

MASTER'S THESIS

IN COMPUTATIONAL LINGUISTICS

Deconstructing Constructed Languages

Author:

Connor KIRBERGER

Supervisors:

Çağrı ÇÖLTEKİN

Christian BENTZ

SEMINAR FÜR SPRACHWISSENSCHAFT
EBERHARD-KARLS-UNIVERSITÄT TÜBINGEN

December 2023

Hiermit versichere ich, dass ich die Arbeit selbständig verfasst, keine anderen als die angegebenen Hilfsmittel und Quellen benutzt, alle wörtlich oder sinngemäß aus anderen Werken übernommenen Aussagen als solche gekennzeichnet habe und dass die Arbeit weder vollständig noch in wesentlichen Teilen Gegenstand eines anderen Prüfungsverfahrens gewesen ist und dass die Arbeit weder vollständig noch in wesentlichen Teilen bereits veröffentlicht wurde sowie dass das in Dateiform eingereichte Exemplar mit den eingereichten gebundenen Exemplaren übereinstimmt.

I hereby declare that this paper is the result of my own independent scholarly work. I have acknowledged all the other authors' ideas and referenced direct quotations from their work (in the form of books, articles, essays, dissertations, and on the internet). No material other than that listed has been used.

Tübingen, September 6, 2024

Firstname Surname

Contents

List of Figures	iv
List of Tables	iv
List of Abbreviations	iv
1 Introduction & Motivation	1
1.1 Scope of Study & Research Question	2
2 Background	3
2.1 History of Constructed Languages	3
2.2 Prior Studies	5
3 Methodology	7
3.1 Data	7
3.1.1 Wikimedia	8
3.2 Data Preprocessing	8
3.3 Libraries and APIs	9
3.4 Feature Extraction	9
3.4.1 TTR & MATTR	9
3.4.2 Morphological Complexity	9
3.4.3 Zipfian Distribution	10
3.4.4 Entropy	10
3.4.5 PCA	10
3.5 Classification	10
3.5.1 Decision Tree	10
3.5.2 Random Forest	10
3.5.3 One-Class SVM	10
3.6 Evaluation of Classifiers	10
4 Results	11
4.1 Results of TTR & MATTR	11
4.2 Results of Morphological Segmentation	11
4.3 Results of PCA	11
4.4 Results of One-Class SVM	11
4.5 Results of Decision Tree	11
4.6 Results of Random Forest	11
5 Discussion	12
6 Conclusion	13
6.1 Future Work	13
7 Acknowledgments	14
8 Appendices	17

Abstract

Write the abstract here.

List of Figures

2.1	Wilson's expression of "dog" in his philosophical language (Goodall 2022)	4
2.2	A taxonomy of constructed languages (Gobbo 2016)	5
4.1	Principal Component Analysis on TTR, MATTR, Word and Char Distribution Entropy	11
4.2	Decision Tree Classifier	12

List of Tables

3.1	Constructed languages and their respective language family influences used in the study.	7
8.1	Lengths of each text after pre-processing.	17

List of Abbreviations

API	Application Programming Interface
NLP	Natural Language Processing
PCA	Principle Component Analysis
TF-IDF	Term Frequency - Inverse Document Frequency
RNN	Recurrent Neural Network
LSTM	Long Short-Term Memory
SVM	Support Vector Machine
XML	eXtensible Markup Language
TTR	Type-Token Ratio
MATTR	Moving-Average Type-Token Ratio
IAL	International Auxiliary Language
SVO	Subject-Verb-Object
SOV	Subject-Object-Verb

1 Introduction & Motivation

Constructed languages—also called artificial languages, invented languages, planned languages, engineering languages, glossopoeia, or more simply as "conlangs" (Ball 2015)—are languages that are consciously and purposefully created for some intended use, usually being defined in antithesis to the spontaneous and organic method in which natural languages arise and develop (Sanders 2016). These variations of the term are often, but not always, used interchangeably, as linguists do not all agree upon a core term due to personal preferences (Adelman 2014), and there are sometimes differences in nuance depending on the context in which they appear. This thesis will mainly refer to them as constructed languages for simplicity.

The intended uses for which they are created can range broadly. Some are made specifically for fictional media, often seen in the genres of fantasy or science-fiction, with some more well-known examples being J. R. R. Tolkien's Elvish languages (e.g., Quenya, Sindarin, Nandorin) found in the world of Middle-earth in his writings, Marc Okrand's Klingon language from the Star Trek universe, and David J. Peter's Dothraki language used in George R. R. Martin's *A Song of Ice and Fire* novels along with their television adaptation, *Game of Thrones* (Jeffrey Punske (editor) 2020). Others are created to function as international auxiliary languages (IALs)—languages planned for the use of international and cross-cultural communication (Gobbo 2016). The most well-known example (based on estimated number of speakers) of these is Esperanto, created in the 19th Century by L. L. Zamenhof. Typically, constructed languages are distinguished and categorized based on these communicative functions, although other taxonomies also exist, as we will discuss further in section 2.

Despite being defined in contrast to one another, however, constructed and natural languages are not necessarily opposite to one another characteristically. Aside from their origins, the boundaries between the two are not always clear when analyzed in greater detail (Goodall 2022). For example, Schubert (1989) argues that some languages which are considered "natural" have some degree of artificiality, such as standardized written German and English differing from their spoken forms, and that the reverse is also true of some languages which are considered "artificial" because they draw from aspects of natural languages. As such, he believes human languages exist on a continuum of the two labels, rather than in the binary distinction—a view echoed by other linguists as well (Novikov 2022).

In many ways, it can be argued that investigation into the disparity between the two is at its core a mere part of the larger debate surrounding what exactly constitutes a language. The search for and defining of language univer-

sals is not only fundamental to the field of linguistic research, but also a topic of widespread debate. At present, many theoretical and foundational contributions to this discussion exist, ranging from Greenberg's proposed universals (Greenberg 1970) to ideas about universal grammars, such as Chomsky's.

Furthermore, while research on constructed languages is far from being novel, it is less common in computational linguistics. Thus, my motivation here was to try something relatively new, by approaching such an investigation using machine learning methods.

1.1 Scope of Study & Research Question

The present work analyzes various linguistic features of both types of languages and seeks to make a comparison on the differences, if any exist, between them. More specifically, this study presents a binary classification task using decision tree and outlier detection models to determine if they can be distinguished from one another, based on the specific features examined and parallel Wikipedia data.

Because of the wide-ranging nature of conducting such a broad analysis, there will of course be many features left unconsidered or excluded, intentionally or otherwise. With this in mind and following the precedent set by other related research on this topic, the focus in this particular study is mainly on features such as the entropy, morphological complexity, and lexical diversity of each language, based on the selected corpora.

The following is a breakdown of the structure of this thesis from here onward: the next section provides relevant background information, including an overview on constructed languages and a comprehensive review of related literature that examines the prior theoretical groundwork laid for exploring linguistic similarities and differences between constructed and natural languages; section 3 covers in detail the methodology taken in this research, from an explanation of the data used to the various experiments performed; section 4 presents the results of the study and discussion of these follows in section 5; lastly, section 6 consists of a conclusion as well as elaboration for possible future work.

2 Background

The vast landscape of linguistic research comprises a myriad of literature delving into the intricacies of languages, both natural and constructed. As this paper is concerned with constructed languages in particular and possible distinctive properties they may have, this section begins with a brief overview of their history and development, which provides some relevant context. Following this is an overview of some related literature, after which I will discuss the various computational methods implemented in this study and provide some background information on how they work.

2.1 History of Constructed Languages

Okrent (2009) states, "The history of invented languages is, for the most part, a history of failure." She may be justified in saying this, depending on one's definition of failure in this context. From past to present, the total number of constructed languages may be as high as a thousand (Libert 2016; Schubert 1989), with hundreds proposed for the purpose of being IALs in Europe alone (Schubert et al. 2001). Yet of these, only Esperanto is commonly considered to be successful in achieving its creator's intended goal of world-wide use as an auxiliary language (or rather that it is by far the most successful), with very few others even coming close, having a conservative estimation of two million speakers (Okrent 2009).

While the construction of languages is possibly as old as human history, they typically were not written down and were limited to in-group communication (Gobbo 2016). The first documented endeavors came out of religious contexts and were likely used as secret languages, intentionally obscured and incomprehensible to lay people. In the 12th century, abbess Hildegard of Bingen described and recorded a lexicon for *Lingua Ignota*, a Latin name meaning "unknown language". While extensive documentation of it (i.e., a grammar) was never found, it possessed a semiotic system based on Latin, German, and Greek. Later in the 14th century, a group of Sufi mystics created *Balaibalan*, a language written in the Ottoman Turkish alphabet and which incorporated features of Persian, Turkish, and Arabic languages (Novikov 2022).

Interest in creating such languages picked up in the 17th century with the rise of so-called philosophical languages. In contrast to the last two, these languages were made to be more precise, less ambiguous, and better allow for philosophical reasoning (compared to natural language), such as by organizing world knowledge into hierarchies (Goodall 2022). Notable figures involved in making these include Francis Lodwick, Gottfried Leibniz, and John Wilkins, the latter of whose being arguably the most well-known and influ-

- (1) special > creature > distributively > substances > animate > species > sensitive > sanguineous > beasts > viviparous > clawed > rapacious > oblong-headed > European > terrestrial > big > docile

Figure 2.1: Wilson's expression of "dog" in his philosophical language (Goodall 2022)

ential. Wilkins created a system of semantic categorization, cataloging all concepts in the universe (Okrent 2009), and then published his proposed language (Wilkins 1968). An example of this hierarchal categorization can be seen in Figure 2.1.

In the 19th and 20th centuries the focus for language construction, especially in Europe, shifted to that of making international auxiliary languages (IALs) intended to better enable communication across language barriers, i.e., people who do not share a similar language (Goodall 2022). Notably, this means they were generally designed to resemble natural language, with choice exceptions being the simplification of certain linguistic features. The surge in need for IALs correlated with the increase in prevalence and accessibility regarding international travel and communication at the time. Such languages were also described as "neutral" (Large 1985), in the sense that individual advantages amongst speakers and learners would, theoretically, not exist due to IALs being second languages to everyone (Gobbo 2016). That being said, many of the most prominent examples (e.g., Volapük, Interlingua, Esperanto, Ido) are derived only from the Indo-European language family (Novikov 2022; Goodall 2022), so such a description might not be apt.

As the constructed languages examined and used in the dataset of the present work are all IALs, it would be beneficial to introduce them in more detail here. Volapük was made in 1879 by Catholic German priest Johann Martin Schleyer, who believed it was given to him by God. Argued to be the first successful constructed language due to amassing so many supporters (Gobbo 2016), it soon died out in favor of Esperanto, which Ludwik Lejzer Zamenhof published in 1887.

Finally,

While they all share the defining feature of having been purposefully created, their other features (e.g., phonetic, morphological, syntactical, lexical, orthographic) can vary immensely depending on factors such as, for example, their intended purpose for use or the other languages they draw from. An example of this was observed by Gobbo (2016) in secret languages, specifically their tendency to have more complicated features, such as morphological irregularities, "in order to preserve their secrecy." Contrast to this are IALs, which have the opposite tendency for the sake of ease of communication and second-language acquisition, reflected in commonly assigned features such as SVO

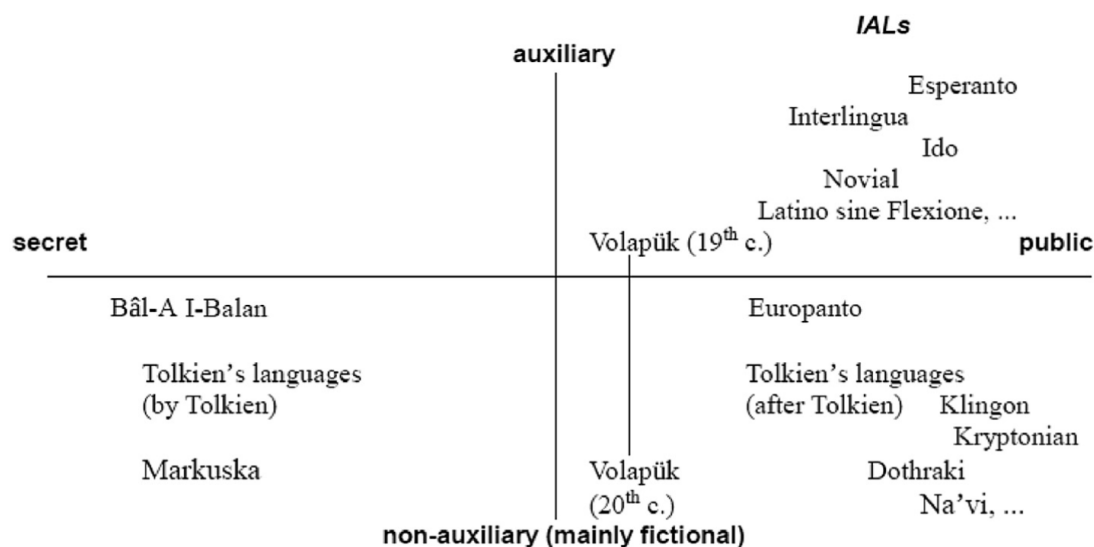


Figure 2.2: A taxonomy of constructed languages (Gobbo 2016)

word orders, head-initial relative clauses, fronted *wh*-phrases, and morphological regularity (Goodall 2022; Gobbo 2016). Section 2.2 further examines research focused on linguistic features of these languages.

In addition to classification based on their intended communicative functions, i.e. as philosophical or international auxiliary languages, there are also taxonomies based on other criterion. For example, another frequently used distinction is that of *a priori* and *a posteriori* (Schreyer and Adger 2021; Gobbo 2008; Schubert 1989; Schubert et al. 2001; Novikov 2022; Adelman 2014; Tonkin 2015). Languages described as being *a priori* are structurally entirely new (Tonkin 2015) and not based on existing languages, whereas so-called *a posteriori* languages are the opposite, drawing from aspects of specific natural languages (Schreyer and Adger 2021). Gobbo (2008) also proposed the dichotomy of *exoteric* (secret) and *esoteric* (public) languages, derived from Bausani (1974). Similar to critiques regarding the distinction between constructed and natural, such dichotomies for categorizing constructed languages are also argued by some linguists to be more accurately described as scales instead, with many languages falling somewhere in the middle (Novikov 2022).

A two-dimensional classification schema for constructed languages containing several notable examples is shown in Figure 2.2 (Gobbo 2016).

2.2 Prior Studies

In contrast to the abundance in literature and cross-linguistic analyses done on natural languages, similar research which also includes constructed languages is relatively sparse. In particular, while there is research that analyzes specific instances of linguistic differences between certain natural and constructed lan-

guages, large-scale cross-linguistic studies which utilize computational methods to classify the two based on linguistic features are practically nonexistent. Consequently, the present study is a somewhat novel approach. However, there is precedent for this research and the specific features examined, as well as computational approaches used, which this section will describe.

As previously mentioned, the creation of IALs such as Esperanto often involved the intentional simplification of particular language features to facilitate language acquisition, for instance having more regularity in their morphological systems. Intuitively, then, one would assume this translates to measurable differences in various aspects of linguistic complexity when compared to natural languages, which often have irregularities as a result of their development and evolution.

3 Methodology

In this section, I explain in greater detail the data and steps taken for preprocessing it, as well as the feature extraction and machine learning methods used in my experiment. Because of the broad nature of this study, several different experiments are done to test various linguistic features. Namely, these features are measurements of morphological complexity, type-token ratio (TTR), moving-average type-token ratio (MATTR), lexical entropy, text entropy, and character and word distribution entropies. Once these were calculated for each language, the task became that of binary classification using these values with the use of two models: a one-class support vector machine (SVM) and a decision tree. Finally, I also include a brief description of the various APIs and libraries used.

3.1 Data

In total, twenty-four languages are analyzed in this study. Six of these are constructed: Esperanto, Interlingua, Lingua Franca Nova, Volapük, Kotava, and Ido. The remaining eighteen are natural languages belonging to five different language families: German, English, Spanish, Polish, Vietnamese, Indonesian, Turkish, Tagalog, Hungarian, French, Finnish, Italian, Dutch, Occitan, Danish, Swedish, Afrikaans, and Icelandic. The language families represented are Austroasiatic, Austronesian, Turkic, Uralic, and Indo-European, with the latter comprising the bulk of the natural languages used. This diversity was intentional,

Language	Language Influences	Family	Word Order
Esperanto			SVO
Interlingua			
Lingua Franca Nova			
Volapük			SVO
Kotava			
Ido			

Table 3.1: Constructed languages and their respective language family influences used in the study.

To briefly introduce the constructed languages used here, Esperanto is the most widely spoken constructed language

The final constructed language used in this study is Lingua Franca Nova, created by Dr. C. George Boeree. Compared to the others discussed so far, it is much more recent in its creation, having first appeared in 1998 online. As a result of being both newer and more niche, there is considerably less existing

research related to it. Its inclusion in my dataset is primarily due to it being used on Wikipedia as an available source language, and thus also having a Wikimedia dump file.

As this study is cross-linguistic in nature, it would naturally be ideal to use parallel text corpora, as this would enable more conclusive comparative and comprehensive analysis. However, finding already existing parallel corpora that also includes constructed languages, particularly less common ones, posed a bit of a challenge.

3.1.1 Wikimedia

The data for this thesis comes from Wikimedia dump files. Wikimedia is a global movement and community founded on shared values, whose goal is to provide free and openly accessible information to everyone in the form of massive collaborative projects (which include, among others, the widely-used Wikipedia and Wiktionary). For a large, cross-linguistic study, massive databases with open-access make for an ideal source for corpora.

The dump files provide detailed, archived snapshots of the content from Wiki repositories for a specified point in time and are available in different formats as well as a wide selection of languages. All dumps used were XML-formatted and from the 2024-07-01 archive, containing articles together with their metadata¹. It is also worth mentioning here that there are some drawbacks to using these dumps for the present study. The files sizes vary considerably depending on the language, with the largest being roughly 22 gigabytes (English) and the smallest around 4 megabytes (Lingua Franca Nova), meaning all files do not contain the exact same articles. Additionally, the open and collaborative nature of Wikimedia means the articles are often authored by a multitude of different people, which can result in inconsistencies in the text such as with writing style. Similarly, it may also produce an imbalance in the amount of information provided across languages, with the same article in one language being considerably more detailed than in another, and inconsistent or low-quality translations, as Novikov (2022) noted to be the case for Wikipedia articles in Volapük. In essence, while Wikimedia was decided as the best available option for the task at hand, there are some unfavorable aspects of it too.

3.2 Data Preprocessing

Preprocessing text data is essential for natural language processing (NLP) tasks. As this study covers a broad range of different languages, including

¹<https://dumps.wikimedia.org/backup-index.html>

constructed languages which typically are low-resource, meticulous effort was made to obtain as close to a comparable set of texts as possible. This involved

Text data was initially extracted from the XML dump files with the use of WikiExtractor, a Python script (Attardi 2015).

Then, regular expressions were used to remove page titles, links, headers, fragments, and other extraneous symbols. The text was then made all lowercase and split by the periods—while also attempting to account for abbreviations—to make separate sentences, and all characters that were not part of the language’s writing system were removed in an attempt to have only that language’s words, without any foreign words that occasionally appear in the scraped Wikipedia texts. Remaining punctuation and numbers were also removed. The end result was a single corpus file corresponding to each language, with each line of the files being a single sentence.

3.3 Libraries and APIs

3.4 Feature Extraction

Since this study uses models for classification, it is necessary to first extract the features...

3.4.1 TTR & MATTR

TTR is a way of measuring lexical diversity. It is calculated using the following formula:

$$TTR = \frac{\sum_{i=1}^n \delta(w_i)}{n}$$

A glaring issue with TTR, however, is that it can vary widely based on a text’s length. The longer a particular text, the higher the likelihood of repetition occurring. There have been several solutions proposed to address this issue, one being MATTR. MATTR is given in the formula

$$MATTR_i = \frac{TTR_1 + TTR_2 + \dots + TTR_i}{i}$$

3.4.2 Morphological Complexity

The morphological systems of each language were analyzed using Morfessor,

3.4.3 Zipfian Distribution

3.4.4 Entropy

In information science, entropy means...

The entropy was calculated for the character and word distributions in each of the corpora, given by the following formula:

$$H = - \sum_{i=1}^n p_i \log_2 p_i$$

3.4.5 PCA

Principal Component Analysis was performed for dimensionality reduction.

3.5 Classification

3.5.1 Decision Tree

Decision Tree Classifier...

3.5.2 Random Forest

3.5.3 One-Class SVM

For outlier detection...

3.6 Evaluation of Classifiers

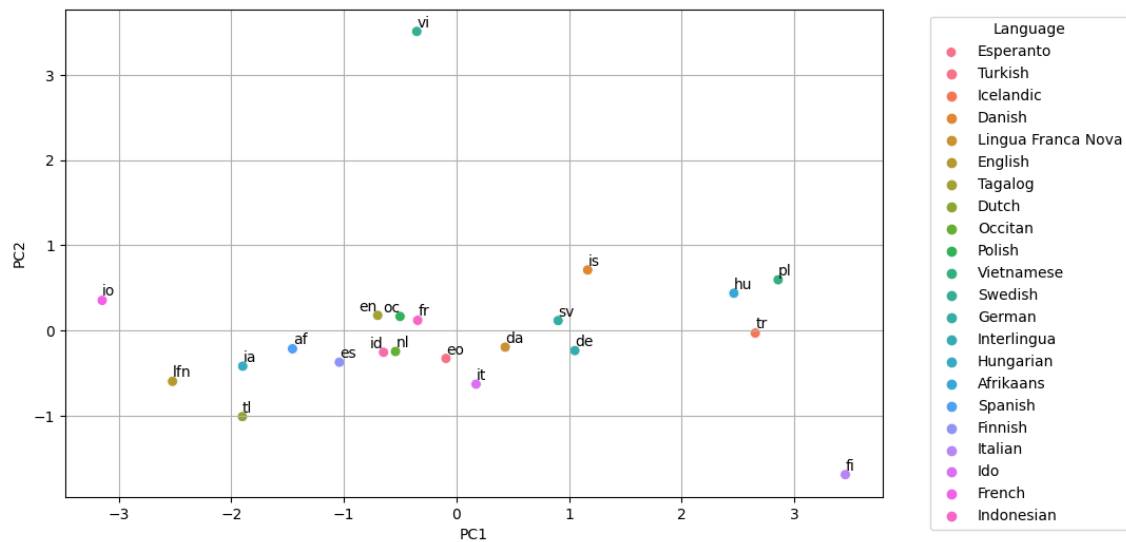


Figure 4.1: Principal Component Analysis on TTR, MATTR, Word and Char Distribution Entropy

4 Results

This section reports the results of each of the methods implemented.

4.1 Results of TTR & MATTR

4.2 Results of Morphological Segmentation

4.3 Results of PCA

A script was used to increase readability of the text in the graph².

4.4 Results of One-Class SVM

4.5 Results of Decision Tree

4.6 Results of Random Forest

²<https://github.com/Phlya/adjustText>

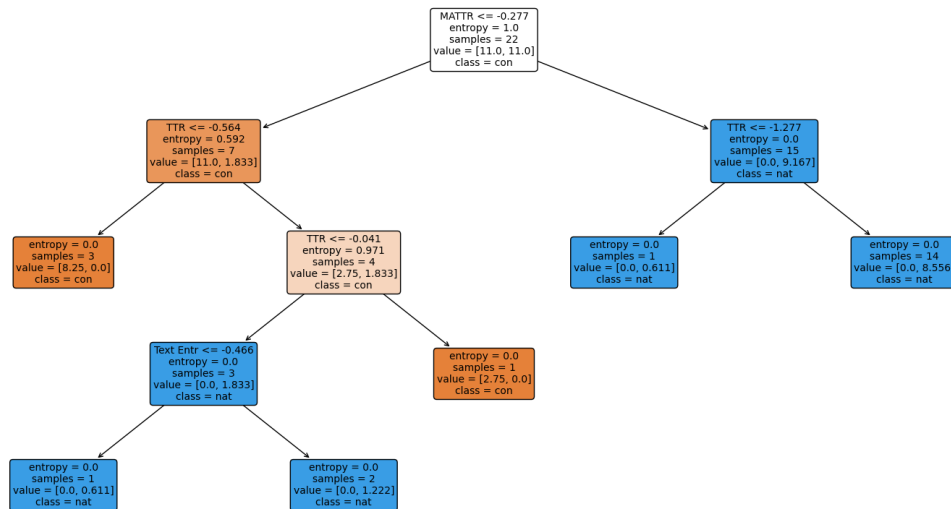


Figure 4.2: Decision Tree Classifier

5 Discussion

6 Conclusion

6.1 Future Work

The research presented in this thesis is far from encompassing all there is to the topic of defining language, and distinguishing between constructed and natural language. At present, this is an area of research with ample room for potential development.

Limiting factors: number of languages and which languages/language families, lack of real parallel corpora, problems associated with low-resource languages, relatively narrow scope of experimentation,

7 Acknowledgments

I would like to thank

References

- Adelman, Michael J. (2014). “Constructed Languages and Copyright: A Brief History and Proposal for Divorce”. In: *Harvard Journal of Law & Technology* 27, p. 543. URL: <https://api.semanticscholar.org/CorpusID:58553165>.
- Attardi, Giuseppe (2015). *WikiExtractor*. <https://github.com/attardi/wikiextractor>.
- Ball, Douglas (2015). “Constructed languages”. In: *The Routledge Handbook of Language and Creativity*. Ed. by Rodney H Jones. Routledge. Chap. 8. DOI: 10.4324/9781315694566.ch8.
- Bausani, A. (1974). *Le lingue inventate: Linguaggi artificiali, linguaggi segreti, linguaggi universali*. Collana di studi umanistici ‘Ulisse’. Ubaldini. ISBN: 9788834003879. URL: <https://books.google.de/books?id=z4GAngEACAAJ>.
- Gobbo, Federico (Jan. 2008). “Planned languages and language planning: The contribution of interlinguistics to cross-cultural communication”. In: *Multilingualism and Applied Comparative Linguistics* 2.
- (Oct. 2016). “Are planned languages less complex than natural languages?” In: *Language Sciences* 60. DOI: 10.1016/j.langsci.2016.10.003.
- Goodall, Grant (Sept. 2022). “Constructed Languages”. In: *Annual Review of Linguistics* 9. DOI: 10.1146/annurev-linguistics-030421-064707.
- Greenberg, Joseph H. (1970). “Language Universals”. In: *Theoretical Foundations*. Berlin, Boston: De Gruyter Mouton, pp. 61–112. ISBN: 9783110814644. DOI: doi:10.1515/9783110814644-003. URL: <https://doi.org/10.1515/9783110814644-003>.
- Jeffrey Punske (editor) Nathan Sanders (editor), Amy V. Fountain (editor) (2020). *Language Invention in Linguistics Pedagogy*. Oxford University Press. ISBN: 0198829876,9780198829874. URL: <http://gen.lib.rus.ec/book/index.php?md5=0B5CF2BFC00DCB569EBA10BD96AD68D4>.
- Large, J.A. (1985). *The Artificial Language Movement*. Language library. B. Blackwell. ISBN: 9780631144977. URL: <https://books.google.de/books?id=xaeCQgAACAAJ>.
- Libert, Alan (Jan. 2016). “On Pragmemes in Artificial Languages”. In: pp. 375–389. ISBN: 978-3-319-43490-2. DOI: 10.1007/978-3-319-43491-9_20.

- Novikov, Philipp (July 2022). “Constructed Languages as Semantic and Semiotic Systems”. In: *RUDN Journal of Language Studies, Semiotics and Semantics* 13, pp. 455–467. DOI: 10.22363/2313-2299-2022-13-2-455-467.
- Okrent, Arika (2009). *In the Land of Invented Languages: Esperanto Rock Stars, Klingon Poets, Loglan Lovers, and the Mad Dreamers who Tried to Build a Perfect Language*. Spiegel & Grau. ISBN: 9780385527880. URL: <https://books.google.de/books?id=E3UE9IoW27AC>.
- Sanders, Nathan (Sept. 2016). “Constructed languages in the classroom”. In: *Language* 92, e192–e204. DOI: 10.1353/lan.2016.0055.
- Schreyer, Christine and David Adger (Mar. 2021). “Comparing prehistoric constructed languages: World-building and its role in understanding prehistoric languages”. In: *Philosophical Transactions of the Royal Society B: Biological Sciences* 376. DOI: 10.1098/rstb.2020.0201.
- Schubert, Klaus (Jan. 1989). “Interlinguistics – Its Aims, Its Achievements, and Its Place in Language Science”. In: pp. 7–44. ISBN: 9783110886115. DOI: 10.1515/9783110886115.7.
- Schubert, Klaus et al. (Jan. 2001). *Planned Languages: From Concept to Reality*.
- Tonkin, Humphrey (Apr. 2015). “Language Planning and Planned Languages: How Can Planned Languages Inform Language Planning?”. In: *Interdisciplinary Description of Complex Systems* 13, pp. 193–199. DOI: 10.7906/indecs.13.2.1.
- Wilkins, John S. (1968). “An essay towards a real character, and a philosophical language, 1668”. In: URL: <https://api.semanticscholar.org/CorpusID:161991811>.

8 Appendices

Here I...

Language	Number of Words	Number of sentences
Icelandic	629995	41847
German	629987	37261
Polish	629997	42138
Ido	629990	43496
Afrikaans	629994	30737
Kotava	617400	48145
Hungarian	629946	39916
Lingua Franca Nova	628683	32188
Danish	629999	38260
Spanish	629978	24886
Interlingua	629996	32229
French	629983	27248
Occitan	629998	33762
Esperanto	629994	33317
Dutch	629997	34627
Turkish	629995	43573
English	629958	29574
Tagalog	629989	29855
Swedish	629998	36370
Vietnamese	629958	21115
Italian	629987	24487
Volapük	629999	55920
Indonesian	629997	34683
Finnish	629994	52637

Table 8.1: Lengths of each text after pre-processing.