<w xpos="VERB" head="0"
↪  feats="Gender=Fem|HebSource=ConvUncertainHead|Number=Sing|Person=3|Tense=Past"
↪  upos="VERB" lemma="דבר" id="2346.4" deprel="root">דברה</w>

<w xpos="VERB" head="0"
↪  feats="Gender=Fem,Masc|Number=Plur|Person=1|Tense=Past"
↪  upos="VERB" lemma="דבר" id="1270.2" deprel="root">דברנו</w>

<w xpos="VERB" head="0"
↪  feats="Gender=Fem,Masc|Number=Plur|Person=3|Tense=Past"
↪  upos="VERB" lemma="דבר" id="368.4" deprel="root">דברו</w>
```

A different lemma, "דיבר" (item 61), is the expected lemma for the verb since it follows the standard third masculine plural conjugation. Interestingly, however, the parser applies this lemma only to attestations of the word with an inserted *yod*, or with a *mem* or *lamed* prefix (present tense or infinitive). All other instances are parsed as the lemma "דבר." Though unexpected and simply wrong, at least this issue is consistent.

```
<w xpos="VERB" head="840.4"
↪  feats="Gender=Fem,Masc|HebBinyan=HITPAEL|Number=Plur|Person=1|Tense=Past"
↪  upos="VERB" lemma="דיבר" id="840.16" deprel="conj">דיברנו</w>

<w xpos="VERB" head="1451.12"
↪  feats="Gender=Masc|HebBinyan=PIEL|Number=Sing|Person=1,2,3|VerbForm=Part|Voice=A
↪  upos="VERB" lemma="דיבר" id="1451.20" deprel="obl">מדבר</w>
```

To complicate matters more, we also find the unexpected lemmas "דיברה" (item 1184),
"שדיבר" (item 2588), and "שדיברה" (item 4106).

Which, based on context (), should clearly be parsed as two separate lemmas, "ש"
and "דיבר."

These are just a few among many examples of the difficulties encountered by the
automatic parser. Though the parsing was carried out by the OPUS team as part
of the corpus's pre-processing stage, it is valuable to at least have an idea of how
it works its magic. I will here explain the basics of the process and some of the
implications entailed.

**4.1.2.2 Automatic parsing**   Automatic parsing refers to the process of having a
computer program create a syntactic tree for a corpus of natural language. Natural
language, as opposed to artificial or constructed language, is notoriously complex
in its structure. Natural language processing (NLP) is an entire field of research,
currently at the forefront of computer science. Parsing can serve many purposes, from
theoretical linguistic research to machine translation or even the creation of artificial
intelligences such as Siri or Alexa. For our purposes, a parsed text is important in
order to use lemmas as the word family level for the CHVL. This decision is discussed
under *Identifying Words* in this thesis.

Two distinct types of syntactic parsers exist, contituency parsers and dependency
parsers. These are based on the two respective linguistic theories of syntax, con-
stituent grammar (sometimes referred to as phrase structure grammar) and depen-
dency grammar.

Constituent grammar is the classic syntax tree structure taught in introductory-level linguistics classes. It is essentially a theory of the logic structure of language as a whole. Dependency grammar is a competing theory that treats words as more directly interconnected to each other. A thorough description of these ideas is outside the scope of this thesis, and is not pertinent to the project. What is important to know is that dependency grammar, and thus dependency parsers, have played an important role in the advancement of NLP and computational linguistics as a whole. The term "automatic parser", therefore, most often refers to an automatic *dependency* parser.

Some parsers proceed in a two-step process of morphological tagging (part of speech) and then dependency parsing (syntactic role and conjugations). In all cases, tokenization must first take place, which refers to splitting the text into individual lemmas.

Most automatic parsers are "trained" using a small corpus that has been manually parsed by a human previously, or at least one that was automatically parsed and then checked and corrected by the researcher. These "gold-standard" pre-parsed corpora are called treebanks, and repositories of them they have been created for many languages. Building on existing databases of knowledge, these many of these parsers use statistical models to determine the most likely syntactic structure and conjugation for each word in each sentence.

Some parsers, however, are instead simply given entirely unparsed corpora and no knowledge of the language's syntactic structure. Working with nothing but the text itself, the program seeks out patterns and begins to create links and relationships that it deems significant.

Unfortunately, though automatic parsers have achieved surprising levels of accuracy in recent years, even the best continue to produce erroneous parsings. Some researchers have claimed as 95% or higher accuracy, including for some Hebrew parsers. When dealing with such a large corpus, such as the Hebrew OpenSubtitles2018 corpus consisting of nearly 200 million tokens, a best-case scenario for a 5% error threshold results in nearly 10 million incorrectly parsed words.

Undoubtedly, this can have a negative impact on the accuracy of lemma frequency counts. Many of the issues found in the CHVL are not due to orthographic ambiguity,

but simply to inaccurate parsing. Some, as shown in the previous section, are even caused by erroneous automatic tokenization (consider the lemma "שדיבר").

The good news is that automatic parsers are continually improving in accuracy. This is a problem that exists across the board, regardless of the corpus being used—unless it is manually parsed and lemmaticized, which is nearly impossible for such large corpora. The tools and techniques outlined in this thesis do not directly deal with the process of parsing.

# 5 Implications for other less commonly taught languages

## 5.1 Easy reproducibility and growth

# Appendix 1: Conversational Hebrew Vocabulary List (CHVL)

|    | LEMMA | FREQUENCY | RANGE | UDP |
|----|-------|-----------|-------|-----|
| 1  | הוא   | 23446109  | 43455 | 0.9480170256 |
| 2  | ל     | 5638813   | 43448 | 0.9420130373 |
| 3  | ה     | 9850733   | 43458 | 0.9292661347 |
| 4  | ב     | 4812778   | 43450 | 0.9292364865 |
| 5  | את    | 6846782   | 43426 | 0.9285176069 |
| 6  | לא    | 5272808   | 43433 | 0.9145688112 |
| 7  | ש     | 3880654   | 43439 | 0.9088900047 |
| 8  | של    | 3892328   | 43445 | 0.9067041511 |
| 9  | על    | 1766990   | 43430 | 0.904286502 |
| 10 | זה    | 5118759   | 43441 | 0.9015544613 |
| 11 | מה    | 2362419   | 43403 | 0.8922532708 |
| 12 | היה   | 2579370   | 43420 | 0.8909904417 |
| 13 | מ     | 1061614   | 43411 | 0.8890067276 |
| 14 | כול   | 1325676   | 43414 | 0.8860074112 |
| 15 | ו     | 1906717   | 43429 | 0.885270638 |
| 16 | יש    | 1069358   | 43376 | 0.8770543442 |
| 17 | עם    | 839575    | 43331 | 0.8668140052 |
| 18 | אם    | 861163    | 43321 | 0.8654587702 |
| 19 | ידע   | 1202416   | 43323 | 0.8586088804 |
| 20 | אבל   | 921757    | 42963 | 0.8519038846 |
| 21 | אמר   | 799835    | 43196 | 0.8515460134 |
| 22 | רק    | 580549    | 43306 | 0.8490225759 |
| 23 | עשה   | 957476    | 43311 | 0.8460669028 |
| 24 | רצה   | 905161    | 43202 | 0.8453711531 |
| 25 | יותר  | 519740    | 43206 | 0.8426501511 |
| 26 | דבר   | 549346    | 43192 | 0.8389916741 |
| 27 | אז    | 785143    | 43202 | 0.8317146818 |
| 28 | חשב   | 585499    | 43062 | 0.8311268353 |
| 29 | ראה   | 464852    | 43120 | 0.8276119303 |
| 30 | אין   | 376940    | 42895 | 0.826471392 |

| | LEMMA | FREQUENCY | RANGE | UDP |
|---|---|---|---|---|
| 31 | איך | 367902 | 42714 | 0.825284943 |
| 32 | זמן | 397979 | 43034 | 0.8227270095 |
| 33 | אחד | 447348 | 43074 | 0.8218430306 |
| 34 | שם | 511696 | 43109 | 0.820013059 |
| 35 | משהו | 424351 | 42768 | 0.8199292075 |
| 36 | צריך | 678461 | 43101 | 0.8173698432 |
| 37 | כך | 538120 | 43151 | 0.8172938964 |
| 38 | כמה | 327641 | 42552 | 0.8144761932 |
| 39 | אל | 548185 | 43249 | 0.8123221549 |
| 40 | עכשיו | 464746 | 42758 | 0.8106640852 |
| 41 | טוב | 947724 | 43291 | 0.8084645436 |
| 42 | יכול | 490428 | 43141 | 0.8064150537 |
| 43 | בא | 419823 | 43050 | 0.8047477721 |
| 44 | כמו | 388089 | 42849 | 0.8041853147 |
| 45 | גם | 321102 | 42702 | 0.8041830813 |
| 46 | כן | 1207533 | 43226 | 0.8041799654 |
| 47 | למה | 433036 | 42608 | 0.8024645921 |
| 48 | מן | 260131 | 42071 | 0.8022138497 |
| 49 | נכון | 394738 | 42700 | 0.8014874793 |
| 50 | מי | 373446 | 42688 | 0.80073736 |
| 51 | אחר | 255588 | 41924 | 0.7996141928 |
| 52 | נראה | 310681 | 42564 | 0.7980211954 |
| 53 | כ | 270578 | 42075 | 0.797586002 |
| 54 | פעם | 286598 | 42191 | 0.7952189434 |
| 55 | איש | 562845 | 42958 | 0.7942488824 |
| 56 | או | 412974 | 42796 | 0.7936468503 |
| 57 | הגיע | 267993 | 41984 | 0.7917680395 |
| 58 | עד | 218184 | 41190 | 0.7917160839 |
| 59 | עצמו | 205086 | 41000 | 0.7894097508 |
| 60 | דיבר | 289788 | 41648 | 0.7883758932 |
| 61 | הרבה | 214954 | 41188 | 0.7877135038 |
| 62 | לפני | 225776 | 41249 | 0.7876190347 |

| | LEMMA | FREQUENCY | RANGE | UDP |
|----|-------|-----------|-------|-----|
| 63 | כבר | 250376 | 41870 | 0.7861533936 |
| 64 | אולי | 316008 | 42239 | 0.7851973472 |
| 65 | דרך | 277379 | 41924 | 0.7851063904 |
| 66 | קרה | 298579 | 42161 | 0.7839807769 |
| 67 | עדיין | 208160 | 40811 | 0.7825677023 |
| 68 | עוד | 260354 | 42041 | 0.7824611913 |
| 69 | ניסה | 201709 | 40669 | 0.7812562648 |
| 70 | הבין | 178461 | 40099 | 0.7787652199 |
| 71 | אף | 224000 | 40829 | 0.7783418455 |
| 72 | עבר | 181166 | 40252 | 0.7771570269 |
| 73 | מישהו | 238416 | 40919 | 0.7759852576 |
| 74 | אפילו | 139866 | 38453 | 0.7743368394 |
| 75 | כאן | 623430 | 41759 | 0.773984761 |
| 76 | שמע | 171278 | 39499 | 0.7729531526 |
| 77 | נתן | 172041 | 39452 | 0.7726047557 |
| 78 | כש | 158697 | 38893 | 0.7723833043 |
| 79 | שוב | 157377 | 39393 | 0.7718752651 |
| 80 | בדיוק | 154089 | 38931 | 0.7716276664 |
| 81 | כדי | 306213 | 41152 | 0.770282592 |
| 82 | אחת | 154130 | 39146 | 0.7695547607 |
| 83 | מקום | 198165 | 40314 | 0.7680001789 |
| 84 | חזר | 202999 | 40579 | 0.7671689493 |
| 85 | יצא | 162483 | 39369 | 0.7662316518 |
| 86 | התחיל | 99971 | 35015 | 0.7652934349 |
| 87 | בטוח | 141725 | 38426 | 0.7652236472 |
| 88 | במקום | 55050 | 27901 | 0.7643551072 |
| 89 | יום | 266260 | 41382 | 0.7642633607 |
| 90 | הספיק | 75388 | 31940 | 0.7641498356 |
| 91 | שב | 95518 | 34110 | 0.7641391656 |
| 92 | באמת | 279723 | 41591 | 0.7624500121 |
| 93 | אחרי | 138136 | 37831 | 0.7622924022 |
| 94 | וה | 54161 | 26863 | 0.7622073428 |

|  | LEMMA | FREQUENCY | RANGE | UDP |
|---|---|---|---|---|
| 95 | שני | 168829 | 39248 | 0.7621461876 |
| 96 | חיים | 256770 | 41514 | 0.7618890655 |
| 97 | תמיד | 138710 | 37943 | 0.7616653194 |
| 98 | לקח | 109842 | 35652 | 0.7609634126 |
| 99 | קשה | 100832 | 34923 | 0.7608985492 |
| 100 | הכיל | 183652 | 34316 | 0.7606625397 |
| 101 | לפחות | 46736 | 25727 | 0.760652816 |
| 102 | כמעט | 54522 | 27109 | 0.7599555508 |
| 103 | קודם | 71584 | 31900 | 0.7598150177 |
| 104 | רגע | 220597 | 40784 | 0.7597664595 |
| 105 | המשיך | 74650 | 30977 | 0.7595020775 |
| 106 | חייב | 271131 | 40994 | 0.7593203249 |
| 107 | הביא | 117845 | 36660 | 0.759057256 |
| 108 | לשם | 43631 | 24044 | 0.7586111453 |
| 109 | קיווה | 65759 | 29695 | 0.7585663791 |
| 110 | מדי | 117371 | 34092 | 0.7581777015 |
| 111 | אחרון | 115785 | 36237 | 0.7580135334 |
| 112 | קרוב | 64765 | 29164 | 0.7578402237 |
| 113 | שמר | 81313 | 31883 | 0.7575578845 |
| 114 | עלה | 63374 | 28020 | 0.7575210474 |
| 115 | קרא | 140013 | 37324 | 0.7572790956 |
| 116 | מלא | 48054 | 25189 | 0.7572739997 |
| 117 | איזה | 179147 | 39606 | 0.7570976824 |
| 118 | שינה | 64559 | 29600 | 0.7570520698 |
| 119 | השאיר | 52040 | 26619 | 0.7569118991 |
| 120 | יחיד | 73740 | 31360 | 0.7568256124 |
| 121 | קצת | 175325 | 38554 | 0.7568215267 |
| 122 | חיכה | 57163 | 27476 | 0.756807387 |
| 123 | איתך | 50853 | 25409 | 0.7565141919 |
| 124 | עמד | 103159 | 34456 | 0.7563573904 |
| 125 | אי | 146976 | 38291 | 0.7559666454 |
| 126 | חשוב | 70535 | 29954 | 0.7559403175 |

|     | LEMMA | FREQUENCY | RANGE | UDP |
| --- | --- | --- | --- | --- |
| 127 | חוץ | 90712 | 32860 | 0.7559332039 |
| 128 | הכיר | 131321 | 36335 | 0.7556745826 |
| 129 | שאל | 106181 | 34201 | 0.7553535119 |
| 130 | נעשה | 70421 | 30093 | 0.7553442662 |
| 131 | מאוחר | 52191 | 25779 | 0.7550953779 |
| 132 | מוכן | 115301 | 35648 | 0.7549545997 |
| 133 | נשמע | 69814 | 30269 | 0.7548470289 |
| 134 | נכנס | 90429 | 33029 | 0.7548084432 |
| 135 | חלק | 87778 | 32785 | 0.754725404 |
| 136 | מבין | 50540 | 25479 | 0.7545120899 |
| 137 | נ | 46983 | 24902 | 0.7544822115 |
| 138 | אמור | 89700 | 33337 | 0.753689558 |
| 139 | קל | 44150 | 24171 | 0.7534149166 |
| 140 | ילך | 59244 | 27026 | 0.7527902619 |
| 141 | בכלל | 60990 | 27698 | 0.7527396448 |
| 142 | אלה | 233644 | 38074 | 0.7527022741 |
| 143 | כלל | 53987 | 25929 | 0.7525384927 |
| 144 | שום | 146411 | 36788 | 0.7521875367 |
| 145 | גרם | 110859 | 35168 | 0.7520426531 |
| 146 | הפסיק | 64442 | 28808 | 0.7520203665 |
| 147 | הפעם | 40967 | 23297 | 0.7519715819 |
| 148 | מתי | 56953 | 26975 | 0.7516978167 |
| 149 | הת | 46255 | 24614 | 0.7515972282 |
| 150 | סיים | 49967 | 25437 | 0.7515551798 |
| 151 | שכח | 45437 | 23959 | 0.7511920907 |
| 152 | איתי | 46761 | 24239 | 0.7509541014 |
| 153 | בין | 79758 | 31121 | 0.7509185448 |
| 154 | עבד | 119322 | 35370 | 0.75091488 |
| 155 | נשאר | 98287 | 33880 | 0.7505987967 |
| 156 | האמין | 133888 | 37080 | 0.7502095967 |
| 157 | בחר | 54812 | 26160 | 0.7502076842 |
| 158 | אכפת | 72821 | 29977 | 0.7500879199 |

|  | LEMMA | FREQUENCY | RANGE | UDP |
|---|---|---|---|---|
| 159 | קיבל | 243889 | 40776 | 0.7500171151 |
| 160 | ישב | 45291 | 23369 | 0.7499018222 |
| 161 | רע | 79634 | 30810 | 0.7497801544 |
| 162 | הוציא | 51586 | 25538 | 0.7494769452 |
| 163 | עזר | 166920 | 38806 | 0.7494677516 |
| 164 | בעיה | 133192 | 36380 | 0.749376319 |
| 165 | הראה | 41990 | 22842 | 0.7493114348 |
| 166 | גדול | 182308 | 39208 | 0.7490466276 |
| 167 | כוונה | 42481 | 23367 | 0.7489349235 |
| 168 | אעשה | 39467 | 22636 | 0.7484708894 |
| 169 | צדק | 57882 | 27334 | 0.7484672402 |
| 170 | שנה | 219155 | 39679 | 0.747974035 |
| 171 | אלא | 41737 | 22893 | 0.747954367 |
| 172 | ביקש | 71177 | 28968 | 0.7478055486 |
| 173 | חסר | 48992 | 24559 | 0.7475294553 |
| 174 | סוף | 87864 | 31625 | 0.7475109849 |
| 175 | תודה | 269458 | 40779 | 0.7473604624 |
| 176 | עובד | 80911 | 30543 | 0.7471789502 |
| 177 | גרוע | 49393 | 25267 | 0.7469436121 |
| 178 | הניח | 100464 | 33988 | 0.7468575081 |
| 179 | השתמש | 76694 | 30709 | 0.7457989917 |
| 180 | מושג | 43162 | 23299 | 0.7452825579 |
| 181 | היום | 161107 | 37991 | 0.7452401042 |
| 182 | בלי | 82463 | 31103 | 0.7451788216 |
| 183 | בבקש | 45787 | 23109 | 0.7450939608 |
| 184 | הפך | 79849 | 30971 | 0.7450116709 |
| 185 | חץ | 55858 | 25677 | 0.7449420523 |
| 186 | הבטיח | 49463 | 24500 | 0.744862888 |
| 187 | ברור | 55491 | 25902 | 0.7448479539 |
| 188 | מזל | 59141 | 26640 | 0.7448047006 |
| 189 | תן | 126600 | 35753 | 0.7444971437 |
| 190 | אופן | 63422 | 26925 | 0.7441921309 |

|     | LEMMA | FREQUENCY | RANGE | UDP |
| --- | --- | --- | --- | --- |
| 191 | לאָן | 61827 | 27273 | 0.7441448176 |
| 192 | מאוד | 242051 | 40437 | 0.7441282577 |
| 193 | הסתכל | 42945 | 22239 | 0.7436579639 |
| 194 | עניין | 115658 | 34716 | 0.7435984846 |
| 195 | איבד | 48731 | 24396 | 0.7435728837 |
| 196 | מעולם | 77133 | 28830 | 0.7435419138 |
| 197 | במשך | 47299 | 23392 | 0.7435258738 |
| 198 | קטן | 148715 | 37651 | 0.7433888266 |
| 199 | רעיון | 60631 | 26575 | 0.7430392985 |
| 200 | הלך | 638998 | 43040 | 0.7429720343 |
| 201 | שתי | 39982 | 21679 | 0.7427308035 |
| 202 | סדר | 834217 | 42733 | 0.7426789343 |
| 203 | החזיק | 54710 | 25545 | 0.7425742773 |
| 204 | עין | 68827 | 28544 | 0.7423705866 |
| 205 | שונה | 45562 | 23096 | 0.7423548847 |
| 206 | מצב | 84165 | 31396 | 0.7422339254 |
| 207 | שה | 34689 | 20798 | 0.7422119424 |
| 208 | הצטער | 190553 | 38552 | 0.7421879954 |
| 209 | חדש | 142727 | 37387 | 0.7420182822 |
| 210 | השיג | 48373 | 23811 | 0.7419727547 |
| 211 | הקשיב | 44182 | 22621 | 0.741928355 |
| 212 | הגיד | 152422 | 35355 | 0.7417440599 |
| 213 | שעה | 121973 | 34939 | 0.7417170913 |
| 214 | מקרה | 98210 | 32986 | 0.7413964143 |
| 215 | שנייה | 58550 | 26144 | 0.7413615068 |
| 216 | עזב | 85029 | 31038 | 0.7412648676 |
| 217 | לבד | 49030 | 24642 | 0.7410606074 |
| 218 | ישן | 62874 | 27460 | 0.7408563335 |
| 219 | ודה | 36919 | 21077 | 0.7407312507 |
| 220 | פנים | 83000 | 31491 | 0.7407147299 |
| 221 | הזדמנות | 42643 | 22444 | 0.7404669742 |
| 222 | רציני | 39596 | 21409 | 0.7402443896 |

|     | LEMMA | FREQUENCY | RANGE | UDP |
| --- | --- | --- | --- | --- |
| 223 | שבוע | 92062 | 30773 | 0.7400989553 |
| 224 | עזרה | 42332 | 23031 | 0.7400542463 |
| 225 | חי | 55083 | 25663 | 0.7399883645 |
| 226 | חיפש | 82840 | 31018 | 0.7395963524 |
| 227 | בהחלט | 45868 | 22666 | 0.7394598613 |
| 228 | שאלה | 67213 | 27781 | 0.7393956671 |
| 229 | אמיתי | 76011 | 29258 | 0.7393012975 |
| 230 | נגמר | 42863 | 22750 | 0.7389565606 |
| 231 | זכר | 77021 | 29885 | 0.7388757029 |
| 232 | בטח | 124636 | 35618 | 0.7386129576 |
| 233 | שניים | 37443 | 20788 | 0.7385213779 |
| 234 | יד | 162906 | 37277 | 0.738349124 |
| 235 | מייד | 43011 | 21960 | 0.7382879283 |
| 236 | אכל | 105627 | 34354 | 0.7381245586 |
| 237 | איפה | 199458 | 38203 | 0.7381112111 |
| 238 | מצא | 206740 | 39632 | 0.738090718 |
| 239 | שלח | 56416 | 25495 | 0.7379336542 |
| 240 | כנראה | 54321 | 24952 | 0.737712106 |
| 241 | פתח | 43581 | 22423 | 0.7377085338 |
| 242 | הנה | 176643 | 38711 | 0.7376343479 |
| 243 | מעל | 45180 | 22991 | 0.737591866 |
| 244 | לעולם | 71132 | 28630 | 0.737578682 |
| 245 | ככה | 50618 | 23796 | 0.7373942767 |
| 246 | חודש | 58403 | 24849 | 0.7372285979 |
| 247 | חזק | 57852 | 25977 | 0.7372138662 |
| 248 | נחמד | 70052 | 27659 | 0.7371124211 |
| 249 | כלום | 93067 | 31268 | 0.7368065449 |
| 250 | לפעמים | 38191 | 21002 | 0.7367866745 |
| 251 | כמובן | 79627 | 29256 | 0.7363298635 |
| 252 | דקה | 76463 | 28622 | 0.7362924936 |
| 253 | קורה | 44598 | 23024 | 0.7362295076 |
| 254 | פרק | 47729 | 27236 | 0.7362027981 |

|     | LEMMA | FREQUENCY | RANGE | UDP |
| --- | --- | --- | --- | --- |
| 255 | מילה | 44487 | 22043 | 0.7358839517 |
| 256 | בעוד | 38921 | 21239 | 0.7358829875 |
| 257 | מספיק | 32456 | 20251 | 0.7357807778 |
| 258 | שאר | 36934 | 20718 | 0.7353728441 |
| 259 | נוסף | 58792 | 25572 | 0.7352874904 |
| 260 | שמח | 77624 | 30078 | 0.7352431471 |
| 261 | יפה | 83352 | 29667 | 0.7351487295 |
| 262 | הציע | 39769 | 21157 | 0.7351081906 |
| 263 | הודה | 65515 | 26896 | 0.7350893805 |
| 264 | סיכוי | 41951 | 22067 | 0.7349590884 |
| 265 | צורה | 44400 | 22438 | 0.7346784759 |
| 266 | הצליח | 59142 | 25896 | 0.7344972776 |
| 267 | חבר | 186844 | 38452 | 0.734147728 |
| 268 | פחות | 33427 | 19982 | 0.7341361879 |
| 269 | לגמרי | 48709 | 23158 | 0.7339458687 |
| 270 | סוג | 57056 | 24733 | 0.7337693895 |
| 271 | חזרה | 58804 | 25849 | 0.7337249196 |
| 272 | אהב | 267134 | 40244 | 0.7336376602 |
| 273 | ירד | 37336 | 20241 | 0.7335996357 |
| 274 | שכן | 36998 | 21089 | 0.7334231512 |
| 275 | לב | 112118 | 34293 | 0.7333369915 |
| 276 | פגע | 47572 | 23466 | 0.7333269908 |
| 277 | כדאי | 73858 | 28723 | 0.733054925 |
| 278 | שלוש | 51076 | 22791 | 0.7330223662 |
| 279 | בתוך | 43425 | 21982 | 0.7328734462 |
| 280 | ליד | 33666 | 20052 | 0.7327245222 |
| 281 | בדק | 55246 | 24773 | 0.7326923297 |
| 282 | עבודה | 179543 | 37349 | 0.7323951341 |
| 283 | מחר | 63071 | 25269 | 0.7323940254 |
| 284 | נמצא | 120923 | 34685 | 0.7323215003 |
| 285 | בית | 327260 | 40888 | 0.7321608087 |
| 286 | הרגיש | 148724 | 36977 | 0.7318832591 |

| | LEMMA | FREQUENCY | RANGE | UDP |
|---|---|---|---|---|
| 287 | בלתי | 44157 | 21790 | 0.7318474784 |
| 288 | אפשר | 111812 | 33563 | 0.7316673429 |
| 289 | עצר | 60763 | 26180 | 0.731573767 |
| 290 | למד | 56993 | 24951 | 0.7315141467 |
| 291 | ראש | 111934 | 33647 | 0.7315138677 |
| 292 | קשר | 146801 | 36266 | 0.7315131651 |
| 293 | דעה | 38500 | 20700 | 0.731489195 |
| 294 | הביתה | 83284 | 29329 | 0.731299241 |
| 295 | מהר | 65270 | 27051 | 0.7312323278 |
| 296 | קח | 46087 | 21743 | 0.7312303244 |
| 297 | פשוט | 276544 | 40438 | 0.7309931605 |
| 298 | סיפר | 134031 | 35660 | 0.7308407835 |
| 299 | אמת | 81307 | 29720 | 0.7305176277 |
| 300 | תראה | 100316 | 31667 | 0.7301189215 |
| 301 | החוצה | 48959 | 22603 | 0.729862538 |
| 302 | די | 92028 | 31259 | 0.7297169818 |
| 303 | שלושה | 41876 | 20973 | 0.7296060878 |
| 304 | רב | 104109 | 32606 | 0.72955377 |
| 305 | סלח | 43823 | 21207 | 0.7295246895 |
| 306 | הצלחה | 32762 | 19333 | 0.7294969843 |
| 307 | סתם | 43930 | 21609 | 0.7294567963 |
| 308 | רגיל | 33532 | 19572 | 0.729249545 |
| 309 | סיבה | 113796 | 34419 | 0.7291835982 |
| 310 | הכי | 92936 | 31255 | 0.7291680678 |
| 311 | למעשה | 51291 | 23015 | 0.7290616868 |
| 312 | התכוון | 135096 | 35138 | 0.7288138228 |
| 313 | נקודה | 44903 | 21657 | 0.7285867446 |
| 314 | בבקשה | 159524 | 36882 | 0.7285685947 |
| 315 | בוקר | 99172 | 30693 | 0.7281264133 |
| 316 | לכן | 49933 | 22825 | 0.7281093816 |
| 317 | אלי | 37067 | 19944 | 0.7280241905 |
| 318 | קנה | 53850 | 23964 | 0.7280068208 |

|     | LEMMA | FREQUENCY | RANGE | UDP |
| --- | --- | --- | --- | --- |
| 319 | תפס | 36698 | 20032 | 0.7278323204 |
| 320 | מוזר | 60733 | 25833 | 0.7277972419 |
| 321 | גש | 35553 | 19488 | 0.7277121856 |
| 322 | בשביל | 141977 | 35703 | 0.7276218143 |
| 323 | עסק | 91416 | 29969 | 0.7268211845 |
| 324 | יחד | 45480 | 21633 | 0.7266399146 |
| 325 | אוכל | 78950 | 29281 | 0.7265271334 |
| 326 | אתן | 42718 | 21524 | 0.7263799243 |
| 327 | כאילו | 81780 | 30119 | 0.7262683135 |
| 328 | מיוחד | 38008 | 20056 | 0.7261426774 |
| 329 | חושבת | 48123 | 22479 | 0.7259277175 |
| 330 | בגלל | 121582 | 33952 | 0.7253824561 |
| 331 | תרא | 63828 | 25651 | 0.725201538 |
| 332 | שילם | 51446 | 22382 | 0.7251261161 |
| 333 | התראה | 51727 | 22837 | 0.7249730967 |
| 334 | בוא | 183452 | 38124 | 0.7249015648 |
| 335 | צעיר | 46982 | 21498 | 0.7248317714 |
| 336 | ביותר | 88518 | 28965 | 0.7245400689 |
| 337 | למעלה | 44407 | 20843 | 0.7244146559 |
| 338 | התקשר | 71098 | 25533 | 0.7242354425 |
| 339 | טעות | 33998 | 19445 | 0.7241870178 |
| 340 | בחור | 68431 | 25592 | 0.7240913472 |
| 341 | ציפה | 28977 | 18426 | 0.7240666724 |
| 342 | זאת | 458405 | 41920 | 0.7238993223 |
| 343 | נהג | 44367 | 20845 | 0.7238395601 |
| 344 | מצטער | 34301 | 19572 | 0.7237249771 |
| 345 | ארוך | 29358 | 18330 | 0.7236610723 |
| 346 | טיפל | 35412 | 19710 | 0.7236317511 |
| 347 | גבוה | 33710 | 18756 | 0.7235971022 |
| 348 | החזיר | 36007 | 20104 | 0.7235659301 |
| 349 | העליי | 77454 | 24793 | 0.723488197 |
| 350 | לאחר | 55784 | 23190 | 0.7231869589 |

|     | LEMMA | FREQUENCY | RANGE | UDP |
| --- | --- | --- | --- | --- |
| 351 | הסכים | 32794 | 18770 | 0.7228958514 |
| 352 | שיחה | 38827 | 19897 | 0.7228810085 |
| 353 | פחד | 51801 | 22795 | 0.7227687501 |
| 354 | כי | 277322 | 38980 | 0.7226128673 |
| 355 | ניתן | 33218 | 18888 | 0.7221544319 |
| 356 | מוקדם | 28974 | 18294 | 0.7219771836 |
| 357 | מת | 209263 | 39030 | 0.7219628992 |
| 358 | יכולת | 32277 | 19032 | 0.7219402362 |
| 359 | צד | 34626 | 19226 | 0.7218224721 |
| 360 | נורא | 38503 | 19735 | 0.7218197452 |
| 361 | חכה | 68492 | 25911 | 0.7216722253 |
| 362 | תדאג | 31904 | 19026 | 0.7213861779 |
| 363 | למען | 38847 | 19906 | 0.721304746 |
| 364 | כפי | 40865 | 19918 | 0.7200534595 |
| 365 | אתמול | 45515 | 21452 | 0.7199803257 |
| 366 | ערב | 57377 | 23196 | 0.7199797507 |
| 367 | מספר | 72795 | 25912 | 0.7198548978 |
| 368 | בעל | 76549 | 27783 | 0.7197942705 |
| 369 | חמש | 38896 | 19160 | 0.7193984301 |
| 370 | מאשר | 28195 | 17633 | 0.7179471355 |
| 371 | תחת | 31938 | 17993 | 0.7178875479 |
| 372 | כרגע | 39722 | 19976 | 0.7178318886 |
| 373 | שווה | 30611 | 18112 | 0.7175492795 |
| 374 | לילה | 169318 | 35873 | 0.7173176145 |
| 375 | שקט | 34966 | 18622 | 0.7171228943 |
| 376 | הסביר | 28968 | 17836 | 0.7170220435 |
| 377 | העביר | 30214 | 17975 | 0.7167067263 |
| 378 | זו | 389942 | 38399 | 0.7164040078 |
| 379 | י | 42114 | 20366 | 0.7157420446 |
| 380 | מסוגל | 35529 | 18869 | 0.7156238736 |
| 381 | הוריד | 32852 | 18273 | 0.7153206863 |
| 382 | חושב | 92894 | 28956 | 0.7152910962 |

|  | LEMMA | FREQUENCY | RANGE | UDP |
|-----|-------|-----------|-------|-----|
| 383 | שאמר | 25388 | 17308 | 0.7152421478 |
| 384 | צריכה | 46925 | 21649 | 0.7152074257 |
| 385 | הבחור | 55510 | 22331 | 0.7150010543 |
| 386 | ללא | 61731 | 24590 | 0.7148350498 |
| 387 | נעלם | 39306 | 19993 | 0.7147805669 |
| 388 | עובדה | 28205 | 17294 | 0.7147342763 |
| 389 | סיפור | 69616 | 26111 | 0.7147277465 |
| 390 | חדר | 109017 | 31987 | 0.7141118947 |
| 391 | כבוד | 51623 | 21817 | 0.7139722358 |
| 392 | נגע | 45566 | 21069 | 0.7135320108 |
| 393 | בחייך | 48561 | 20788 | 0.7134925247 |
| 394 | סליחה | 72672 | 25859 | 0.7134712582 |
| 395 | לגבי | 63139 | 23852 | 0.7133698597 |
| 396 | מטה | 46518 | 20269 | 0.7133580286 |
| 397 | רוח | 55806 | 23413 | 0.7131913761 |
| 398 | בקרוב | 29234 | 18039 | 0.7130681253 |
| 399 | האליי | 38658 | 18630 | 0.7128121409 |
| 400 | דלת | 60595 | 24076 | 0.7127045891 |
| 401 | הכין | 32161 | 18648 | 0.7126869215 |
| 402 | דאג | 32783 | 18349 | 0.7124880102 |
| 403 | אית | 26211 | 17189 | 0.7124481146 |
| 404 | שיחק | 58699 | 23577 | 0.7120624872 |
| 405 | אפשרי | 30242 | 17658 | 0.7120244767 |
| 406 | אדם | 200100 | 38089 | 0.7118554468 |
| 407 | אצל | 31511 | 17589 | 0.7113908312 |
| 408 | לפ | 31149 | 17418 | 0.7112366987 |
| 409 | ממש | 190693 | 37369 | 0.7108608252 |
| 410 | נהדר | 77117 | 26285 | 0.7107238238 |
| 411 | נגד | 38107 | 18684 | 0.7105908384 |
| 412 | רחוק | 28613 | 17301 | 0.7105602372 |
| 413 | ביחד | 40371 | 19338 | 0.710509941 |
| 414 | כאב | 42011 | 19849 | 0.7098332821 |

|     | LEMMA | FREQUENCY | RANGE | UDP |
| --- | --- | --- | --- | --- |
| 415 | כיוון | 33585 | 17992 | 0.7097300266 |
| 416 | רוב | 30663 | 17127 | 0.7082659531 |
| 417 | ארוחה | 35455 | 18292 | 0.7078307699 |
| 418 | אלך | 26957 | 17337 | 0.7076901203 |
| 419 | אישי | 27779 | 16623 | 0.7075250817 |
| 420 | מעבר | 28859 | 16944 | 0.7074280407 |
| 421 | עלול | 30722 | 17521 | 0.707404867 |
| 422 | תורגם | 35600 | 19872 | 0.7073848364 |
| 423 | הופיע | 28375 | 16743 | 0.707236713 |
| 424 | בנה | 30064 | 17246 | 0.707115285 |
| 425 | נסע | 51262 | 20354 | 0.7064758117 |
| 426 | עולם | 106212 | 31680 | 0.7063031666 |
| 427 | זהו | 88493 | 28003 | 0.7057294006 |
| 428 | שלום | 134482 | 34158 | 0.705414933 |
| 429 | משך | 26666 | 16544 | 0.7054148652 |
| 430 | ערך | 27833 | 16806 | 0.7052770346 |
| 431 | מאחורי | 25711 | 16436 | 0.7046718771 |
| 432 | שנא | 31478 | 17318 | 0.7046126398 |
| 433 | בגד | 31315 | 17232 | 0.7043811871 |
| 434 | יצר | 38263 | 18486 | 0.7037014243 |
| 435 | א | 47853 | 20598 | 0.7036041426 |
| 436 | חברה | 112610 | 31444 | 0.7035715971 |
| 437 | תוכנית | 77143 | 26625 | 0.7034012997 |
| 438 | ניצח | 43784 | 19665 | 0.7033827166 |
| 439 | כתב | 68326 | 23615 | 0.7032628794 |
| 440 | תמונה | 52029 | 20845 | 0.703052461 |
| 441 | הזכיר | 22880 | 15934 | 0.7028798911 |
| 442 | מוות | 67543 | 24487 | 0.702866482 |
| 443 | בערך | 26705 | 16236 | 0.7028281894 |
| 444 | גר | 37657 | 17985 | 0.7025931732 |
| 445 | אמצע | 23722 | 15917 | 0.7025409909 |
| 446 | מתחת | 27256 | 16600 | 0.7025043494 |

| | LEMMA | FREQUENCY | RANGE | UDP |
|---|---|---|---|---|
| 447 | לחץ | 33426 | 17264 | 0.7024540764 |
| 448 | לעזאזל | 68624 | 24177 | 0.7023060468 |
| 449 | טלפון | 66712 | 22700 | 0.7022575358 |
| 450 | הגן | 34085 | 17830 | 0.7021851713 |
| 451 | התאים | 24649 | 16114 | 0.702067948 |
| 452 | הכניס | 24328 | 15835 | 0.7020622143 |
| 453 | התמודד | 27809 | 16712 | 0.7018327426 |
| 454 | נפגש | 26092 | 16054 | 0.7015122968 |
| 455 | עבור | 85841 | 24695 | 0.7007542386 |
| 456 | בן | 270834 | 40029 | 0.7002527794 |
| 457 | מצחיק | 31729 | 16994 | 0.7001774739 |
| 458 | מעט | 28455 | 16294 | 0.7000108207 |
| 459 | זוכר | 30267 | 16942 | 0.6996247424 |
| 460 | קיים | 28007 | 16353 | 0.6995253097 |
| 461 | הציל | 44028 | 19759 | 0.6992853964 |
| 462 | הזמין | 27482 | 16304 | 0.6992744887 |
| 463 | למרות | 26167 | 16043 | 0.6992328041 |
| 464 | אקח | 23209 | 15773 | 0.6991071693 |
| 465 | איתן | 29262 | 16540 | 0.69897796 |
| 466 | לפי | 29227 | 16438 | 0.698903123 |
| 467 | סימן | 31504 | 16983 | 0.6985069202 |
| 468 | לבש | 29275 | 16629 | 0.6984565372 |
| 469 | ספק | 25278 | 15635 | 0.6979374643 |
| 470 | בת | 93599 | 28383 | 0.6975583872 |
| 471 | במיוחד | 22176 | 15147 | 0.6971721905 |
| 472 | מחדש | 27450 | 16186 | 0.6971600739 |
| 473 | התחלה | 21710 | 15161 | 0.6967019141 |
| 474 | רחוב | 39941 | 17507 | 0.6963750102 |
| 475 | משפחה | 104525 | 29823 | 0.6963376646 |
| 476 | הערב | 45069 | 18671 | 0.6961500605 |
| 477 | אלוהים | 210842 | 35618 | 0.6959731886 |
| 478 | קצר | 22297 | 15089 | 0.6959566078 |

|     | LEMMA | FREQUENCY | RANGE | UDP |
| --- | --- | --- | --- | --- |
| 479 | עזאזל | 100331 | 28613 | 0.695819585 |
| 480 | הצטרך | 22495 | 15287 | 0.6955729046 |
| 481 | טיפש | 28604 | 16103 | 0.695523551 |
| 482 | אסור | 29549 | 16451 | 0.6955129843 |
| 483 | החלטה | 28762 | 15978 | 0.6955030906 |
| 484 | זה | 28681 | 16105 | 0.6953489373 |
| 485 | הודעה | 32800 | 16692 | 0.6952133673 |
| 486 | יופי | 38738 | 18075 | 0.6948367181 |
| 487 | גבר | 119470 | 31398 | 0.6946615752 |
| 488 | נקרא | 22878 | 15250 | 0.6942184777 |
| 489 | סביבה | 24446 | 15259 | 0.6941596047 |
| 490 | אור | 30970 | 16257 | 0.6939727559 |
| 491 | חוק | 38433 | 17001 | 0.6938093837 |
| 492 | אח | 62249 | 22129 | 0.6937258948 |
| 493 | גנב | 36703 | 17273 | 0.6935358098 |
| 494 | משרד | 53856 | 19641 | 0.6933660069 |
| 495 | החליט | 22521 | 14991 | 0.6932626846 |
| 496 | מערכת | 45814 | 18981 | 0.6930330541 |
| 497 | נפל | 24793 | 15183 | 0.6929591896 |
| 498 | מושלם | 25446 | 15526 | 0.692932052 |
| 499 | שתיים | 29266 | 15645 | 0.6927257115 |
| 500 | הוטרף | 26467 | 15580 | 0.6923941155 |
| 501 | ן | 23039 | 15104 | 0.6917846177 |
| 502 | העדיף | 21524 | 14882 | 0.6914479499 |
| 503 | ספר | 126021 | 32397 | 0.6911662598 |
| 504 | מהלך | 27210 | 15503 | 0.6908194406 |
| 505 | קטע | 31794 | 16392 | 0.6900902886 |
| 506 | טעם | 24740 | 15311 | 0.6900390323 |
| 507 | ניסיון | 22094 | 14770 | 0.6900377888 |
| 508 | בתור | 26690 | 15608 | 0.6900352933 |
| 509 | מוצא | 22210 | 14629 | 0.6897565967 |
| 510 | נהנה | 24566 | 15334 | 0.6895847836 |

65

|     | LEMMA | FREQUENCY | RANGE | UDP |
|-----|-------|-----------|-------|-----|
| 511 | מין | 31053 | 15961 | 0.6895422947 |
| 512 | שירות | 28723 | 15767 | 0.6894819475 |
| 513 | צעד | 26512 | 15122 | 0.6892665452 |
| 514 | נפלא | 33808 | 16076 | 0.6891480192 |
| 515 | גוף | 32953 | 16491 | 0.6889814462 |
| 516 | קול | 43516 | 17953 | 0.6888180192 |
| 517 | אדיר | 42280 | 17816 | 0.6887415417 |
| 518 | חדשות | 25677 | 15380 | 0.6887395018 |
| 519 | תפקיד | 30626 | 15537 | 0.6885656229 |
| 520 | צהריים | 25563 | 14887 | 0.6884474358 |
| 521 | אראה | 21223 | 14588 | 0.6883200942 |
| 522 | תשובה | 23630 | 14878 | 0.688264663 |
| 523 | חלה | 60204 | 21244 | 0.6882338866 |
| 524 | סבל | 23693 | 14809 | 0.6881210183 |
| 525 | זקוק | 28408 | 15699 | 0.6879757716 |
| 526 | גמור | 25377 | 14927 | 0.6879061051 |
| 527 | מים | 46803 | 18838 | 0.6878891266 |
| 528 | עיר | 84067 | 26088 | 0.6878717009 |
| 529 | רצינות | 24555 | 15123 | 0.6876706831 |
| 530 | והיי | 388785 | 36676 | 0.6876699482 |
| 531 | שן | 22289 | 15141 | 0.6876507207 |
| 532 | החליף | 22008 | 14641 | 0.6874974022 |
| 533 | בצד | 22200 | 14506 | 0.687494966 |
| 534 | גידי | 29006 | 15672 | 0.6874567503 |
| 535 | ברוך | 24841 | 14984 | 0.6873558061 |
| 536 | נעל | 30612 | 15835 | 0.6872947396 |
| 537 | הוביל | 24282 | 14940 | 0.6872409139 |
| 538 | צורך | 21780 | 14581 | 0.6870485965 |
| 539 | צוות | 60182 | 21080 | 0.6870229797 |
| 540 | ברח | 26786 | 15514 | 0.6866285704 |
| 541 | כוח | 65992 | 23354 | 0.6865035252 |
| 542 | נשק | 44828 | 18112 | 0.6864818854 |

| | LEMMA | FREQUENCY | RANGE | UDP |
|---|---|---|---|---|
| 543 | שׁוֹלחֹן | 26796 | 14963 | 0.686306115 |
| 544 | ככל | 21610 | 14377 | 0.6860225856 |
| 545 | הורה | 34611 | 16310 | 0.6859805942 |
| 546 | מטרה | 33106 | 15951 | 0.6859462123 |
| 547 | פנימה | 23950 | 14506 | 0.6854226527 |
| 548 | גילה | 22066 | 14658 | 0.685287813 |
| 549 | הרשה | 21051 | 14105 | 0.6852823893 |
| 550 | חם | 25170 | 14682 | 0.6850446951 |
| 551 | מיטה | 28113 | 15212 | 0.6849563886 |
| 552 | ארבע | 25315 | 14298 | 0.6845052354 |
| 553 | אהבה | 48349 | 18292 | 0.6842906115 |
| 554 | ילד | 286526 | 39003 | 0.6842788489 |
| 555 | ישר | 21503 | 13788 | 0.684210625 |
| 556 | זוג | 29494 | 15370 | 0.6841469783 |
| 557 | ותק | 48434 | 18007 | 0.6839113203 |
| 558 | תוך | 22093 | 14170 | 0.6838908592 |
| 559 | נוח | 22057 | 14384 | 0.6837318084 |
| 560 | חשבתי | 19504 | 14490 | 0.6835461456 |
| 561 | מדינה | 39099 | 15768 | 0.6835264385 |
| 562 | סמך | 24286 | 14787 | 0.6834645923 |
| 563 | מול | 22733 | 14323 | 0.6831493689 |
| 564 | אב | 36237 | 16121 | 0.6830663561 |
| 565 | זכות | 25428 | 14256 | 0.6829699097 |
| 566 | כלומר | 42567 | 16988 | 0.6826894041 |
| 567 | יקר | 21816 | 14095 | 0.6826519838 |
| 568 | שחרר | 27115 | 15146 | 0.6826342435 |
| 569 | מידע | 37642 | 16139 | 0.6826077034 |
| 570 | ענה | 21309 | 14788 | 0.6821787425 |
| 571 | חשבון | 25497 | 14268 | 0.6813941773 |
| 572 | מאה | 24481 | 14247 | 0.6812393257 |
| 573 | הפריע | 19994 | 13961 | 0.6808298315 |
| 574 | אוויר | 29391 | 14943 | 0.6806467383 |

|     | LEMMA | FREQUENCY | RANGE | UDP |
| --- | --- | --- | --- | --- |
| 575 | פגישה | 31627 | 15003 | 0.6806137431 |
| 576 | הצטרף | 21644 | 14144 | 0.6804024728 |
| 577 | שלומך | 25313 | 13675 | 0.6803961452 |
| 578 | מוח | 33414 | 15815 | 0.6802982858 |
| 579 | אלו | 79895 | 21603 | 0.6801781522 |
| 580 | הדבר | 21069 | 14108 | 0.6800283772 |
| 581 | כדור | 58850 | 20345 | 0.6798777039 |
| 582 | הרגל | 23986 | 14088 | 0.6798187118 |
| 583 | מיני | 24270 | 14421 | 0.6796423144 |
| 584 | תקשיב | 24742 | 14245 | 0.679063459 |
| 585 | נעים | 22824 | 13761 | 0.6788524482 |
| 586 | מבט | 20802 | 13651 | 0.6787799978 |
| 587 | צפה | 20797 | 14047 | 0.6787504505 |
| 588 | מתוק | 30675 | 15261 | 0.6786662669 |
| 589 | חבל | 21785 | 13954 | 0.6786475839 |
| 590 | נשא | 23551 | 13857 | 0.6782422631 |
| 591 | הרג | 129020 | 29925 | 0.6782185649 |
| 592 | לתוך | 23236 | 14037 | 0.6781059414 |
| 593 | שייך | 22585 | 14040 | 0.6780237312 |
| 594 | הרס | 21537 | 14270 | 0.6779511674 |
| 595 | לבן | 26079 | 14048 | 0.6776427474 |
| 596 | שעשה | 18128 | 13540 | 0.6775737869 |
| 597 | המ | 27799 | 14494 | 0.6775318892 |
| 598 | הסתובב | 19667 | 13377 | 0.677064829 |
| 599 | בחורה | 28881 | 14657 | 0.6767060865 |
| 600 | פעולה | 24189 | 13885 | 0.6766103988 |
| 601 | נרגע | 26088 | 14117 | 0.6766056949 |
| 602 | סגר | 20917 | 13586 | 0.6763634772 |
| 603 | גדל | 20911 | 13562 | 0.6763225834 |
| 604 | היקח | 19466 | 13467 | 0.6760658873 |
| 605 | תפסיק | 23856 | 13940 | 0.6759811966 |
| 606 | אגיד | 20937 | 13702 | 0.6758237864 |

|     | LEMMA | FREQUENCY | RANGE | UDP |
| --- | --- | --- | --- | --- |
| 607 | מסוכן | 21915 | 13987 | 0.6758127032 |
| 608 | זרק | 21314 | 13618 | 0.6755505569 |
| 609 | תהי | 21587 | 13941 | 0.6755314781 |
| 610 | עשר | 23977 | 13697 | 0.6755251882 |
| 611 | חכם | 20336 | 13425 | 0.6752554862 |
| 612 | גברת | 64140 | 20248 | 0.6751813371 |
| 613 | דולר | 60965 | 18338 | 0.6751702352 |
| 614 | שינוי | 21790 | 13717 | 0.67490225 |
| 615 | ביצע | 23118 | 13684 | 0.6747570407 |
| 616 | שלם | 18514 | 13358 | 0.6747451015 |
| 617 | עלייך | 25562 | 14499 | 0.6747042683 |
| 618 | זונה | 40225 | 15443 | 0.6746088468 |
| 619 | מנהל | 31646 | 14360 | 0.6745423267 |
| 620 | חתיכה | 22389 | 13661 | 0.6744828553 |
| 621 | חומר | 27716 | 14014 | 0.6744228935 |
| 622 | רכב | 35362 | 15119 | 0.6742473557 |
| 623 | אחראי | 21545 | 13581 | 0.6736649998 |
| 624 | חתך | 22388 | 13556 | 0.6735182669 |
| 625 | ניהל | 19826 | 12930 | 0.6730627155 |
| 626 | משטרה | 62585 | 18340 | 0.6727316177 |
| 627 | צוחק | 27135 | 13873 | 0.6726623337 |
| 628 | עוזב | 19926 | 13289 | 0.6724123 |
| 629 | עונה | 21128 | 16314 | 0.6723301941 |
| 630 | רופא | 45488 | 16659 | 0.672024932 |
| 631 | מכר | 23577 | 13216 | 0.6719295775 |
| 632 | השתנה | 20353 | 13421 | 0.6718465454 |
| 633 | מפה | 29533 | 14189 | 0.6709086698 |
| 634 | עץ | 34139 | 14898 | 0.670832602 |
| 635 | כלא | 43491 | 15589 | 0.6708308059 |
| 636 | מהיר | 20903 | 13155 | 0.6705956003 |
| 637 | כרטיס | 29702 | 13967 | 0.6705263381 |
| 638 | אצטרך | 18209 | 13342 | 0.6704775286 |

|     | LEMMA | FREQUENCY | RANGE | UDP |
|-----|-------|-----------|-------|-----|
| 639 | פשע | 35547 | 14565 | 0.6703819783 |
| 640 | קבוצה | 40544 | 15813 | 0.6703108733 |
| 641 | המון | 22918 | 13256 | 0.6702285421 |
| 642 | כלשהו | 21146 | 13207 | 0.6701005457 |
| 643 | סכנה | 21596 | 13566 | 0.669462576 |
| 644 | מתוך | 19147 | 12809 | 0.6693526876 |
| 645 | שנוכל | 18269 | 13339 | 0.66928734 |
| 646 | קו | 25216 | 13233 | 0.6692700889 |
| 647 | הלוואה | 19687 | 13220 | 0.6690227782 |
| 648 | מסר | 21693 | 13107 | 0.6689513536 |
| 649 | יחסים | 26687 | 13752 | 0.6687970945 |
| 650 | מכונית | 78708 | 21273 | 0.6686906036 |
| 651 | וב | 17646 | 12447 | 0.6686623457 |
| 652 | ארץ | 40008 | 15669 | 0.6686397384 |
| 653 | הגיוני | 19235 | 13055 | 0.6685594793 |
| 654 | דם | 62041 | 20297 | 0.6685565848 |
| 655 | הדה | 21093 | 13101 | 0.668513652 |
| 656 | כיף | 25538 | 14072 | 0.6682454932 |
| 657 | עשוי | 22701 | 13352 | 0.6682395193 |
| 658 | העריך | 17831 | 12680 | 0.6679289203 |
| 659 | שליטה | 21269 | 13108 | 0.6678991286 |
| 660 | זכה | 24942 | 13316 | 0.667862798 |
| 661 | רמה | 20702 | 12815 | 0.6676860515 |
| 662 | אוי | 44899 | 16449 | 0.6674228528 |
| 663 | אפשרות | 19160 | 12787 | 0.6673955164 |
| 664 | שמחה | 20060 | 13077 | 0.6673158968 |
| 665 | פתוח | 18684 | 12827 | 0.6671985811 |
| 666 | שיר | 42679 | 15685 | 0.6670326267 |
| 667 | חופשי | 19901 | 12818 | 0.6670280318 |
| 668 | כסף | 135843 | 28087 | 0.6670047228 |
| 669 | רשימה | 26386 | 13403 | 0.6668214889 |
| 670 | פרטי | 20520 | 12665 | 0.6666343952 |

| | LEMMA | FREQUENCY | RANGE | UDP |
|---|---|---|---|---|
| 671 | אש | 37360 | 15071 | 0.6665901443 |
| 672 | חמוד | 25354 | 13620 | 0.666167358 |
| 673 | צא | 23703 | 13087 | 0.6659784686 |
| 674 | שש | 22812 | 12825 | 0.6658906289 |
| 675 | עתיד | 25814 | 13385 | 0.6655604852 |
| 676 | הלו | 28281 | 13386 | 0.6655441795 |
| 677 | נושא | 20248 | 12484 | 0.6654212333 |
| 678 | שחור | 29409 | 13552 | 0.6653889217 |
| 679 | משחק | 81151 | 24065 | 0.6653530879 |
| 680 | תיק | 36739 | 14271 | 0.6653430249 |
| 681 | גיל | 22207 | 12759 | 0.6652898669 |
| 682 | פעל | 19709 | 12705 | 0.6651821943 |
| 683 | איתה | 18666 | 12929 | 0.6649849948 |
| 684 | קפה | 27340 | 13261 | 0.6646339454 |
| 685 | עקב | 20069 | 12874 | 0.6640737607 |
| 686 | היכן | 48255 | 16013 | 0.6640217227 |
| 687 | החלק | 18140 | 12666 | 0.6639917196 |
| 688 | אזור | 24762 | 13129 | 0.663918207 |
| 689 | שטח | 24612 | 13205 | 0.6638577069 |
| 690 | חייך | 18082 | 12577 | 0.6637210853 |
| 691 | לחלוטין | 19290 | 12506 | 0.6634580465 |
| 692 | וואו | 43989 | 15605 | 0.6634340093 |
| 693 | עמוק | 18686 | 12460 | 0.6632497577 |
| 694 | נלחם | 29162 | 13749 | 0.6632355502 |
| 695 | גמר | 19924 | 12307 | 0.6630216875 |
| 696 | תגיד | 17813 | 12195 | 0.6629976726 |
| 697 | כאשר | 56166 | 15772 | 0.6623824349 |
| 698 | אחרת | 15921 | 12093 | 0.6622270427 |
| 699 | מסוים | 18260 | 12276 | 0.6616631782 |
| 700 | זקן | 25313 | 12838 | 0.6616223306 |
| 701 | דובר | 22191 | 12523 | 0.6606712046 |
| 702 | נצטרך | 17872 | 12496 | 0.6604553639 |

|     | LEMMA | FREQUENCY | RANGE | UDP |
| --- | --- | --- | --- | --- |
| 703 | תוצאה | 20866 | 12297 | 0.66018425 |
| 704 | לישון | 19558 | 12295 | 0.6597971878 |
| 705 | כבד | 18444 | 11980 | 0.6592738659 |
| 706 | חשש | 18628 | 12296 | 0.6591882722 |
| 707 | מעניין | 17128 | 12116 | 0.6590920683 |
| 708 | גב | 27010 | 12727 | 0.6590021389 |
| 709 | מצוין | 22819 | 12431 | 0.6589085092 |
| 710 | שקר | 20460 | 12361 | 0.6588475702 |
| 711 | ובכן | 188842 | 27119 | 0.6588335 |
| 712 | תקופה | 18466 | 11999 | 0.6588048374 |
| 713 | האשים | 17959 | 12138 | 0.6586872276 |
| 714 | ביי | 32434 | 12971 | 0.6586047704 |
| 715 | בדיקה | 30789 | 13232 | 0.6582379081 |
| 716 | תני | 19112 | 12338 | 0.6581687605 |
| 717 | התקרב | 17162 | 12090 | 0.6580893063 |
| 718 | פרט | 17794 | 11925 | 0.6577512017 |
| 719 | אידיוט | 21741 | 12174 | 0.6575986185 |
| 720 | מעמד | 18726 | 11840 | 0.6575800008 |
| 721 | פגש | 16502 | 11762 | 0.6572396301 |
| 722 | שהייה | 15521 | 11779 | 0.657085437 |
| 723 | הוכיח | 18670 | 11951 | 0.6568365934 |
| 724 | הבחורה | 23272 | 12398 | 0.6565952224 |
| 725 | שכנע | 17282 | 11976 | 0.6564856379 |
| 726 | רץ | 19090 | 11630 | 0.6564817288 |
| 727 | כה | 23858 | 12452 | 0.6562881351 |
| 728 | צבע | 22920 | 12078 | 0.6561303903 |
| 729 | חלום | 28440 | 12785 | 0.655806186 |
| 730 | חנות | 25012 | 12676 | 0.6558036879 |
| 731 | דין | 43969 | 14317 | 0.6555802175 |
| 732 | מחיר | 20145 | 11950 | 0.6554877141 |
| 733 | עדיף | 16851 | 11850 | 0.6551502507 |
| 734 | דירה | 31147 | 13077 | 0.6551130059 |

|     | LEMMA | FREQUENCY | RANGE | UDP |
| --- | --- | --- | --- | --- |
| 735 | הפחיד | 17331 | 11866 | 0.6550067758 |
| 736 | מתנה | 20836 | 12112 | 0.654993717 |
| 737 | מלחמה | 43193 | 14580 | 0.6549903666 |
| 738 | הגנה | 22583 | 12072 | 0.6544375918 |
| 739 | מרחק | 18451 | 11900 | 0.6543854146 |
| 740 | אדום | 23426 | 11969 | 0.6543603976 |
| 741 | שלט | 18549 | 12031 | 0.6540224021 |
| 742 | רגל | 18311 | 11455 | 0.6540123152 |
| 743 | עורך | 29916 | 12270 | 0.653985097 |
| 744 | תחנה | 25747 | 12378 | 0.6538947506 |
| 745 | סרט | 52329 | 16423 | 0.6537311197 |
| 746 | שבר | 17298 | 11673 | 0.653518665 |
| 747 | קדימה | 198854 | 34380 | 0.6534927601 |
| 748 | זיהה | 18212 | 11751 | 0.6529948422 |
| 749 | הוגן | 17267 | 11626 | 0.6528978377 |
| 750 | כלב | 43138 | 14561 | 0.6528336538 |
| 751 | שומר | 18922 | 11831 | 0.6528197914 |
| 752 | לכי | 22331 | 12188 | 0.6527805514 |
| 753 | ויתר | 17187 | 11732 | 0.6526623736 |
| 754 | לאחרונה | 16225 | 11795 | 0.6526518282 |
| 755 | יחידה | 23205 | 12092 | 0.6526375596 |
| 756 | תא | 26021 | 12524 | 0.6524951103 |
| 757 | אבא | 158901 | 29216 | 0.6524220653 |
| 758 | ביטחון | 20138 | 11698 | 0.6520770684 |
| 759 | יכל | 18556 | 11541 | 0.6517396914 |
| 760 | עלי | 15428 | 11528 | 0.651617516 |
| 761 | כלי | 20356 | 11784 | 0.6514161345 |
| 762 | בעצם | 18483 | 11350 | 0.6513389503 |
| 763 | משפט | 39156 | 13243 | 0.6513068635 |
| 764 | הסתיים | 16495 | 11512 | 0.6512381499 |
| 765 | חך | 18850 | 11749 | 0.6512331414 |
| 766 | עוזר | 18391 | 11674 | 0.6511786639 |

| | LEMMA | FREQUENCY | RANGE | UDP |
|---|---|---|---|---|
| 767 | אתר | 23175 | 12246 | 0.6510946499 |
| 768 | נפגע | 17792 | 11568 | 0.6510894301 |
| 769 | פ | 18936 | 11502 | 0.6509383997 |
| 770 | פה | 179198 | 32392 | 0.650795368 |
| 771 | נפטר | 17128 | 11360 | 0.6507900754 |
| 772 | חלון | 19129 | 11414 | 0.6507829235 |
| 773 | שלב | 19177 | 11500 | 0.6506977908 |
| 774 | אחי | 63906 | 17859 | 0.6505429884 |
| 775 | תלוי | 14728 | 11185 | 0.6502090406 |
| 776 | אמא | 167065 | 26720 | 0.6500226527 |
| 777 | בעצמך | 14318 | 11272 | 0.6500174782 |
| 778 | ההוא | 19265 | 11584 | 0.6499870609 |
| 779 | כעת | 43252 | 14172 | 0.6499634779 |
| 780 | שכר | 18688 | 11152 | 0.6496699478 |
| 781 | עצמך | 14313 | 11145 | 0.6493888189 |
| 782 | רצח | 60144 | 16162 | 0.6493052613 |
| 783 | הבחר | 16890 | 11356 | 0.6491876544 |
| 784 | בר | 18472 | 11298 | 0.6491741908 |
| 785 | מר | 146704 | 26570 | 0.6490224129 |
| 786 | תכנן | 15857 | 11352 | 0.6488646482 |
| 787 | שיעור | 22104 | 11679 | 0.6487376366 |
| 788 | הא | 34636 | 12766 | 0.6486768104 |
| 789 | התרחק | 17474 | 11527 | 0.6486502522 |
| 790 | מותק | 25611 | 12259 | 0.648618291 |
| 791 | שותף | 22502 | 11556 | 0.6485146259 |
| 792 | נדבר | 15215 | 11122 | 0.6484347421 |
| 793 | שמונה | 18191 | 11034 | 0.6483853238 |
| 794 | הלאה | 16442 | 11215 | 0.6479885944 |
| 795 | הצעה | 20113 | 11235 | 0.6478501144 |
| 796 | תסתכל | 17573 | 11244 | 0.6476732706 |
| 797 | ראייה | 24456 | 11606 | 0.6471709215 |
| 798 | הדע | 14982 | 11125 | 0.6471310956 |

| | LEMMA | FREQUENCY | RANGE | UDP |
|-----|-------|-----------|-------|-----|
| 799 | מדוע | 37332 | 12801 | 0.6470783704 |
| 800 | הידי | 16651 | 11046 | 0.6468596261 |
| 801 | עשית | 14287 | 11109 | 0.6466812122 |
| 802 | מחשבה | 15714 | 11067 | 0.6466228783 |
| 803 | אקדח | 39948 | 12963 | 0.6465681297 |
| 804 | מאושר | 18661 | 11038 | 0.6463500223 |
| 805 | סביב | 15711 | 10810 | 0.6460287958 |
| 806 | גישה | 17476 | 11205 | 0.645804639 |
| 807 | מהירות | 18824 | 10975 | 0.6454430352 |
| 808 | פנה | 15694 | 10759 | 0.6450294635 |
| 809 | שוטר | 56235 | 14924 | 0.6448564384 |
| 810 | מידה | 15478 | 10674 | 0.6443406095 |
| 811 | מיליון | 26116 | 10942 | 0.644080189 |
| 812 | נחש | 16005 | 11043 | 0.644000921 |
| 813 | קר | 17343 | 10900 | 0.6439416391 |
| 814 | מרכז | 19237 | 10929 | 0.6436605734 |
| 815 | נקי | 16549 | 10838 | 0.6435246041 |
| 816 | קבע | 14837 | 10636 | 0.6434463895 |
| 817 | זהיר | 16261 | 10768 | 0.6431411339 |
| 818 | העלה | 14300 | 10511 | 0.6431050256 |
| 819 | הסתדר | 14436 | 10726 | 0.6430832736 |
| 820 | ארבעה | 17056 | 10579 | 0.642847399 |
| 821 | שער | 23225 | 11155 | 0.6427641541 |
| 822 | ראוי | 16150 | 10716 | 0.6426654687 |
| 823 | נת | 14024 | 10760 | 0.6425963739 |
| 824 | אדמה | 22675 | 11340 | 0.6425716519 |
| 825 | מוכר | 14943 | 10646 | 0.6424313062 |
| 826 | ם | 18094 | 10857 | 0.6423498462 |
| 827 | תינוק | 49123 | 14886 | 0.6421549492 |
| 828 | העמיד | 14541 | 10781 | 0.6416158382 |
| 829 | אגב | 14601 | 10783 | 0.6415822338 |
| 830 | רצון | 17003 | 10723 | 0.6415436783 |

|     | LEMMA | FREQUENCY | RANGE | UDP |
| --- | --- | --- | --- | --- |
| 831 | דרש | 15086 | 10705 | 0.6413020835 |
| 832 | סגור | 14788 | 10525 | 0.6412690057 |
| 833 | הכה | 16928 | 10523 | 0.6411804032 |
| 834 | יוצא | 13603 | 10314 | 0.6411760238 |
| 835 | בניין | 24245 | 11215 | 0.6411259746 |
| 836 | בקושי | 13857 | 10689 | 0.6409341156 |
| 837 | עשי | 14841 | 10993 | 0.6409186484 |
| 838 | כניסה | 15394 | 10563 | 0.640841239 |
| 839 | ביקר | 15213 | 10469 | 0.6405967235 |
| 840 | ירה | 21259 | 10718 | 0.6405054387 |
| 841 | כוכב | 33861 | 12372 | 0.64040029 |
| 842 | דעתך | 14270 | 10610 | 0.6403974806 |
| 843 | רגש | 17663 | 10886 | 0.6402669368 |
| 844 | רעב | 17084 | 10755 | 0.6401895042 |
| 845 | איום | 16945 | 10677 | 0.6401645514 |
| 846 | אירוע | 17203 | 10547 | 0.6401587224 |
| 847 | השנה | 17729 | 10483 | 0.6401000299 |
| 848 | נשבע | 15614 | 10431 | 0.640096061 |
| 849 | אינו | 97424 | 22002 | 0.6400922382 |
| 850 | כעס | 16194 | 10766 | 0.6399348394 |
| 851 | הושלם | 14931 | 10624 | 0.6397692214 |
| 852 | המשך | 14399 | 10453 | 0.6397504879 |
| 853 | יגע | 15599 | 10380 | 0.6396627983 |
| 854 | התעורר | 15781 | 10636 | 0.6395637873 |
| 855 | בפני | 15101 | 10488 | 0.6395623345 |
| 856 | תקווה | 16258 | 10662 | 0.6394894516 |
| 857 | אדוני | 119659 | 23260 | 0.6394781788 |
| 858 | ניו | 31706 | 11132 | 0.6394041344 |
| 859 | רואה | 14931 | 10948 | 0.6392676481 |
| 860 | ים | 26783 | 11409 | 0.6392602266 |
| 861 | היטב | 15348 | 10445 | 0.63907521 |
| 862 | טיפול | 21154 | 10912 | 0.6388078307 |

| | LEMMA | FREQUENCY | RANGE | UDP |
|---|---|---|---|---|
| 863 | מפתח | 20255 | 10730 | 0.6386307043 |
| 864 | אכן | 18504 | 10630 | 0.6382189052 |
| 865 | קלט | 17929 | 10591 | 0.6381770053 |
| 866 | תראו | 16401 | 10783 | 0.6381297552 |
| 867 | מחלקה | 21061 | 10556 | 0.6381147361 |
| 868 | כמוני | 13449 | 10338 | 0.6381142233 |
| 869 | דוד | 34311 | 12007 | 0.6380696722 |
| 870 | מעשה | 15444 | 10331 | 0.6375348014 |
| 871 | אסף | 14132 | 10389 | 0.6375294472 |
| 872 | מאית | 14315 | 10630 | 0.637363769 |
| 873 | שטות | 15059 | 9762 | 0.6370325165 |
| 874 | גאה | 14788 | 10180 | 0.6368535345 |
| 875 | תת | 14572 | 10301 | 0.6368165865 |
| 876 | אלף | 22105 | 10277 | 0.6362856564 |
| 877 | הודיע | 13848 | 10053 | 0.6362325456 |
| 878 | ידיד | 18120 | 10365 | 0.636000479 |
| 879 | היסטוריה | 16055 | 10043 | 0.6359981157 |
| 880 | צחק | 15117 | 10313 | 0.6359893169 |
| 881 | מאחור | 14145 | 10141 | 0.6359808507 |
| 882 | האם | 252114 | 31767 | 0.6355576557 |
| 883 | משימה | 26215 | 11277 | 0.6353326807 |
| 884 | אורח | 15323 | 10094 | 0.6350002439 |
| 885 | וידא | 13419 | 10373 | 0.6347930762 |
| 886 | חרא | 35830 | 11147 | 0.6347731563 |
| 887 | תיקן | 15845 | 10546 | 0.6345010271 |
| 888 | ע | 17728 | 10061 | 0.6343152887 |
| 889 | תעש | 14043 | 10224 | 0.6341616893 |
| 890 | בירה | 18982 | 10361 | 0.6336301933 |
| 891 | בחינה | 14402 | 9702 | 0.6333294761 |
| 892 | מס | 27950 | 10714 | 0.6331734526 |
| 893 | התגעגע | 16436 | 10209 | 0.6331526696 |
| 894 | שלישי | 14900 | 9673 | 0.6331134505 |

|     | LEMMA | FREQUENCY | RANGE | UDP |
| --- | --- | --- | --- | --- |
| 895 | ניקה | 14371 | 10171 | 0.6330208976 |
| 896 | מקומי | 15015 | 9921 | 0.6330172436 |
| 897 | תחושה | 14659 | 9992 | 0.6329762137 |
| 898 | תהה | 13079 | 9984 | 0.6329708853 |
| 899 | התנהג | 13402 | 10000 | 0.6329102726 |
| 900 | ר | 18144 | 10008 | 0.6328678992 |
| 901 | קרב | 18170 | 10203 | 0.6327703065 |
| 902 | תאונה | 20467 | 10322 | 0.6324579178 |
| 903 | מילא | 13112 | 9922 | 0.6323616322 |
| 904 | נתחיל | 13094 | 9890 | 0.6320485045 |
| 905 | הפה | 15224 | 9761 | 0.6319536943 |
| 906 | עצור | 18357 | 9969 | 0.6318472859 |
| 907 | הציג | 13058 | 9316 | 0.6310836173 |
| 908 | הריח | 15936 | 10142 | 0.6310753406 |
| 909 | טובה | 12856 | 9796 | 0.630920668 |
| 910 | השג | 14189 | 9964 | 0.6307658955 |
| 911 | ברירה | 14294 | 10242 | 0.6305407031 |
| 912 | נו | 17461 | 10079 | 0.6303642267 |
| 913 | אתקשר | 13871 | 9724 | 0.630079445 |
| 914 | חוסר | 13243 | 9667 | 0.6300372872 |
| 915 | השקר | 14766 | 9914 | 0.6299231011 |
| 916 | מישהי | 15678 | 10164 | 0.6296962365 |
| 917 | שר | 27584 | 10350 | 0.6294577543 |
| 918 | הבנה | 12504 | 9531 | 0.6294190374 |
| 919 | חטף | 16497 | 9866 | 0.629398524 |
| 920 | משמעות | 14356 | 9692 | 0.6293963852 |
| 921 | טען | 14950 | 9634 | 0.6293123723 |
| 922 | שדה | 18519 | 9777 | 0.6292853814 |
| 923 | סם | 32318 | 10829 | 0.6292684006 |
| 924 | אבי | 25210 | 10450 | 0.6290176556 |
| 925 | צפון | 17658 | 9785 | 0.6290015587 |
| 926 | הפסיד | 15106 | 9441 | 0.6289462904 |

| | LEMMA | FREQUENCY | RANGE | UDP |
|---|---|---|---|---|
| 927 | עבורך | 16207 | 9891 | 0.6289375764 |
| 928 | התחל | 11877 | 9398 | 0.6288659047 |
| 929 | ייתכן | 19498 | 10136 | 0.6287532454 |
| 930 | מפני | 22710 | 10575 | 0.6286980446 |
| 931 | אלייך | 15140 | 9960 | 0.6286547419 |
| 932 | גבול | 15935 | 9761 | 0.6285210688 |
| 933 | הרים | 12802 | 9306 | 0.6283404914 |
| 934 | תנועה | 14982 | 9386 | 0.6282947313 |
| 935 | חקירה | 24262 | 10087 | 0.628291565 |
| 936 | ראשי | 14696 | 9411 | 0.6281780833 |
| 937 | סיפק | 13210 | 9530 | 0.6281450779 |
| 938 | אשר | 22578 | 9766 | 0.6280505889 |
| 939 | מקור | 15699 | 9702 | 0.6280228726 |
| 940 | מלון | 24568 | 9865 | 0.6276081781 |
| 941 | יורק | 27791 | 9829 | 0.6273454869 |
| 942 | דומה | 12924 | 9468 | 0.6271437481 |
| 943 | רשם | 13080 | 9321 | 0.626766924 |
| 944 | התחה | 15025 | 9242 | 0.6267395572 |
| 945 | ודאי | 20206 | 9844 | 0.6266106799 |
| 946 | כביש | 18887 | 9677 | 0.6265302389 |
| 947 | מנה | 15573 | 9390 | 0.6264745692 |
| 948 | סיכון | 14296 | 9540 | 0.6263396616 |
| 949 | עשרה | 15509 | 9185 | 0.6262430927 |
| 950 | לימד | 13514 | 9280 | 0.626033317 |
| 951 | הכול | 52201 | 11674 | 0.6258631736 |
| 952 | חש | 14861 | 9585 | 0.6258606337 |
| 953 | בחירה | 13755 | 9557 | 0.6256709006 |
| 954 | לפה | 19366 | 9858 | 0.6252465957 |
| 955 | עצוב | 14366 | 9531 | 0.625212899 |
| 956 | התנצל | 13273 | 9438 | 0.6249713179 |
| 957 | הסתיר | 13497 | 9576 | 0.6249458234 |
| 958 | לאט | 16978 | 9232 | 0.6247997152 |

| | LEMMA | FREQUENCY | RANGE | UDP |
|---|---|---|---|---|
| 959 | נהרג | 15526 | 9546 | 0.6247885972 |
| 960 | שיקר | 14289 | 9589 | 0.6246517235 |
| 961 | התייחס | 12043 | 9036 | 0.6246440976 |
| 962 | מכירה | 16418 | 9343 | 0.6245483334 |
| 963 | דיווח | 14769 | 9390 | 0.6244039585 |
| 964 | הו | 108693 | 17359 | 0.624233775 |
| 965 | טלוויזיה | 19027 | 9566 | 0.6240775861 |
| 966 | ריצה | 14132 | 9286 | 0.6240761445 |
| 967 | דפק | 14239 | 9071 | 0.6238463174 |
| 968 | נולד | 13391 | 9172 | 0.6237391282 |
| 969 | לשעבר | 15093 | 9252 | 0.6236662063 |
| 970 | אמריקני | 21671 | 9337 | 0.6233666638 |
| 971 | רחב | 12844 | 9012 | 0.6231106402 |
| 972 | תרופה | 24470 | 10022 | 0.6229912611 |
| 973 | מאחר | 12768 | 9169 | 0.6228875959 |
| 974 | זר | 13710 | 9107 | 0.6227707514 |
| 975 | התחתן | 19464 | 9118 | 0.6225162241 |
| 976 | פרץ | 14472 | 9273 | 0.6224521242 |
| 977 | קלות | 11632 | 8976 | 0.6223452963 |
| 978 | חמישה | 13755 | 8933 | 0.6223190967 |
| 979 | שישה | 13188 | 8819 | 0.6221384595 |
| 980 | שת | 11984 | 9151 | 0.6220896285 |
| 981 | גוש | 12111 | 8928 | 0.6215892099 |
| 982 | קפץ | 12880 | 9033 | 0.6215815402 |
| 983 | הרגשה | 12244 | 9093 | 0.6214685059 |
| 984 | משוגע | 12826 | 8867 | 0.6214370418 |
| 985 | זבל | 17386 | 9105 | 0.6213662061 |
| 986 | לקוח | 20913 | 9252 | 0.6210193268 |
| 987 | קרע | 13143 | 9067 | 0.6207150099 |
| 988 | ול | 11528 | 8693 | 0.6206236099 |
| 989 | חוקי | 13720 | 8773 | 0.6206086392 |
| 990 | נמשך | 11977 | 8918 | 0.6201722638 |

|  | LEMMA | FREQUENCY | RANGE | UDP |
|------|-------|-----------|-------|-----|
| 991 | החבב | 15419 | 9143 | 0.6200478095 |
| 992 | רשמי | 12053 | 8877 | 0.620044703 |
| 993 | גודל | 12512 | 8786 | 0.6199315469 |
| 994 | חן | 15280 | 8825 | 0.6197141211 |
| 995 | משקה | 13544 | 8853 | 0.6195560105 |
| 996 | חופש | 13379 | 8735 | 0.6194960178 |
| 997 | אצבע | 14209 | 8982 | 0.6193206225 |
| 998 | שורה | 12503 | 8574 | 0.6193160177 |
| 999 | הרוויח | 12824 | 8729 | 0.6192660882 |
| 1000 | לם | 12282 | 8955 | 0.61925743 |

# Appendix 2: Scripts

## APPENDIX 2.1: HEBREWLEMMACOUNT.PY

```python
#! /usr/bin/env python3
# -*- coding: utf-8 -*-

import re
import os
import gzip
from collections import defaultdict


############################################################
# ----------------- INITIALIZE VARIABLES ----------------- #
############################################################

# Define path for topmost directory to search. Make sure this points
    to
# the correct location of your corpus.
corpus_path = './OpenSubtitles2018_parsed_single'

# Initialize dictionaries
lemma_by_corpus_dict = {}
lemma_totals_dict = {}
token_count_dict = {}
lemma_DPs_dict = defaultdict(float)
lemma_UDPs_dict = defaultdict(float)

total_tokens_int = 0
table_list = []

```

```python
# Set size of final list
list_size_int = 5000



#############################################################
# ------------------ DEFINE FUNCTIONS ------------------ #
#############################################################



# Open XML file and read it.
def open_and_read(file_loc):
    with gzip.open(file_loc, 'rt', encoding='utf-8') as f:
        read_data = f.read()
    return read_data



# Search for lemma and add counts to "frequency{}".
def find_and_count(doc):
    corpus = str(f)[38:-4]
    match_pattern = re.findall(r'lemma="[א-ת]+"', doc)
    for word in match_pattern:
        if word[7:-1] in lemma_by_corpus_dict:
            count = lemma_by_corpus_dict[word[7:-1]].get(corpus, 0)
            lemma_by_corpus_dict[word[7:-1]][corpus] = count + 1
        else:
            lemma_by_corpus_dict[word[7:-1]] = {}
            lemma_by_corpus_dict[word[7:-1]][corpus] = 1



#############################################################
# -------------------- OPEN AND READ -------------------- #
#############################################################
```

83

```
61   # Open and read all files. If calculating only for a specific
     ↪  language,
62   # comment out this code and uncomment the large block that follows.
63   #
64   for dirName, subdirList, fileList in os.walk(corpus_path):
65       if len(fileList) > 0:
66           f = dirName + '/' + fileList[0]
67           find_and_count(open_and_read(f))
68
69   ############################################################
70   # ---------------- LANGUAGE-SPECIFIC BLOCK ----------------
71   #
72   # This large block of code is for creating a list using only movies
     ↪  #
73   # with a specific primary language (in this case, Hebrew). Be sure
     ↪  to #
74   # uncomment the relevant lines of code, and to comment out the block
     ↪  #
75   # above. #
76   #
77   #
78   # Create list of IDs for movies with Hebrew as primary language. #
79   # This makes use of a text file that must already exist with this
     ↪  list. #
80   #
81   # Hebrew_IDs_list = []
82   # with open('./Hebrew_originals.txt', 'r', encoding='utf-8') as f:
83   #     read_data = f.read()
84   #     Hebrew_IDs_list = re.findall(r'\s\stt[0-9]+\t', read_data)
85   # Hebrew_IDs_list = [line[4:-1] for line in Hebrew_IDs_list]
86   #
87   #
88   # Delete extra 0s at the beginning of Hebrew movie IDs. #
```

84

```python
89  #
90  # for item in Hebrew_IDs_list:
91  #     if item[0] == '0':
92  #         Hebrew_IDs_list[Hebrew_IDs_list.index(item)] = item[1:]
93  # for item in Hebrew_IDs_list:
94  #     if item[0] == '0':
95  #         Hebrew_IDs_list[Hebrew_IDs_list.index(item)] = item[1:]
96  #
97  #
98  # Open and read files for movies with Hebrew as the primary
    ↪  language. #
99  #
100 # for dirName, subdirList, fileList in os.walk(corpus_path):
101 #     if len(fileList) > 0:
102 #         f = dirName + '/' + fileList[0]
103 #         folders = re.split('/', dirName)
104 #         if folders[len(folders)-1] in Hebrew_IDs_list:
105 #             find_and_count(open_and_read(f))
106 #
107 # ------------- END OF LANGUAGE-SPECIFIC BLOCK -------------
108 ############################################################
109
110
111 ############################################################
112 # -------------------- CALCULATIONS -------------------- #
113 ############################################################
114
115 # Calculate token count per corpus
116 for lemma in lemma_by_corpus_dict:
117     for corpus in lemma_by_corpus_dict[lemma]:
118         token_count_dict[corpus] = token_count_dict.get(
119             corpus, 0) + lemma_by_corpus_dict[lemma][corpus]
120
```

```python
121  # Calculate total frequencies per lemma
122  for lemma in lemma_by_corpus_dict:
123      lemma_totals_dict[lemma] =
     ↪  sum(lemma_by_corpus_dict[lemma].values())
124
125  # Calculate total token count
126  for corpus in token_count_dict:
127      total_tokens_int = total_tokens_int +
     ↪  token_count_dict.get(corpus, 0)
128
129  # Calculate DPs
130  for lemma in lemma_by_corpus_dict.keys():
131      for corpus in lemma_by_corpus_dict[lemma].keys():
132          lemma_DPs_dict[lemma] = lemma_DPs_dict[lemma] + abs(
133              (token_count_dict[corpus] /
134               total_tokens_int) -
135              (lemma_by_corpus_dict[lemma][corpus] /
136               lemma_totals_dict[lemma]))
137  lemma_DPs_dict = {lemma: DP/2 for (lemma, DP) in
     ↪  lemma_DPs_dict.items()}
138
139  # Calculate UDPs
140  lemma_UDPs_dict = {lemma: 1-DP for (lemma, DP) in
     ↪  lemma_DPs_dict.items()}
141
142
143  ###########################################################
144  # -------------- SORT LIST AND CREATE TABLE -------------- #
145  ###########################################################
146
147  # Sort entries by UDP
148  UDP_sorted_list = [(k, lemma_UDPs_dict[k]) for k in sorted(
149      lemma_UDPs_dict, key=lemma_UDPs_dict.__getitem__,
```

```
150    reverse=True)]

151

152    # Create list of tuples with all values (Lemma, Frequency, Range,
       ↪  UDP)
153    for k, v in UDP_sorted_list[:list_size_int]:
154        table_list.append((k, lemma_totals_dict[k], sum(
155            1 for count in lemma_by_corpus_dict[k].values() if count >
               ↪  0),
156            v))

157

158    ############################################################
159    # ---------------- SORT-BY-FREQUENCY BLOCK -----------------
160    #
161    # Sort entries by raw frequency (total lemma count). To sort the
       ↪  final #
162    # list by frequency instead of UDP, comment out the above code
       ↪  within the #
163    # "SORT LIST AND CREATE TABLE" section, and also uncomment the
       ↪  relevant #
164    # lines of code in this block. #
165    #
166    #
167    # Sort entries by raw frequency #
168    #
169    # frequency_sorted_list = [(k, lemma_totals_dict[k]) for k in
       ↪  sorted(
170    #     lemma_totals_dict, key=lemma_totals_dict.__getitem__,
171    #     reverse=True)]
172    #
173    #
174    # Create list of tuples with all values (Lemma, Frequency, Range,
       ↪  UDP) #
175    #
```

87

```
176  # for k, v in frequency_sorted_list[:list_size_int]:
177  #     table_list.append((k, v, sum(
178  #         1 for count in lemma_by_corpus_dict[k].values() if count >
     ↪ 0),
179  #         lemma_UDPs_dict[k]))
180  #
181  # ------------ END OF SORT-BY-FREQUENCY BLOCK ------------
182  ##########################################################
183
184  # Calculate list size for 80% coverage and set that as the list
     ↪ size. Note
185  # that if the initial list_size_int (set near the beginning of the
     ↪ script)
186  # provides less than the desired coverage, it will default to that
     ↪ instead.
187  #
188  # added_freq_int = 0
189  # count = 0
190  # for k, v in UDP_sorted_list:
191  #     if added_freq_int / total_tokens_int < 0.8:
192  #         added_freq_int = added_freq_int + lemma_totals_dict[k]
193  #         count = count + 1
194  #     else:
195  #         break
196  # list_size_int = count
197
198  # Write final tallies to CSV file
199  result = open('./export/HebrewWordList2.csv', 'w')
200  result.write('LEMMA, FREQUENCY, RANGE, UDP\n')
201  for i in range(list_size_int):
202      result.write(str(table_list[i][0]) + ', ' +
203                   str(table_list[i][1]) + ', ' +
204                   str(table_list[i][2]) + ', ' +
```

```
205                    str(table_list[i][3]) + '\n')
206  result.close()
207
208  # Print final tallies. Uncomment this code to see the results
209  # printed instead of writing them to a file.
210  #
211  # for i in range(list_size_int):
212  #     print('Lemma: ' + table_list[i][0] +
213  #           '\tFrequency: ' + str(table_list[i][1]) +
214  #           '\tRange: ' + str(table_list[i][2]) +
215  #           '\tUDP: ' + str(table_list[i][3]))
```

## Appendix 2.2: OMDb-fetch.py

```python
#! /usr/bin/env python3
# -*- coding: utf-8 -*-

# import re
from sys import argv
import os
import glob
import omdb

# year = '1996'
script, year, id_start = argv

dirs = []
p = []


for name in glob.glob(
        '../OpenSubtitles2018_parsed/parsed/he/' + year + '/*/'):
    p.append(name)
# p = Path('../OpenSubtitles2018_parsed/parsed/he')
# p = list(p.glob('[198-199]*/*/*.xml'))

p = [os.path.basename(os.path.dirname(str(i))) for i in p]

for i in p:
    if i not in dirs:
        dirs.append(i)

for i in dirs:
    while len(i) < 7:
        dirs[dirs.index(i)] = '0' + i
```

```python
32          i = '0' + i

33

34 dirs.sort()

35

36 # for i in dirs:
37 #     print('tt' + i)

38

39 print('# ' + year + '\n' +
40       'IMDb ID\tTitle\tYear\tLanguage(s)')

41

42

43 omdb.set_default('apikey', '906517b3')

44

45 for i in dirs:
46     if id_start != '':
47         if i > id_start:
48             print('tt' + i + '\t', end="", flush=True)
49             doc = omdb.imdbid('tt' + i)
50             # if doc['language'] == 'Hebrew':
51             print(doc['title'] + '\t' +
52                   doc['year'] + '\t' +
53                   doc['language'])
54     else:
55         print('tt' + i + '\t', end="", flush=True)
56         doc = omdb.imdbid('tt' + i)
57         # if doc['language'] == 'Hebrew':
58         print(doc['title'] + '\t' +
59               doc['year'] + '\t' +
60               doc['language'])
```

# Appendix 2.3: single_file_extract.py

```python
#! /usr/bin/env python3
# -*- coding: utf-8 -*-

import shutil
import os

source = '../OpenSubtitles2018_parsed'
destination = './OpenSubtitles2018_parsed_single'

# Copy the directory tree into a new location
shutil.copytree(source, destination,
    ignore=shutil.ignore_patterns('*.*'))

# Copy the first file in each folder into the new tree
for dirName, subdirList, fileList in os.walk(source):
    for fname in fileList:
        if fname == '.DS_Store':
            fileList.remove(fname)
    if len(fileList) > 0:
        del fileList[1:]
        src = dirName + '/' + fileList[0]
        dst = destination + dirName[27:] + '/'
        shutil.copy2(src, dst)
```

Appendix 3: Movies used

# References

Albert, A., MacWhinney, B., Nir, B., & Wintner, S. (2013). The Hebrew CHILDES corpus: Transcription and morphological analysis. *Language Resources and Evaluation*, *47*(4), 973–1005. https://doi.org/10.1007/s10579-012-9214-z

Al-Surmi, M. (2012). Review: Quaglio (2009). Television dialogue: The sitcom Friends vs. Natural conversation. Philadelphia: John Benjamins. *Corpora*, *7*(1). https://doi.org/10.3366/corp.2012.0022

Amir, N., Silber-Varod, V., & Izre'el, S. (2004). Characteristics of Intonation Unit Boundaries in Spontaneous Spoken Hebrew and Acoustic Correlates. In B. Bernard & I. Marlien (Eds.), *Speech Prosody 2004, Nara, Japan, March 23-26, 2004: Proceedings* (pp. 677–680). Nara, Japan.

Anglin, J. M., Miller, G. A., & Wakefield, P. C. (1993). Vocabulary development: A morphological analysis. *Monographs of the Society for Research in Child Development*, *58*(10), i–186. https://doi.org/10.2307/1166112

Bauer, L., & Nation, P. (1993). Word families. *International Journal of Lexicography*, *6*(4), 253–279. https://doi.org/10.1093/ijl/6.4.253

Biber, D. (1988). *Variation across speech and writing*. Cambridge, UK: Cambridge University Press.

Biber, D. (1993). Representativeness in corpus design. *Literary and Linguistic Computing*, *8*(4), 243–257. https://doi.org/10.1093/llc/8.4.243

Biber, D. (1995). *Dimensions of Register Variation: A Cross-Linguistic Comparison*. Cambridge University Press.

Brezina, V., & Gablasova, D. (2015). Is there a core general vocabulary? Introducing the new general service list. *Applied Linguistics*, *36*(1), 1–22. https://doi.org/10.1093/applin/amt018

Brysbaert, M., & New, B. (2009). Moving beyond Kučera and Francis: A critical evaluation of current word frequency norms and the introduction of a new and improved word frequency measure for American English. *Behavior Research Methods*, *41*(4), 977–990. https://doi.org/10.3758/BRM.41.4.977

Carroll, J. B., Davies, P., & Richman, B. (1971). *The American heritage word frequency book*. Boston: Houghton Mifflin.

Chung, T. M., & Nation, P. (2003). Technical vocabulary in specialised texts. *Reading in a Foreign Language*, *15*(2), 103.

Chung, T. M., & Nation, P. (2004). Identifying technical vocabulary. *System*, *32*(2), 251–263.

https://doi.org/10.1016/j.system.2003.11.008

*Collins Cobuild English grammar.* (2005). Glasgow: HarperCollins.

Cowie, A. P. (2009). *The Oxford History of English Lexicography.* Oxford Univ.

Coxhead, A. (2000). A new academic word list. *TESOL Quarterly, 34*(2), 213–238. https://doi.org/10.2307/3587951

Coxhead, A. (2016). Reflecting on Coxhead (2000), "a new academic word list". *TESOL Quarterly, 50*(1), 181–185. https://doi.org/10.1002/tesq.287

Dang, T. N. Y., & Webb, S. (2016). Evaluating lists of high-frequency words. *ITL - International Journal of Applied Linguistics, 167*(2), 132–158. https://doi.org/10.1075/itl.167.2.02dan

Dekel, N. (2010). *A matter of time: Tense, mood and aspect in Spontaneous Spoken Israeli Hebrew.* Utrecht: LOT.

Droop, M., & Verhoeven, L. (2003). Language proficiency and reading ability in first- and second-language learners. *Reading Research Quarterly, 38*(1), 78–103.

Francis, W. N., Kučera, H., & Mackie, A. W. (1982). *Frequency analysis of English usage: Lexicon and grammar.* Boston: Houghton Mifflin.

Gardner, D. (2007). Validating the construct of word in applied corpus-based vocabulary research: A critical survey. *Applied Linguistics, 28*(2), 241–265. https://doi.org/10.1093/applin/amm010

Gilner, L. (2011). A primer on the general service list. *Reading in a Foreign Language, 23*(1), 65.

Goldberg, Y. (2011, November). *Automatic Syntactic Processing of Modern Hebrew* (PhD thesis). Ben-Gurion University, Beer-Sheva, Israel.

Goldberg, Y., & Elhadad, M. (2009). Hebrew dependency parsing: Initial results. In *Proceedings of the 11th International Conference on Parsing Technologies* (pp. 129–133). Paris: Association for Computational Linguistics.

Goldberg, Y., & Elhadad, M. (2010). Easy-First Dependency Parsing of Modern Hebrew. In *Proceedings of the NAACL HLT 2010 First Workshop on Statistical Parsing of Morphologically-Rich Languages* (pp. 103–107). Los Angeles, CA, USA: Association for Computational Linguistics.

Goldberg, Y., & Elhadad, M. (n.d.). Two Syntactic Parsers for Modern Hebrew and a large automatically parsed corpus.

Gretz, S., Itai, A., MacWhinney, B., Nir, B., & Wintner, S. (2015). Parsing Hebrew CHILDES transcripts. *Language Resources and Evaluation, 49*(1), 107–145. https://doi.org/10.1007/s10579-013-9256-x

Guthmann, N., Krymolowski, Y., Milea, A., & Winter, Y. (2008). Automatic Annotation of Morpho-Syntactic Dependencies in a Modern Hebrew Treebank. *LOT Occasional Series*, *12*, 77–90.

Hirsh, D., & Nation, P. (1992). What vocabulary size is needed to read unsimplified texts for pleasure? *Reading in a Foreign Language*, *8*, 689–696.

Hoek, J., Evers-Vermeul, J., & Sanders, T. (2015). The role of expectedness in the implicitation and explicitation of discourse relations.

Hu, M., & Nation, P. (2000). Unknown vocabulary density and reading comprehension. *Reading in a Foreign Language*, *13*(1).

Itai, A., & Segal, E. (2003). A Corpus Based Morphological Analyzer for Unvocalized Modern Hebrew.

Izre'el, S. (2004). Transcribing Spoken Israeli Hebrew: Preliminary Notes. In D. D. Ravid & H. B.-Z. Shyldkrot (Eds.), *Perspectives on Language and Language Development: Essays in Honor of Ruth A. Berman* (pp. 61–72). Kluwer: Dodrecht. https://doi.org/10.1007/1-4020-7911-7_6

Izre'el, S., Auran, C., Bertrand, R., Chanet, C., Colas, A., Di Cristo, A., … Vion, M. (2005). Intonation Units and the Structure of Spontaneous Spoken Language: A View from Hebrew. In *Proceedings of the IDP05 International Symposium on Discourse-Prosody Interfaces* (p. 20).

Izre'el, S., Hary, B., & Rahav, G. (2001). Designing CoSIH: The Corpus of Spoken Israeli Hebrew. *International Journal of Corpus Linguistics*, *6*(2), 171–197. https://doi.org/10.1075/ijcl.6.2.01izr

Jakubíček, M., Kilgarriff, A., Kovář, V., Rychlý, P., & Suchomel, V. (2013). The TenTen Corpus Family. In *7th International Corpus Linguistics Conference CL 2013* (pp. 125–127). Lancaster.

Jeon, E. H. (2011). Contribution of morphological awareness to second-language reading comprehension. *The Modern Language Journal*, *95*(2), 217–235. https://doi.org/10.1111/j.1540-4781.2011.01179.x

Kilgarriff, A., & Grefenstette, G. (2003). Introduction to the special issue on the web as corpus. *Computational Linguistics*, *29*(3), 333–347.

Kobeleva, P. P. (2012). Second language listening and unfamiliar proper names: Comprehension barrier? *RELC Journal*, *43*(1), 83–98. https://doi.org/10.1177/0033688212440637

Kučera, H., & Francis, W. N. (1967). *Computational analysis of present-day American English*. Brown University Press.

Lison, P., & Tiedemann, J. (2016). OpenSubtitles2016: Extracting large parallel corpora from movie and TV subtitles. *Proceedings of the 10th International Conference on Language Resources and Evaluation (LREC 2016)*, 7.

Matsushita, T. (2012). In What Order Should Learners Learn Japanese Vocabulary? A Corpus-based Approach.

Mettouchi, A., Lacheret-Dujour, A., Silber-Varod, V., & Izre'el, S. (2007). Only Prosody? Perception of speech segmentation in Kabyle and Hebrew. In *Intefaces discours prosodie : Actes du 2ème Symposium international & Colloque Charles Bally* (pp. 207–218).

Milton, J. (2009). *Measuring second language vocabulary acquisition*. Bristol, UK ; Buffalo N.Y.: Multilingual Matters.

Mochizuki, M., & Aizawa, K. (2000). An affix acquisition order for EFL learners: An exploratory study. *System*, *28*, 291–304. https://doi.org/10.1016/S0346-251X(00)00013-0

Nagy, W., Anderson, R. C., Schommer, M., Scott, J. A., & Stallman, A. C. (1989). Morphological families in the internal lexicon. *Reading Research Quarterly*, *24*(3), 262–282. https://doi.org/10.2307/747770

Nagy, W. E., Diakidoy, I.-A. N., & Anderson, R. C. (1991). The development of knowledge of derivational suffixes. *Center for the Study of Reading Technical Report; No. 536.*

Nation, I. (1982). Beginning to learn foreign vocabulary: A review of the research. *RELC Journal*, *13*(1), 14–36. https://doi.org/10.1177/003368828201300102

Nation, I. (2006). How large a vocabulary is needed for reading and listening? *Canadian Modern Language Review*, *63*(1), 59–82.

Nation, I. (2016). *Making and using word lists for language learning and testing*. Amsterdam: John Benjamins Publishing Company. https://doi.org/10.1075/z.208

Nation, I. S. P. (1990). *Teaching & learning vocabulary* (1 edition). Boston, Mass: Heinle ELT.

Nation, I. S. P. (2013). *Learning vocabulary in another language* (2nd ed.). Cambridge: Cambridge University Press.

Nation, I. S. P., & Webb, S. (2010). *Researching and analyzing vocabulary* (1 edition). Boston, MA: Heinle ELT.

Nation, P. (2000). Learning vocabulary in lexical sets: Dangers and guidelines. *TESOL Journal*, *9*(2), 6–10. https://doi.org/10.1002/j.1949-3533.2000.tb00239.x

Nation, P. (2004). A study of the most frequent word families in the British National Corpus. In P. Bogaards & B. Laufer (Eds.), *Language Learning & Language Teaching* (Vol. 10, pp. 3–13). Amsterdam: John Benjamins Publishing Company. https://doi.org/10.1075/lllt.10.03nat

Nation, P., & Coxhead, A. (2014). Vocabulary size research at Victoria University of Wellington, New Zealand. *Language Teaching*, *47*(03), 398–403. https://doi.org/10.1017/S0261444814000111

Nation, P., & Kyongho, H. (1995). Where would general service vocabulary stop and special purposes vocabulary begin? *System, 23*(1), 35–41. https://doi.org/10.1016/0346-251X(94)00050-G

New, B., Brysbaert, M., Veronis, J., & Pallier, C. (2007). The use of film subtitles to estimate word frequencies. *Applied Psycholinguistics, 28*(4), 661–677. https://doi.org/10.1017/S014271640707035X

Quaglio, P. (2009). *Television dialogue: The sitcom Friends vs. Natural conversation.* John Benjamins Publishing.

Read, J. (1988). Measuring the vocabulary knowledge of second langauge learners. *RELC Journal, 19*(2), 12–25. https://doi.org/10.1177/003368828801900202

Restrepo Ramos, F. D. (2015). Incidental vocabulary learning in second language acquisition: A literature review. *PROFILE Issues in Teachers' Professional Development, 17*(1), 157–166. https://doi.org/10.15446/profile.v17n1.43957

Schmitt, N. (2010). *Researching vocabulary: A vocabulary research manual.* Houndmills, Basingstoke, Hampshire; New York: Palgrave Macmillan.

Schmitt, N., Jiang, X., & Grabe, W. (2011). The percentage of words known in a text and reading comprehension. *The Modern Language Journal, 95*(1), 26–43. https://doi.org/10.2307/41262309?ref=no-x-route:cb78a69b6dc8bf1478b58d47243b1248

Schmitt, N., & Meara, P. (1997). Researching vocabulary through a word knowledge framework. *Studies in Second Language Acquisition, 19*(1), 17–36.

Schmitt, N., & Schmitt, D. (2014). A reassessment of frequency and vocabulary size in L2 vocabulary teaching. *Language Teaching, 47*(04), 484–503. https://doi.org/10.1017/S0261444812000018

Schmitt, N., & Zimmerman, C. B. (2002). Derivative Word Forms: What Do Learners Know? *TESOL Quarterly, 36*(2), 145–171. https://doi.org/10.2307/3588328

Sima'an, K., Itai, A., Winter, Y., Altman, A., & Nativ, N. (2001). Building a tree-bank of modern Hebrew text. *Traitement Automatique Des Langues, 42*(2), 247–380.

Sorell, C. J. (2012). Zipf's law and vocabulary. In C. Chapelle (Ed.), *The Encyclopedia of Applied Linguistics.* Hoboken, NJ, USA: John Wiley & Sons, Inc.

Sorell, C. J. (2013). *A study of issues and techniques for creating core vocabulary lists for English as an international language* (Unpublished Dissertation). Victoria University of Wellington, Wellington, New Zealand.

The history of Collins COBUILD. (n.d.). https://www.collinsdictionary.com/cobuild/.

Tiedemann, J. (2012). Parallel Data, Tools and Interfaces in OPUS, 5.

Tiedemann, J. (2016). Finding Alternative Translations in a Large Corpus of Movie Subtitles, 5.

Tyler, A., & Nagy, W. (1989). The acquisition of English derivational morphology. *Journal of Memory and Language*, *28*(6), 649–667. https://doi.org/10.1016/0749-596X(89)90002-8

Tyler, A., & Nagy, W. (1990). Use of derivational morphology during reading. *Cognition*, *36*(1), 17–34. https://doi.org/10.1016/0010-0277(90)90052-L

van Zeeland, H., & Schmitt, N. (2013). Lexical coverage in L1 and L2 listening comprehension: The same or different from reading comprehension? *Applied Linguistics*, *34*(4), 457–479. https://doi.org/10.1093/applin/ams074

Wang, M., Cheng, C., & Chen, S.-W. (2006). Contribution of morphological awareness to Chinese-English biliteracy acquisition. *Journal of Educational Psychology*, *98*(3), 542–553. https://doi.org/10.1037/0022-0663.98.3.542

Ward, J., & Chuenjundaeng, J. (2009). Suffix knowledge: Acquisition and applications. *System*, *37*(3), 461–469. https://doi.org/10.1016/j.system.2009.01.004

Webb, S. A., & Chang, A. C.-S. (2012). Second language vocabulary growth. *RELC Journal*, *43*(1), 113–126.

Yael, M. (2014). The Haifa Corpus of Spoken Hebrew. http://weblx2.haifa.ac.il/~corpus/corpus_website/.

Zipf, G. K. (1935). *The psycho-biology of language.* Cambridge, Mass.: M.I.T. Press.

Zipf, G. K. (1949). *Human behavior and the principle of least effort: An introduction to human ecology.* New York: Hafner.