# Turkronicles: Diachronic Resources for the Fast Evolving Turkish Language

Togay Yazar[1], Mucahid Kutlu[2], İsa Kerem Bayırlı[3]

[1]Department of Computer Engineering, TOBB University of Economics and Technology, Ankara, Turkey.
[2]Department of Computer Science and Engineering, Qatar University, Doha, Qatar.
[3]TOBB University of Economics and Technology, Ankara, Turkey.

Contributing authors: u.yazar@etu.edu.tr; mucahidkutlu@qu.edu.qa; isakerem@gmail.com;

**Abstract**

Over the past century, the Turkish language has undergone substantial changes, primarily driven by governmental interventions. In this work, our goal is to investigate the evolution of the Turkish language since the establishment of Türkiye in 1923. Thus, we first introduce Turkronicles which is a diachronic corpus for Turkish derived from the Official Gazette of Türkiye. Turkronicles contains 45,375 documents, detailing governmental actions, making it a pivotal resource for analyzing the linguistic evolution influenced by the state policies. In addition, we expand an existing diachronic Turkish corpus which consists of the records of the Grand National Assembly of Türkiye by covering additional years. Next, combining these two diachronic corpora, we seek answers for two main research questions: How have the Turkish vocabulary and the writing conventions changed since the 1920s? Our analysis reveals that the vocabularies of two different time periods diverge more as the time between them increases, and newly coined Turkish words take the place of their old counterparts. We also observe changes in writing conventions. In particular, the use of circumflex noticeably decreases and words ending with the letters "-b" and "-d" are successively replaced with "-p" and "-t" letters, respectively. Overall, this study quantitatively highlights the dramatic changes in Turkish from various aspects of the language in a diachronic perspective.

**Keywords:** Diachronic Corpora, Diachronic Analysis, Turkish Corpus, Frequency Analysis

# 1 Introduction

Languages undergo perpetual transformations over time. This evolution can be attributed both to natural factors such as semantic bleach and simplification and to cultural factors such as technological advancements and social developments. These changes might reduce the utilization of language models for historical texts due to differences in the meaning or spelling of words in their training data. Therefore, it is important to understand the historical evolution of languages.

Turkish language had a noticeably different path in terms of how it has evolved in the last century compared to other languages. In particular, after the establishment of the Republic of Türkiye[1] in 1923, cultural and technological modernization was the most urgent agenda of the reform program adopted by the new government. The Turkish language underwent two radical changes in the context of this new campaign towards modernization: one concerning its alphabetic system and one its lexical repertoire. In 1928, the Perso-Arabic script[2] in which Turkish had been written during the period of the Ottoman Empire was given up in favor of a version of the Latin script consisting of 29 letters. The second major change in the Turkish language was an attempt to "simplify" and "purify" the Turkish language by replacing words of Persian and Arabic origin, which were numerous during the Ottoman period, with words of Turkish origin (i.e., with words that are either historically Turkish or derivable by the rules of Turkish morpho-phonology). This process was concomitant with the establishment of the Turkish Language Association, widely known as TDK[3], in 1932 and can be understood to be part of the attempts at the crystallization of a new national identity.

In this work, we investigate how the Turkish language has changed since 1920s. In order to conduct our study, we first developed a diachronic corpus for Turkish. Specifically, we crawled issues of the Official Gazette of Türkiye (OGT) and the records of Grand National Assembly of Türkiye between 1920 and 2022. Given that both resources contain documents about governmental actions such as laws, regulations and discussions surrounding them, we think that our diachronic corpus is an important resource for analyzing the evolution of the Turkish language and the government's role in this transformation. Our corpus contains 45,375 documents, 842M words and 211K unique words. Using our corpus, we seek answers to the following research questions.

**RQ-1: How has Turkish vocabulary changed since 1920?** We divide our dataset into ten-year periods and compare the words used in each time period using different methodologies. We find that the words in two time periods diverge more as the time difference increases. While the frequency of the newly coined Turkish words increases over time, the frequency of their counterpart words with Arabic of Persian origins decreases. Around 75% of the words existing in 1920s have not been used between 2010 and 2019.

**RQ-2: How has the writing conventions changed since 1920s?** We observe the use of circumflex noticeably has decreased noticeably compared to the 1920s and 1930s. In addition, the last letter of the words changed over time based on Turkish phonology. In particular, we found that the use of words that end with "-b" (e.g., "kitab", which means book) decreases over time compared to its versions in which the last letter is "-p" (i.e., "kitap"). However, we observe a different pattern for words ending with "-d/t" letters: The percentage of words

---

ending with "-d" compared to the corresponding words with "-t", e.g., Ahmed vs. Ahmet in 2010s, is similar to their percentage in 1920s, although it has a decreasing trend since 1990s.

The main contributions of our work are as follows.

- We create the largest Turkish diachronic resources comprised of a diachronic corpus of formal Turkish documents, different kinds of diachronic word embeddings, a digitalized modern-old Turkish counterparts dictionary, and a python library that enables diachronic analysis.
- We explore the language change in Turkish since the 1920s using our corpus.
- We share our code and data to enable further research studies[4].

The rest of the paper is organized as follows. We first provide background information about Turkish for non-Turkish speakers in Section 2. Next, we discuss the related work in Section 3. In Section 4, we explain how we constructed our diachronic resources. We present our analysis and discuss our findings in Section 5. In Section 6 we discuss the limitations of our work and conclude in Section 7.

## 2 Background

Turkish belongs to the southwestern/Oghuz branch of the Turkic language family which includes languages such as Uigur, Uzbek, Kazakh, and Kyrgyz (Johanson, 1998). Its characteristic phonological feature is the assimilatory process of vowel harmony whereby a vowel shows agreement with the preceding vowel in terms of frontness and, to a more limited extent, roundness (Lees, 1961). Moreover, Turkish obeys various phonotactic constraints such as the absence of adjacent vowels inside words (with the exception of the "loan" words) as well as the ban on word-final voiced stop consonants such as [b], [d] and [g] (again, with the exception of some mono-morphemic words such as *ad*, meaning "name"). At the morphological level, Turkish is known as an agglutinative language in which inflectional suffixes attach to a nominal or verbal stem one by one, creating a structure similar to beads on a string. Syntactically, the default word order in Turkish is Subject-Object-Verb (SOV); however, other word order permutations are also acceptable under various prosodic and information-structural conditions, especially in spoken registers (see Lewis (2000), Underhill (1976); Kornfilt (1997); Göksel and Kerslake (2004) for extensive overviews of Turkish grammar).

## 3 Related Work

In this section, we discuss the studies in the literature from two different perspectives parallel to the contributions of this study.

### 3.1 Turkish Corpora

The NLP resources for Turkish are highly limited compared to English. However, the digitization of printed materials enabled the development of various Turkish corpora. METU Corpus (Say et al., 2002), and Turkish National Corpus (TNC) (Aksan et al., 2012) are general-purpose, genre-balanced, Turkish corpora. Both accommodate text resources from different

---

[4]URL is hidden due to the double-blind review process

genres, and the period of the text files is post-1990. However, they differ in the size of word count and annotation style. The former has 2 million tokens and XCES annotations, while the latter has 50 million words and provides part-of-speech tag annotations.

There are also larger Turkish corpora available such as BOUN Corpus (Sak et al., 2008), with 500 million tokens derived from web. Moreover, METU-Sabancı Tree-bank (Say et al., 2002) and IMST Turkish Dependency Tree-bank (Sulubacak et al., 2016) provide richer syntactic annotations such as morphological features, and dependency relations. However, none of these datasets clearly reflect language change, and hence, do not allow diachronic analysis.

To the best of our knowledge, there exists only one diachronic corpus for Turkish which consists of documents of parliamentary sessions between 1920-2015 (Güngör et al., 2018), named Corpus of Turkish Grand National Assembly. Since the reports are exact transcribed versions of the speeches made by the deputies, the corpus can reflect the historical evolution of the modern Turkish language. In our work, we developed a more comprehensive corpus by extending the temporal scope of parliamentary records up to 2022 and crawling issues of the Official Gazette of Türkiye published between 1922 and 2022. We believe that the resulting corpus provides better insights into the linguistic dynamics of Turkish and the political discourse of Türkiye throughout its history.

## 3.2 Diachronic Analysis

There are several studies that consider various aspects of the language change in a diachronic perspective in the literature. Michel et al. (2011) and Lieberman et al. (2007) focus on the evolutionary dynamics of English. They explore the grammatical changes through the history and attempt to reveal long-term patterns in linguistic change and the effects of cultural shifts. Their studies are mainly built on quantitative analysis of the frequency of the words across different time periods. As a different methodological approach, Pechenick et al. (2015) utilize information theory methods, such as Jensen-Shannon Divergence, to examine the evolution of the English language by exploiting the Google Books dataset introduced by Michel et al. (2011).

Many researchers use word embeddings to study semantic change throughout the years and reveal various facets and features of semantic change (Hamilton et al., 2016, Szymanski, 2017). In addition, the meaning-bearing nature of the word embeddings allows the evaluation of the validity of linguistic hypotheses about semantic change, such as the Law of Parallel Change and the Law of Differentiation (Xu and Kemp, 2015).

Although there exist several corpus-based diachronic analyses for English in the literature, the studies for Turkish are limited. These studies diachronically examine various aspects of the language. Salan and KABADAYI (2022) and Vahit (2003) provide extensive reviews on sound change in word-initial vowels of Turkish words and find instances of the phenomena by inspecting the dictionaries of different languages. Sultanzade (2012) qualitatively and quantitatively examine the valency change of a list of verbs in the Book of Dedekorkut by comparing them with modern Turkish correspondents. Aksan (1965) and Bahattin (2003) study semantic change in Turkish and primarily investigate the development of individual words to explain the mechanisms causing the semantic shift. To our knowledge, the study of Güngör et al. (2018) is the only work that analyzes the Turkish language in a corpus-driven diachronic way. They investigate the language change in their corpus by examining the frequency changes of near-synonym words and topic distributions using Latent Dirichlet Allocation. In our work,

we use different computational approaches to observe the change in the lexicon and writing conventions using a larger corpus.

# 4 Turkronicles

| Date | Original |
|------|----------|
| February 7, 1921 | (1) Hâkimiyet bilâ-kayd ü şart milletindir. İdare usulü, halkın mukadderatını bizzat ve bil-fiil idare etmesi esasına müstenittir. (2) Türkiye Devleti, Büyük Millet Meclisi tarafından idare olunur ve hükûmeti "Türkiye Büyük Millet Meclisi Hükûmeti" unvanını taşır. <br><br> *(1) Sovereignty unconditionally belongs to the people. The administration is based on the principle that the people themselves directly and actively manage their own destiny. (2) The State of Turkey is governed by the Grand National Assembly and its government bears the title of "Government of the Grand National Assembly of Turkey".* |
| October 4, 2020 | (1) Bu yönetmeliğin amacı, TOBB Ekonomi ve Teknoloji Üniversitesi Laboratuvar Okullarındaki eğitim-öğretim, yönetim, kayıt-kabul, devam-devamsızlık, nakil ile öğrenci başarısının tespiti ve işleyişe yönelik usul ve esasları düzenlemektir. (2) Laboratuvar okulları ile Üniversitenin öğrenci, öğretim elemanı ve öğretmenleri birbirlerinin dersleri ile kültür, sanat, spor ve sosyal faaliyetlerine katılarak müşterek etkinlikler gerçekleştirirler. <br><br> *(1) The purpose of this code is to regulate the procedures and principles regarding education, management, registration-acceptance, attendance-absence, transfer, and determination of student success, as well as the operation of TOBB University of Economics and Technology Laboratory Schools. (2) Students, faculty members, and teachers from the Laboratory School and the university participate in each other's courses, as well as cultural, artistic, sporting, and social activities.* |

**Table 1** Example sentences occurred from the official gazette of Türkiye from 1920 and 2022. The English translations are given in *italics*.

As the Turkish state takes an active role in changing the Turkish language, the official statements made by the government might be one of the best ways to observe how the Turkish language is affected by the state. Therefore, we create Turkronicles which is the first diachronic dataset using the official gazette of Türkiye (OGT), named "Resmi Gazete" and the official records of the Grand National Assembly of Türkiye (RGNAT). In this section, we first provide brief information about OGT and RGNAT. Next, we explain how we crawled the data and prepared it for analysis.

## 4.1 Official Gazette of Türkiye

OGT was founded on 7 October 1920 to inform about governmental actions and other topics such as statesmen's opinions about various issues. Its publication frequency has varied, from weekly to more sporadic schedules. Today, it is published every day except the holidays and Sundays based on the regulations implemented in 2009.

The content of OGT is a reflection of the administration process of the Türkiye. The issues of OGT mostly contain state-related news such as:

- Laws,

- Decisions of the Turkish Grand National Assembly and its internal regulations

- International treaties,

- Procedures about the dismissal, election, appointment, substitution, or resignation of the authorities such as vice president and high judicial members, ministers,

- Decisions regarding assignments, dismissals, and terminations of duty made by the president,

- Interior administrative decisions such as administrative jurisdiction changes and decisions regarding the establishment of municipalities.

As official documents often take place in the gazette, yielding a formal language with almost no typos and grammar errors. The documents might also include non-Turkish texts due to the obligation to publish international agreements. In addition, the first 1053 issues are originally written with the Ottoman Turkish alphabet. However, with the reform of the Turkish alphabet in 1928, OGT began using Latin letters[5]. Furthermore, they might contain non-sentence structures such as charts, tables, and others. **Table 1** provides example sentences extracted from OGT.

## 4.2 Extended Corpus of Grand National Assembly of Türkiye

As mentioned above, we have extended the temporal coverage of Güngör et al. (2018)'s corpus from 1920-2015 to 1920-2022. RGNAT consist of texts about all activities that occur during a session of the Grand National Assembly of Türkiye, including any kind of speeches, inspections, voting, noise, debates, schedules, agendas, reports, letters, and suggestions, which have been systematically documented since 1920. Since the meeting schedule of the parliament is variable from year to year, the publication frequency of these documents is irregular compared to OGT.

RGNAT and OGT have considerable overlap in terms of structural elements, such as charts and tables, and the topics covered. However, unlike OGT, parliamentary records capture a range of language styles from formal to colloquial, depending on the speaker and context.

It is worth mentioning that, as in OGT, documents that belong to 1920-1928 are written in the Ottoman-Turkish alphabet. However, translated versions of these documents are available in the official website.

## 4.3 Crawling

All the published documents of OGT and RGNAT are available at www.resmigazete.gov.tr and https://www.tbmm.gov.tr/Tutanaklar/TutanakMetinleri, respectively. To collect the documents, we used Scrapy[6] which is an open-source web scraping Python library. Using Scrapy's HTML parsing engine, we store the content of the documents extracted from the related web pages and their metadata. The metadata contains the publisher, publication date, file name,

---

[5]The issues with the Ottoman Turkish alphabet have been translated into modern Turkish in 2020 to honor the $100^{th}$ anniversary of OGT.

[6]https://scrapy.org

download link, and volume info. Eventually, we crawled issues of OGT published between 7 February 1921 and 31 December 2022, yielding 31,999 issues. For RGNAT, we collected 13,376 documents published between 23 April 1920 and 1 August 2022.

## 4.4 Preprocessing

Most of the downloaded files are in PDF format. We convert these pdf files into plain text files using PyPDF[7] tool to easily process the documents. The issues of OGT with issue numbers between 24,092 and 28,500 are shared as texts directly, instead of PDF files. For these documents, we directly extract the text content of these issues from the associated web pages.

We manually investigated the text extracting performance of PyPDF from the documents. We observed that the tool is generally successful but its performance is degraded in three cases: i) poorly scanned documents, ii) physically damaged documents, and iii) documents with non-sentence elements such as tables and charts.

In order to eliminate the errors introduced by the extraction tool and prepare the documents for further analysis, we perform the following pre-processing steps.

i) We reduce multiple consecutive spaces to a single space and convert different types of space characters (e.g., tab and non-breaking space characters) to a single regular space.

ii) We remove characters that cause problems in tokenization and/or do not have a representation in UTF encoding such as \xa0 and \xad.

iii) We use NLTK tool for tokenization.

iv) In our manual processes, we observe that the incorrect extraction from PDF documents yield very uncommon words. Thus, in order to eliminate those noisy ones, we remove words that contain any non-alphabetic characters and words that appear less than a particular threshold value. However, we realized that the quality of pdf files and, thereby, the performance of PyPDF changed over the years. Therefore, instead of using a single threshold value for all documents, we divide the data into ten-year time periods and we set the threshold value to $\left\lceil \frac{N}{10000000} \right\rceil$ where $N$ is the number of tokens of the time period under consideration. We filtered out the words below the frequency threshold. We observed that this type of filtering mechanism was effective in removing noisy words.

v) As Turkish is an agglutinative and morphologically rich language, analysis of surface-level words might be misleading. In order to detect lemmas, we use a morphological parser proposed by Öztürel et al. (2019). If the morphological parser cannot find a stem for a particular word, we apply the F5 method (i.e., using the first five letters as stems) which is shown as an effective stemming method by Can et al. (2006).

## 4.5 TDK Dictionary

We extract the modern Turkish equivalents of old Ottoman words from the 'Türkçeden Osmanlıcaya Cep Kılavuzu' (TDK, 1935) (Pocket Guide from Turkish to Ottoman Turkish), published in 1935 by the Turkish Language Association (Türk Dil Kurumu). This book aids speakers of contemporary Turkish in understanding and translating into Ottoman Turkish. It lists a wide range of modern Turkish terms alongside their Ottoman equivalents, presented in a certain format. Each entry is organized such that different senses of polysemous words,

---

[7]https://pypi.org/project/pypdf

synonyms, terms, abbreviations, and French equivalents are separated and indicated with specific markers, which are explained on the book's opening page. First, we converted the PDF version of the book into a plain text file through *tesseract* OCR (Optical Character Recognition), and subsequently, we employed regular expressions to generate a JSON file containing the entries. There are a total of 8647 newly coined Turkish words and their old counterparts available.

## 4.6 Ngrams

We have extracted unigrams, bigrams, and trigrams from each file, along with their frequencies. Preprocessing steps were applied to the tokens, as in the construction of other resources. Both surface-level forms and lemma frequencies can be readily accessible. In addition, we've organized the n-grams into the time periods and marked them with timestamps. One can analyze the distribution of these ngrams to examine the change in the usage of individual words or phrases through time. Moreover, we provide the association strength of unigrams using the PPMI measure. This measure provides insight into the strength of association between unigrams, allowing for a deeper understanding of language usage patterns and relationships.

These resources enable users to easily conduct queries pertaining to frequency changes over time, discover linguistic patterns, and test hypotheses about the Turkish language across different time periods.

## 4.7 Lingan: A Python Library for Linguistic Analysis

We developed a Python library to conduct diachronic analyses and facilitate the reproducibility of our experiments.

### 4.7.1 Components of Lingan

Lingan is fundamentally based on three different layers of abstractions: *Data*, *Container*, and *Operation*. The *Data* component is a representation of actual data that is used in linguistic analysis. Practical equivalents of this component include data types such as word embeddings, vocabulary, and ngrams. Members of this class contain both static and derived features of the relevant data. Generally, the responsibility of this component is to manage the data.

*Container* represents collectively created text data. Technically speaking, the *Container* component is a common interface for the hierarchical structure of the container classes used to wrap objects representing data, meaning that data objects should be wrapped by a Corpus object to perform diachronic analysis. Within Lingan, there are two classes derived from this interface, namely *DiachronicCorpus* and *Corpus*. *Corpus* class also includes attributes to interact with text files. For example, *TokenProcessor*, which sequentially transforms tokens, and *TextProcessor*, which is in charge of processing streamed raw text data, play an active role within this class. On the other hand, *DiachronicCorpus* is a composite object that consists of many *Container* objects, such as *Corpus* and *DiachronicCorpus*, marked with a timestamp of their time periods. *Corpus* and *DiachronicCorpus* together constitute a tree structure. An example of such a tree is depicted in Figure 1.

The Operation component is responsible for the algorithms performed on the hierarchical structure of a *Container*. Additionally, the *Data* objects should be generated by subclasses of *Operation*. This is a design decision ensuring system-wide integrity and consistency.
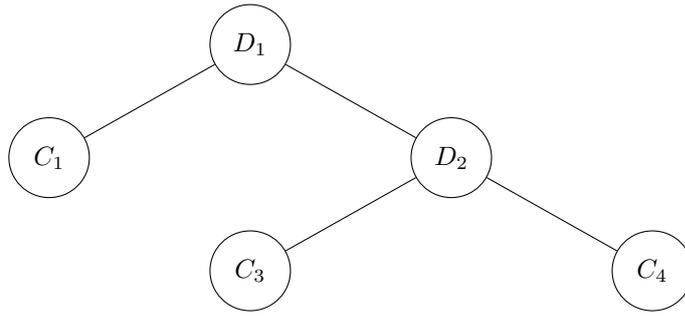
**Fig. 1** An example of a hierarchical tree structure consisting of Corpus and DiachronicCorpus objects. Nodes with $D$ and $C$ represent *DiachronicCorpus*, respectively. $D_2$ contains two *Corpus* objects, $C_3$ and $C_4$. $C_1$ and $D_2$ together compose $D_1$.

Also, operations are separated from the data structure, *Container*, to enhance flexibility and maintainability, e.g., without modifying the composite structure, one can easily define new operations. However, the design of *Operation* enforces the classes derived from *Operation* to implement functions specific to the *DiachronicCorpus* and *Corpus* types. In Lingan, we have implemented the functionality used in the section 5. The functions are readily available to the users of the library. We share a sample piece of code below to show the easy usage of the predefined operations. The task to be performed is to calculate the relative frequency of a specific word.

```
1  corpus_1930 = Corpus.load("path/to/corpus_1930")
2  corpus_1980 = Corpus.load("path/to/corpus_1980")
3
4  dia_corpus = DiachronicCorpus(corpora=[corpus_1930, corpus_1980])
5  frequency_series = dia_corpus.perform(Frequency(word = "belge", normalize=True))
6  print(frequency_series)
```

**Fig. 2** An example usage of Lingan. This code piece computes the relative frequency of *belge* (document) across time periods through pre-defined function *Frequency*.

The code assumes that *Corpus* objects are serialized and saved locally. In the first two lines of the program, corpus of 1930–1939 and 1980–1989, the variables *corpus_1930* and *corpus_1980*. In line 4, a diachronic corpus is initialized by constructing the tree structure mentioned above. In line 5, a new *Frequency* instance, which is the concrete object responsible for calculating the relative frequency of a given word, is created and passed as an argument to the *perform* method of *dia_corpus*. The result is stored in the variable *frequency_series*.

### 4.7.2 Available Operations

We have implemented various data structures, containers, and operations that can be utilized in a diachronic analysis. Out-of-the-box functionality of the library is listed in Table 2 with their names, data structures, on which the operations are performed, and short descriptions. In implementing these functions and other parts of the framework, libraries such as *numpy* and *scipy* were extensively used. Note also that the entire operations in the framework are

not listed in 2; There are also different types of operations such as *CreatePPMIMatrix, CreateSVDEmbeddings, CreateNgrams, CreateVocabulary* which create data models defined in the framework.

### 4.7.3 Extensibility

Due to the flexible nature of Lingan, our framework can be easily extended in terms of the fundamental components of the architecture. Users can define their custom data types and operations. Custom types should conform to the related interface. The proper way to achieve this is to extend relevant classes and interfaces. For clarity, we provide a showcase example where the task is to calculate the total number of sentences in a diachronic corpus. However, the data structure and the function for this task are not defined in Lingan. One should first define the data model:

```python
class Sentences(Data):
    def __init__(self, sentences: list[str], *args, **kwargs):
        super().__init__(*args, **kwargs)
        self.sentences = sentences
```

**Fig. 3** Defining a new *Data* component to model sentences in the corpus.

The data model inherits from *Data* interface to specify that *Sentences* is a newly defined data structure in the framework. Next, the logic for counting sentences is implemented in a subclass of *Operation*.

```python
class NumberOfSentences(Operation):
    def on_corpus(self, corpus: Corpus) -> int:
        return len(corpus.sentences)

    def on_diachronic_corpus(self, diachronic_corpus: DiachronicCorpus) -> int:
        total = 0
        for c in diachronic_corpus:
            total += c.perform(self)
        return total
```

**Fig. 4** Defining a new *Operation* to calculate the total number of sentences on a diachronic corpus structure.

*NumberOfSentences* contains two methods: *on_corpus* and *on_diachronic*. These are abstract methods from *Operation* interface that every subclass should implement to interact with the composite *DiachronicCorpus* structure. NumberOfSentences class traverses each element of the Corpus in a tree structure one by one, and defines its operation recursively according to the type of each element. Here in line 7, there is a method named *perform* and bounded to *Corpus* object. This is one part of the double dispatch mechanism of the framework: *perform* method of currently processed *CorpusContainer* object is invoked with an *Operation* instance as its argument, and inside the *perform* method CorpusContainer object *c* calls one of *on_diachronic_corpus* and *on_corpus* according to its type and passes itself as an argument, e.g., if the element is an instance of the Corpus class, then it invokes the *on_corpus* method of the operation.

| Operation | Datastructure | Description |
|---|---|---|
| *Exists(word, time_range)* | *Vocabulary* | Checks whether *word* exists in a diachronic corpus within the time range *time_range*. |
| *Frequency(word, time_range)* | Vocabulary | Returns a time series where each member is the frequency of *word* in each corpus within the *time_range* |
| *MergeVocabulary(time_range)* | Vocabulary | Merges the vocabularies in the *Diachronic-Corpus* and returns the composed vocabulary |
| *FilterFrequency(threshold, time_range)* | *Vocabulary* | Filters the vocabulary of each corpus, i.e, removes the words whose frequency is below *threshold* |
| *VocabularySimilarity(time_range)* | Vocabulary | Returns a matrix where each element results from Jaccard Index between the vocabularies of different time periods. |
| *VocabularyDistance(time_range)* | Vocabulary | Returns a matrix where each element represents Jensen-Shannon divergence between different time periods |
| *MorphemFrequency(morpheme, time_range)* | *Vocabulary* | Returns the usage frequency of *morpheme* within the *time_range*, e.g., the usage of $\hat{a}$ across the time periods |
| *WordsWithMorpheme(morpheme, time_range)* | *Vocabulary* | Returns the words in each time periods that contains *morpheme* |
| *WordsEndWith(suffix, time_range)* | *Vocabulary* | Returns the words in each time period that ends with *suffix* |
| *WordsStartsWith(prefix, time_range)* | *Vocabulary* | Returns the words in each time period that starts with *prefix* |
| *UniqueWordCount(time_range)* | *Vocabulary* | Returns a time series where each element is the total number of unique words in the vocabulary of each time period. |
| *CommonWords(time_range)* | *Vocabulary* | Returns a set of words that exist in every time period |
| *AverageWordLength(time_range)* | *Vocabulary* | Returns an array of numbers where each element represents the average length of unique words in each time period |
| *NgramCount(time_range)* | *Ngrams* | Returns a time series data of the total number of ngrams in each time period. |
| *CoFrequency(u, v, time_range)* | *Embeddings* | Returns a time series of cooccurence frequency of the word *u* and *v* within the *time_range*. |
| *Collocations(word, k, time_range)* | *Embeddings* | Returns an array where each element is a set of words. These sets with size *k* correspond to the words with the highest collocation value with the *word* for each period |
| *Association(u, v, time_range)* | Embeddings | Returns a time series array. Each element in this array corresponds to the collocation value between the words *u* and *v* in the respective time period. |
| *Similarity(u, v, time_range)* | Embeddings | Returns a time series array where each element is the similarity between *u* and *v* in respective time periods |
| *AlignedMostSimilar(word, k, target_period, base_period)* | Embeddings | Returns *k* most similar words to *word* by aligning the embedding space of *target_period* to that of *base_period* |
| *SemanticChange(word, time_range)* | Embeddings | Returns a time series array where each element is the distance from the vector of *word* in the starting period. |
| *MostSimilar(word, k, time_range)* | Embeddings | Returns an array where each element corresponds to a specific time period and the *k* closest words to *word* in that time period |

**Table 2** The list of currently available operations in the library

### 4.8 Embeddings

We provide three types of diachronic embeddings: PPMI (Positive Pointwise Mutual Information), SVD (Singular Value Decomposition) of PPMI, and CBOW (Continuous Bag of Words) embeddings. First, preprocessed text files are grouped into 10-year time periods according to their publication date. For each period, we count term-to-term cooccurrences to construct the PPMI matrix. The following formula is used to fill the elements of the matrix:

$$PPMI(u,v) = \max\left(\log\frac{p(u,v)}{p(u)\cdot p_\alpha(v)}, 0\right) \tag{1}$$

where $p(u)$ and $p(u,v)$ are the marginal probability of word $u$ and the joint probability of words $u,v$ respectively. It is well known that PPMI is very sensitive to infrequent events. So, $p_\alpha(v)$, smoothed unigram distribution (Mikolov et al., 2013) of word $u$, with $\alpha = 0.75$ is used to alleviate such negative effects. Also, an unweighted context window with size 2 was employed to relate the target word to the context words. After that, SVD factorization of PPMI matrices has been taken. In the SVD approach, we calculated the vector representations of the words by $W = U\Sigma^{1/2}$ where $U$ is left singular vectors and $\Sigma$ is the singular values of the PPMI matrix. The size of the embeddings is 300. Note that, unlike the classical SVD implementation, we take the square root of $\Sigma$ which has been shown to improve the quality of the SVD embeddings (Levy et al., 2015).

Finally, CBOW embeddings (Mikolov et al., 2013) are created using *gensim* (Rehurek and Sojka, 2010) library for each time period. We use context window size of 2, $alpha = 0.75$, and embedding size = 300 for the CBOW algorithm as in SVD embeddings. Furthermore, CBOW-specific parameters such as the number of negative words and downsample rate are chosen to be 5 and 0,00001.

The non-unique nature of SVD and the randomized processes involved in CBOW embeddings, prevent direct comparison of embeddings from different periods (Hamilton et al., 2016). Word embeddings for the same word, trained at distinct times and with the same parameters, can still be different from each other. Specifically, two embedding spaces may be rotated, translated, or dilated via a transformation matrix $R$. Therefore, we align the embedding matrices $W_{t_1}$ to $W_{t_2}$ using the Orthogonal Procrustes Problem (Schönemann, 1966). The solution to the Orthogonal Procrustes problem tries to find a transformation matrix R that minimizes the Frobenius norm of the squared difference of two embedding matrices:

$$\arg\min_R \|W_{t_1}R - W_{t_2}\|_F^2 \tag{2}$$

There is also an orthogonality constraint on the optimization of R such that R should satisfy $R^T R = I$. R can be obtained by first taking the SVD factorization of $M = W_{t_2}{}^T W_{t_1}$ and then multiplying the left singular vectors $U$ and right singular vectors $V$, $R = UV^T$. This methodology makes diachronic analysis available on embeddings. Therefore, we provide aligned embeddings of consecutive time periods and a transformation matrix for each pair of time periods, along with synchronic embeddings for further analysis.

## 5 Data Analysis

In this section, we provide an analysis of our corpus to get more insight about it and how it can be utilized in future studies. In particular, we first provide statistical features of the corpus

(Section 5.1). Next, we diachronically analyze the corpus to understand how the Turkish language has changed over time (Section 5.2).

## 5.1 General Statistics

**Table 3** provides the general statistics about Turkronicles. The dataset contains 45,375 documents and its size is 8.5 GB in text format. Before our filtering step, the dataset consists of around 849 million tokens and the number of unique stems is 1,961,044. After the filtering process, the total number of tokens and unique stems were reduced to 842,957,298, and 211,775, respectively.

.

| | |
|---|---:|
| The number of documents | 45375 |
| The number of words before filtering | 849,335,014 |
| The number of words after filtering | 842,957,298 |
| The number of unique surface level words | 10,689,405 |
| The number of unique stems | 1,961,044 |
| The number of unique stems after filtering | 211,775 |
| Average token count per document | 18,718 |

**Table 3** Descriptive statistics of Turkronicles.

## 5.2 Vocabulary Change Across Years

We first divide the documents into 10-year time period and compare each time period from different aspects. In all our calculations, we use lemmas (and stems for the words we the F5 stemming method).

Firstly, we compare vocabulary size across time periods. **Figure 5** shows the number of words for each ten-year time period. We observe that the vocabulary size is balanced in almost all time periods. Interestingly, the vocabulary size in 1940-1949 is higher than the others. The vocabulary size for 2020-2022 is less than others due to the limited number of documents for this time period.

Next, we turn our attention to the vocabulary distance across different time-periods. In particular, we first create a separate list of unique words for each 10-year period. Next, we compute both Jaccard similarity and Jensen-Shannon Divergence (JSD) between documents in different time intervals. We chose Jaccard due to its high interpretability and JSD to better show how the vocabulary changes over time.

**Figure 6** shows the Jaccard similarity scores for the vocabulary of each 10-year period in heatmap format. We observe that the Jaccard similarity between two consecutive time periods is less than 50% in around half of the cases. Furthermore, the similarity between documents in the 1990s and documents in the 1920s is approximately 0.2%. To illustrate this huge vocabulary change, we rewrite the first sentence shown in Table 1 using modern
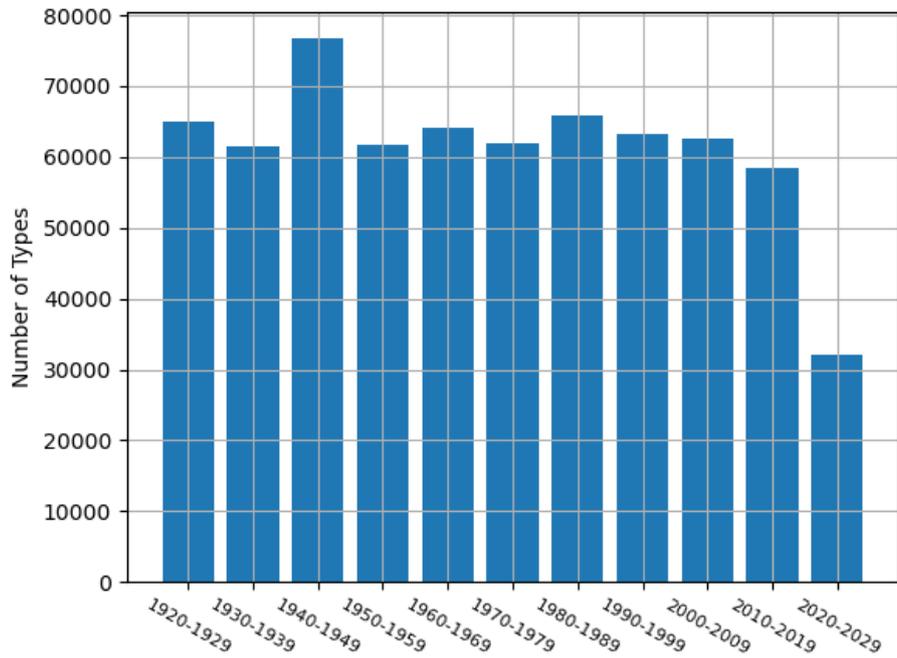
13

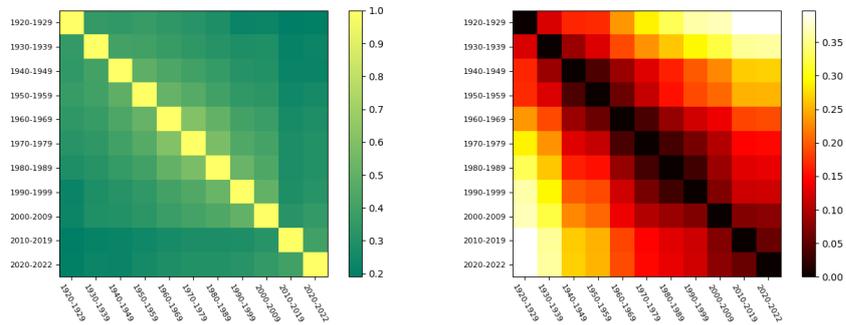**Fig. 5** The number of unique lemmas/stems for each 10-years time period.



**Fig. 6** Jaccard similarity values of the vocabularies of 10-year periods.

Turkish words[8]:

---

[8] the original words are written in parentheses, while their modern equivalents are highlighted in red

*Egemenlik (Hâkimiyet) kayıtsız şartsız (bilâ-kayd ü şart) milletindir. Yönetim (idare) şekli (usulü), halkın yazgısını (mukadderatını) doğrudan doğruya (bizzat) ve fiilen/gerçekten (bilfiil) yönetmesi (idare etmesi) esasına dayanmaktadır (müstenittir).* **Figure ??** shows the JSD scores between every pair of 10-year time periods. The first salient aspect of the heatmap is that as the distance between the compared time periods increases, the divergency increases in parallel.

In order to further investigate the vocabulary change, we rank the words based on their contribution to the JSD score. **Figure 7** shows the 60 words that cause the most divergency between documents written in 1930-1939 and 1980-1989. We are not able to show the other comparisons due to the space limitation. In the figure, the red and blue bars represent words that are more characteristic in 1930-1939 and 1980-1989, respectively. The length of the bar indicates the magnitude of the contribution of each word to the overall JSD score.

| 1930-1939 | 1980-1989 | Meaning |
|---|---|---|
| vekil | bakan | minister |
| sene | yıl | year |
| umumi | genel | general |
| reis | başkan | president |
| heyet, encümen | kurul | committee |
| vesika | belge | document |
| icra | uygula | perform |
| mucip, lazım | gerek | required |
| aza | üye | member |
| idare (et) | yönet | manage |
| sayı | numara | number |
| layiha | tasarı | pleading |

**Table 4** The words that have similar or identical meaning but were prevalent in different time periods.

Our results imply dramatic changes in Turkish vocabulary in the last 100 years.

In Figure 7, we observe that newly coined Turkish terms replaced the corresponding Arabic and Persian-origin words in the period of 1980-1989. **Table 4** lists all the word pairs that have the same or similar meaning but appear in different time periods in Figure 7. For instance, both the words "*reis*" and "*başkan*" mean "president" in Turkish, but the word *reis* is one of the most divergent words of the 1930-1939 period while the word *başkan* is one of the most divergent words of the 1980-1989 period. Furthermore, the words *gerek* and *lazım* were used as replacements for the word *mucip* in 1980-1989 period. Moreover, the word *kurul* was used to replace two different words, *heyet* and *encümen*.

We also observe that some words appear as divergent due to content or style change in the documents being compared. For example, English words such as *the* and *of* are in the list of divergent words because international agreements and contracts are included in the
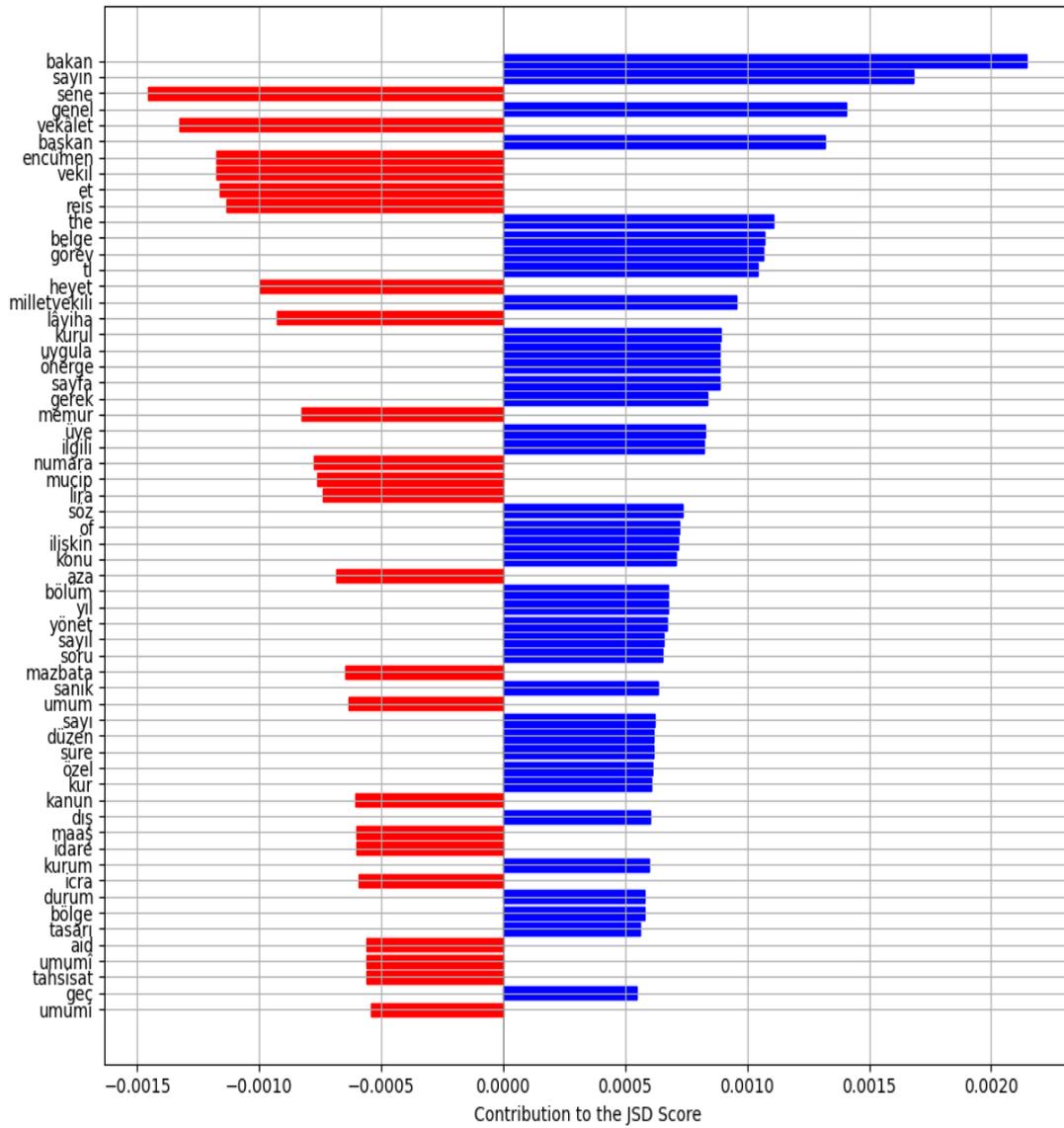
**Fig. 7** The first 60 words from 1930-1939 and 1980-1989 ordered by the individual contributions to the Jensen-Shannon divergence. The sign of the values indicates the corpus in which individual words are relatively frequent; the bars to the left represent the words that are more common in 1930-1939, while the bars to the right correspond to the words that are more frequent in 1980-1989.
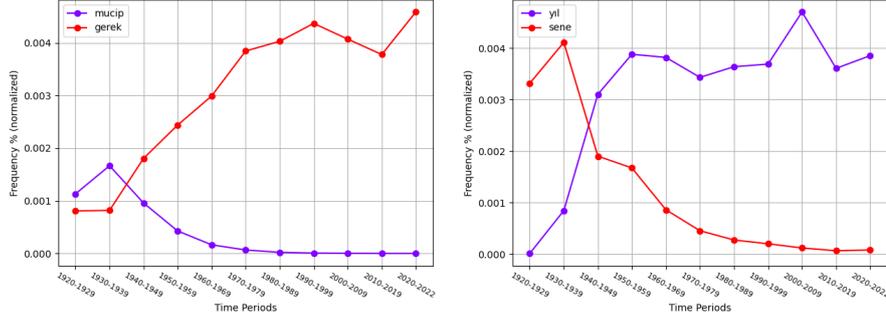
16

**Fig. 8** Normalized frequency of words that have the same meaning: gerek vs. mucip (required) and yıl vs. sene (year).

documents of the 1980-1989 period. Regarding the style change, the words *TL* which is the abbreviated version of *Türk Lirası* (Turkish Lira) and *lira* are two distinct terms from different time periods, indicating a shift towards adopting the abbreviation *TL* for *lira*.

To shed more light on how words are introduced to replace Arabic-Persion origin words, we focus on two specific pairs of words with identical meanings: i) *gerek* vs. *mucip* (required) and ii) *yıl* vs. *sene* (year). We calculate the frequency of each word and then normalize their frequency by the total number of tokens for each time period. **Figure 8** shows the normalized frequency of these words. While both the words *gerek* and *mucip* (both means "required") existed in the 1920s, *mucip*, which is an Arabic-origin word, is used more frequently than the word *gerek*. However, in the following time periods the word *gerek* becomes more popular than *mucip*, and the word *mucip* does not appear in the documents since the 1980s.

In our second example, we observe another interesting case. The word *yıl* does not exist in the 1920s in our corpus. However, it becomes popular to the extent that it is more prevalent than the word *sene* in all documents written after 1930s.

In our last analysis of vocabulary change, we investigate the presence of words used in the 1920s across subsequent time periods. **Figure 9** shows the number of words used in 1920s for the subsequent time periods, i.e., *survived* words. As expected, the number of survived words decreases as the time difference increases. Considering that there are around 65,000 words used in documents of 1920s (See Figure 5), around half of the words are not used in the subsequent years.

## 5.3 Vocabulary Change on Diachronic Embeddings

We perform a diachronic analysis on word embeddings to further emphasize the change in the vocabulary, particularly the effect of the purification of Turkish. To be consistent with JSD analysis, we chose the same time intervals 1930-1939 and 1980-1989 using the same words listed in the second column of Table 4. First, we train a CBOW model for both of the time periods. The configuration of the parameters of the models is specified in 4.8. After the models are trained, the word embeddings of 1980-1989 are aligned to 1930–1939 using the solution of the Orthogonal Procrustes method. Similarity scores between any two word vectors are calculated using cosine distance.
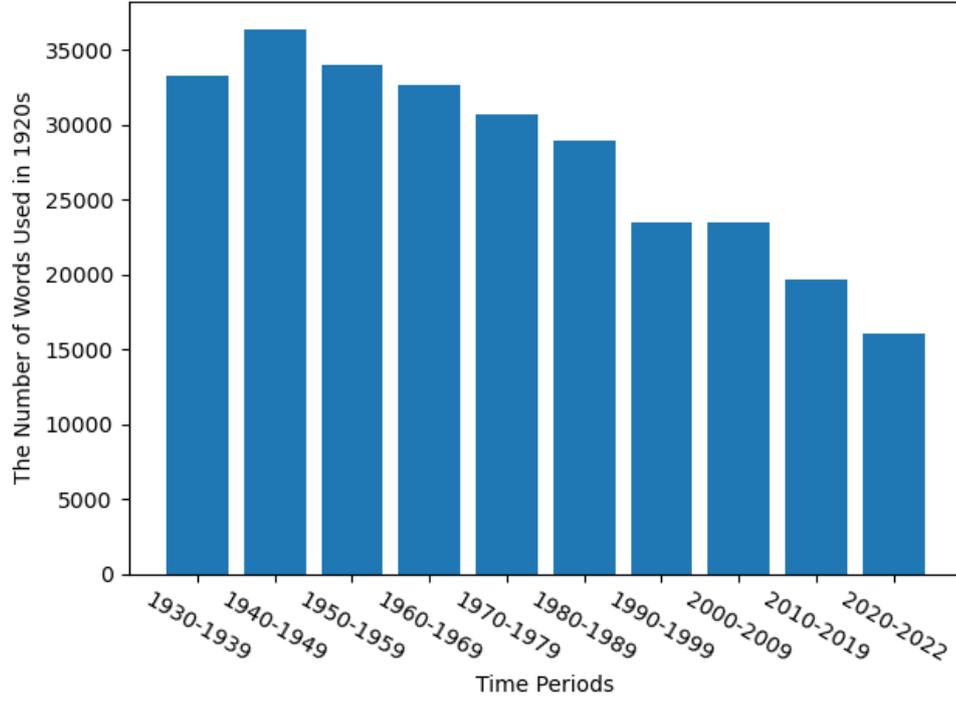
**Fig. 9** The number of words used in 1920s for each time-period.

| 1980-1989 | 1930-1939 |
|---|---|
| bakan | vekâlet, iktıs, îktıs, vekâl, içtimaî, maarif, tktıs, iktisat, nafıa, îcra |
| yıl | sene, yıl, yılma, ayı, takvim, aylık, seneye, iptida, *dörd*, katıl |
| belge | vesika, ibraz, istek, vesai, makbu, musaddak, mütea, vesaik, talih, makbuz |
| gerek | icabe, iktiza, göre, icap, ayrıç, kanunî, zarurî, kabîl, lüzum, önce |
| başkan | reis, müteşekkil, seçim, zatte, müsteşar, inha, vekâlet, mütal, riyaset, seç |
| genel | müdür, umum, işle, umumî, idare, teşki, denizyolu, genel, havayolu, îdare |
| kurul | seçim, heyet, îcra, ödev, baro, seçilir, yönetim, inha, teşekkül, seç |
| uygula | tatbik, göre, gözet, dışmd, hüküm, önce, tatbikat, icabl, istisnaî, tatbi |
| üye | seç, seçilir, seçilmiş, seçim, üye, intihap, seçi, zatte, müntahap, müntehap |
| yönet | talimatname, talimat, nizamname, teşkilat, bölüm, teşki, izahname, sayıl, ilgili, plânl |
| numara | numara, yazı, değiştiri, aşağı, ilişik, gösteri, sayı, mezkûr, ün, yaz |
| tasarı | lâyih, lâyiha, eneüm, encüm, değişik, mütenazır, tadil, bazı, encümen, ncü |

**Table 5** The words that have similar or identical meaning but were prevalent in different time periods.

Table 5 shows some of the most characteristic divergent words of 1980-1989 and their 10 most similar words in 1930–1939. It can be observed that the old Arabic equivalents of the modern Turkish words are included in the closest 10 words as a result of the alignment of the word vectors. Underlined words are the near-synonyms of the words in the second column. Words colored with red stand for OCR errors, and words colored with blue are words that the morphological analyzer fails to stem; hence, they are the result of the F5 stemming process. We made this distinction since no word in the dictionary starts with words with red color, while words with blue color can be found as a prefix of a word in the vocabulary. Therefore, words blue words can be estimated from their similar words and the context. For example, in the first row of the table, the most similar 10 words of the old Arabic version of the modern Turkish word *bakan* (minister) is in the last column. *iktıs* is colored blue because it is probably the stemmed form of *iktısad-* (economy). Note also that *iktisat* which is an Arabic-originated word, is presented among the most similar words of *bakan*. *iktisat* and *iktısad* corresponds the same meaning, *ECONOMY*, and *iktisat* is another form of *iktısad*. Also, there is a third option for the meaning *ECONOMY*: ekonomi. The word *ekonomi* which originated from French, économie, was later introduced to replace *iktısad* (TDK, 1935) and prevails over the other two words, e.g in 1980-1989, the relative frequencies of *iktisat* and *ekonomi* are respectively $1.33 \times 10^{-5}$ and $16.9 \times 10^{-5}$, and *iktısad* has become extinct. Furthermore, the surrounding words of *iktisat* in 1930-1939, *âli, celal, program (program), abdülhalik, vekâlet (ministry), nafia (development), vekil (minister), millî (national), sıhhat (wellness), mustafa*, which is extracted with the help of PPMI matrix of 1930-1939, belong to the governmental context and *ali, celal, abdülhalik, mustafa* are the first names of the authorities of the economy of Turkey. Whereas surrounding words of *iktisat* in 1980-1989 are *teori (theory), matematik (math), siyasal (political), kongre (congress), maliye (finance), politika (policy), teşebbüs (attempt), kıymet (value), teşekkül (organization)*. It can be observed that *iktisat* has become used in the academic context as well. Another example of the lexical change can be found by aligning the word vector of *adet* (unit), which is not presented in the table 4, to the vector space of 1930-1939. The most similar word from 1930-1939 is *aded*. These results imply that such types of lexical change can be effectively discovered by examining diachronic word embeddings.

Diachronic embeddings have another interesting property. Although, a word or a concept is not present in a vocabulary, aligning a word from a different time period may reveal related words. To be more specific, words similar to the aligned word are often related to the relevant concept. As an example from our corpus, the word *televizyon* (television) is not present in the vocabulary of 1930-1939, and TELEVISION was not a well-known concept in these years. We aligned the vector of *televizyon* from 1980-1989 to the vector space of 1930-1939. The most similar 10 words of the aligned vector of *televizyon* are radyo (radio), sinema (cinema), tiyatro (theater), mecmua (magazine), arsıulusal (international), rehber (guide), broşür (brochure), reklâm (advertisement), telsiz (radio), konser (concert). Most of the similar words belong to the concept MEDIA as does *televizyon*.

## 5.4 Change in Writing Conventions

We conduct two different analyses for changes in writing conventions. We first explore how the word endings have changed and then we focus on the usage of circumflexes.

The words originating from Turkish do not end with the letters "-b", "-c", "-d", and "-g" with few exceptions. However, there are several loanwords from Arabic and Persian that

end with one of these letters. We observe that these names have been written in two different ways in which the last letter is changed from "-b/c/d/g" to "p/ç/t/k" such as *Ahmet* vs. *Ahmed*. In order to observe the transformation of these loanwords, we first detect words in which a single letter is different based on this phonetical rule, e.g., Ahmet vs. Ahmed[9]. Next, we count the words that end with -d and -b and their counterparts ending with the letter -t and -p for each 10-year time period. We ignore words ending with -c/-ç and -g/-ğ letters due to their low prevalence. Subsequently, we calculate the ratio of words ending with -b/-d letters compared to the words ending with -p/-t letters. The results are shown in **Figure 10**.
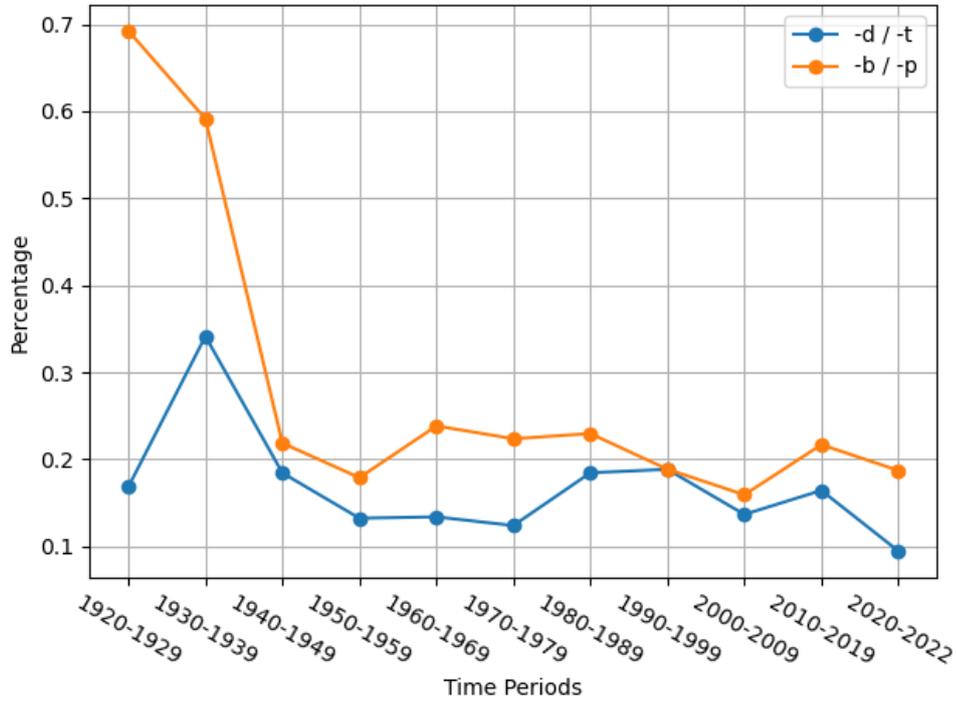


**Fig. 10** The ratio of words ending with -b/-d letters compared to the words ending with -p/-t letters in each time period.

We observe that the ratio for both letter pairs is less than 1, suggesting that ending words with "-t" and "-p" letters is more common than using "-d" and "-b" letters in all time periods. However, this might be because of the morphological analysis tool we use. In particular, as it is developed based on modern Turkish grammar rules, it might identify stems as if they end with "-t" and "-p" in some cases. Therefore, it is important to focus on the trend instead of actual values.

---

[9]We removed the word *et* (which is an auxiliary verb) because of its extremely high prevalence compared to others.

The percentage of words ending with "-b" significantly decreases between 1920s and 1940s. Afterwards, it fluctuates between the levels of 0.16 and 0.24. However, we observe a different pattern for words ending with "-d". Interestingly, the percentage of words ending with "-d" letter first increases (from 1920s to 1930s) and then decreases (from 1930s to 1940s). In the following time periods, it fluctuates between 0.1 and 0.2. These fluctuations might stem from the limitations of the corpus or due to people's resistance to reforms in language. For instance, it is still common that many people in Türkiye give names ending with letter "-d" to their children.

An interesting issue with the spelling changes in the Turkish language is the urban legend about the removing of letters with a circumflex (^). In particular, some letters like "-a", "-ı", and "-u" are written with a circumflex in some words, e.g., *kâğıt* (paper), *abidevî* (monumental), *şûra* (council). While these letters are not removed from the official alphabet, many people on social media platforms claim that it was first removed but brought back later on. Even fact-checking websites had to verify the veracity of the claim[10]. In our corpus, we also focus on this urban legend and count the number of letters with a circumflex. The results are shown in **Figure 11**. We observe that letters with a circumflex have varying frequencies over time but they are continued to be used. However, we also notice that their frequencies have significantly dropped, which might be the reason to have such an urban legend.

# 6 Limitations

While our work provides valuable insights into how the Turkish language has changed since the 1920s, there are particular issues that need to be taken into account when analyzing our results. Our corpus mainly represents the language used by authorities on topics about governing. Therefore, it does not represent the whole characteristics of the Turkish language. That said, one of the main reasons for the relatively rapid evolution of the Turkish language is the government's policies aimed at language reform such as proposing new words to replace Arabic or Persian origin words and changing the alphabet. Therefore, our corpus might be one of the best resources to investigate the intervention of government in the Turkish language reform.

In our study, we use tools to automatize text extraction from the PDF files and detect lemmas of words. The tools we use might introduce noise and affect our results. Especially, in the diachronic analysis of embeddings, we aligned the consecutive time periods and measured the semantic change. We see that the most semantically displaced words are the noise words that survived throughout time periods. Also, in the word similarity task, we see that a set of k-similar words to a word contains OCR errors, and the errors follow some specific patterns, e.g. *-ü* has been scanned as *-ii* or *-e* has been taken as *-c*. Such errors mislead the embedding models in the training phase.. Therefore, we take action to reduce noise and its impact. To further mitigate this problem, we also share our code and dataset, increasing the reproducibility of our findings and enabling further research on this dataset. Nevertheless, the possible impact of the tools we use should be taken into account when analyzing our findings.
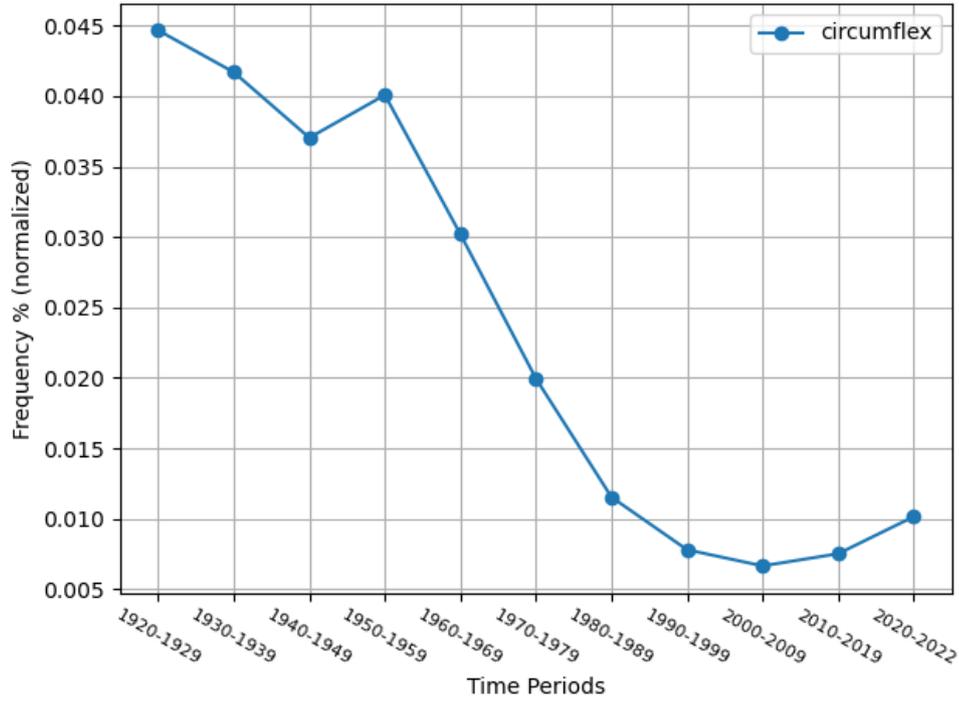
---

[10]https://www.malumatfurus.org/sapka-isaretinin-kaldirildigi-iddiasi/

**Fig. 11** The frequency of the usage of the circumflex in our corpus.

# 7 Conclusion

The Turkish language has encountered multifaceted transformation over the last century, underscored by state-driven initiatives such as changing the alphabet and replacement of loan-words with Turkish-origin words. In order to enable future studies on this interesting linguistic transformation of Turkish, we introduce Turkronicles; a toolkit that comprises various types of Turkish diachronic resources such as raw text corpus extracted from the Official Gazette of Türkiye and the records of Grand National Assembly of Türkiye over a century-long, diachronically aligned embeddings of different kinds, collocation matrices, a digitized dictionary of modern-old correspondent Turkish words, and Python library that allows diachronic analysis. Next, we conduct a comprehensive diachronic analysis using our corpus to investigate language reform in Turkish. In particular, we first explore how the vocabulary has changed since 1920. Next, we investigate how the spellings of words have changed.

Based on our comprehensive analysis, our findings are as follows. i) The vocabulary has dramatically changed throughout the years such that almost half of the words used in the 1920s were not used in the 2010s. We observe that the frequency of loanwords decreases while the frequency of words used for replacement increases throughout the years. Regarding the changes in spelling, our analysis reveals a noticeable decline in the use of circumflex compared to the 1920s and 1930s. Furthermore, there has been a shift in the final letters of

22

words over time, influenced by Turkish phonology. Specifically, words ending in "-b" have decreased over time in favor of versions ending in "-p". However, a distinct pattern emerges for words ending with "-d" or "-t" letters: The proportion of words ending with the letter "-d" compared to those ending with the letter "-t" remains similar to the proportions seen in the 1920s, although there has been a declining trend since the 1990s.

We think that Turkronicles paves the way for targeted studies on specific linguistic phenomena in Turkish, such as the evolution of certain grammatical structures, lexical borrowing, or semantic drift over time. While the introduction of a diachronic corpus for Turkish fills a significant gap in linguistic research, we plan to extend our corpus by the inclusion of texts from other sources, such as newspapers, literary works, and public broadcasts from corresponding periods in the future. Once we build a larger diachronic corpus from various sources, we plan to extend our analysis on Turkish language reform and compare the differences across different data sources. In addition, we plan to develop a software enabling easy access and analyze of the corpus for researchers.

# References

Johanson, L.: Turkic languages, (1998)

Lees, R.B.: The phonology of modern standard turkish. (No Title) (1961)

Lewis, G.: Turkish Grammar. Oxford University Press, ??? (2000)

Underhill, R.: Turkish Grammar. MIT press Cambridge, MA, ??? (1976)

Kornfilt, J.: Turkish. Routledge, London & New York (1997)

Göksel, A., Kerslake, C.: Turkish: A Comprehensive Grammar. Routledge, ??? (2004)

Say, B., Zeyrek, D., Oflazer, K., Özge, U.: Development of a corpus and a treebank for present-day written turkish. In: Proceedings of the Eleventh International Conference of Turkish Linguistics, pp. 183–192 (2002). Eastern Mediterranean University

Aksan, Y., Aksan, M., Koltuksuz, A., Sezer, T., Mersinli, Ü., Demirhan, U.U., Yilmazer, H., Atasoy, G., Öz, S., Yildiz, I., *et al.*: Construction of the turkish national corpus (tnc). In: LREC, pp. 3223–3227 (2012)

Sak, H., Güngör, T., Saraçlar, M.: Turkish language resources: Morphological parser, morphological disambiguator and web corpus. In: Advances in Natural Language Processing: 6th International Conference, GoTAL 2008 Gothenburg, Sweden, August 25-27, 2008 Proceedings, pp. 417–427 (2008). Springer

Sulubacak, U., Eryiğit, G., Pamay, T.: Imst: A revisited turkish dependency treebank. In: Proceedings of TurCLing 2016, the 1st International Conference on Turkic Computational Linguistics (2016). Ege University Press

Güngör, O., Tiftikci, M., Sönmez, Ç.: A corpus of grand national assembly of turkish parliament's transcripts. In: Proceedings of the Eleventh International Conference on Language Resources and Evaluation (2018)

Michel, J.-B., Shen, Y.K., Aiden, A.P., Veres, A., Gray, M.K., Team, G.B., Pickett, J.P., Hoiberg, D., Clancy, D., Norvig, P., *et al.*: Quantitative analysis of culture using millions of digitized books. science **331**(6014), 176–182 (2011)

Lieberman, E., Michel, J.-B., Jackson, J., Tang, T., Nowak, M.A.: Quantifying the evolutionary dynamics of language. Nature **449**(7163), 713–716 (2007)

Pechenick, E.A., Danforth, C.M., Dodds, P.S.: Characterizing the google books corpus: Strong limits to inferences of socio-cultural and linguistic evolution. PloS one **10**(10), 0137041 (2015)

Hamilton, W.L., Leskovec, J., Jurafsky, D.: Diachronic word embeddings reveal statistical laws of semantic change. arXiv preprint arXiv:1605.09096 (2016)

Szymanski, T.: Temporal word analogies: Identifying lexical replacement with diachronic word embeddings. In: Proceedings of the 55th Annual Meeting of the Association for Computational Linguistics (volume 2: Short Papers), pp. 448–453 (2017)

Xu, Y., Kemp, C.: A computational evaluation of two laws of semantic change. In: CogSci (2015)

Salan, M., KABADAYI, O.: Çağdaş türk yazı dillerinde art zamanlı söz başı ünlü düşmesi üzerine. Türk Dili Araştırmaları Yıllığı-Belleten (73 (Haziran)), 85–107 (2022)

Vahit, T.: Türkçede ön seste ünlü düşmesi örnekleri. Türk Dili Araştırmaları Yıllığı-Belleten **49**(2001), 223–233 (2003)

Sultanzade, V.: The syntactic valency of some verbs in the book of dede korkut: diachronic differences. bilig **61**, 223 (2012)

Aksan, D.: Türk anlam bilimine giriş–anlam değişmleri. Türk Dili Araştırmaları Yıllığı-Belleten **13**, 167–184 (1965)

Bahattin, S.: Anlam değişmeleri üzerine artzamanlı bir inceleme. Gazi Üniversitesi Gazi Eğitim Fakültesi Dergisi **23**(1) (2003)

Öztürel, A., Kayadelen, T., Demirsahin, I.: A syntactically expressive morphological analyzer for turkish. In: Proceedings of the 14th International Conference on Finite-state Methods and Natural Language Processing, pp. 65–75 (2019)

Can, F., Kocberber, S., Balcik, E., Kaynak, C., Ocalan, H.C., Vursavas, O.M.: First large-scale information retrieval experiments on turkish texts. In: Proceedings of the 29th Annual International ACM SIGIR Conference on Research and Development in Information Retrieval, pp. 627–628 (2006)

TDK: Türkçeden Osmanlıcaya Cep Kılavuzu. Turkish Language Association, Ankara (1935)

Mikolov, T., Chen, K., Corrado, G., Dean, J.: Efficient estimation of word representations in vector space. arXiv preprint arXiv:1301.3781 (2013)

Levy, O., Goldberg, Y., Dagan, I.: Improving distributional similarity with lessons learned from word embeddings. Transactions of the association for computational linguistics **3**, 211–225 (2015)

Rehurek, R., Sojka, P.: Software framework for topic modelling with large corpora (2010)

Schönemann, P.H.: A generalized solution of the orthogonal procrustes problem. Psychometrika **31**(1), 1–10 (1966)