UNIVERSITY OF TARTU

Institute of Computer Science

Computer Science Curriculum

Lauri Lüüsi

# Political Stance Detection
# in Estonian News Media

**Bachelor's Thesis (9 ECTS)**

Supervisors:  Uku Kangur, MSc

Roshni Chakraborty, PhD

Tartu 2024

# Political Stance Detection in Estonian News Media

**Abstract:**

As online news media continues to grow, it is becoming increasingly important to check that the way information is presented is fair and unbiased. The goal of this thesis is to explore and identify features that are indicative of political stance in Estonian news media. This is done by describing the relevance of political stance detection and outlining the existing works and approaches in this field. Additionally, an exhaustive analysis of the existing features used for stance detection in English is conducted, and novel features specifically for Estonian news media are proposed.

# Poliitiliste hoiakute tuvastamine Eesti uudistekstidest

**Lühikokkuvõte:**

Veebiuudiste hulk on kasvamas, mistõttu on üha olulisem kontrollida, et meedias avaldatud teave oleks esitatud õiglaselt ja kallutamata. Käesoleva töö eesmärk on katsetada ja tuvastada tunnuseid, mis osutavad poliitilistele hoiakutele Eesti uudistekstides. Töös kirjeldatakse poliitiliste hoiakute tuvastamise tähtsust ja tutvustatakse varasemaid uurimistöid ning lähenemisviise. Analüüsitakse ka ingliskeelsetele tekstidele kohandatuid tunnuseid ning pakutakse välja uusi tunnuseid poliitiliste hoiakute tuvastamiseks Eesti uudistekstidest.

# Contents

# 1.    Introduction

Humans are naturally curious and constantly seek knowledge about the world around them. News can spread in multiple ways, from word-of-mouth conversations to television broadcasts. However, with the rise of the internet, this information-seeking behavior has undergone a shift. News media has pivoted away from traditional printed newspapers and moved towards social media and online platforms [1, 2]. In today's modern and interconnected world, online news coverage has become increasingly diverse and widespread, with articles becoming readily available within minutes for almost any notable event or topic [1, 2].

News can diverge from being honest and of high quality in multiple ways, such as by deliberately lying or leaving out context, not fact-checking sources, using clickbait, or being biased [3]. Readers place a great amount of trust into the media since they assume the news to be truthful and accurate, not speculative and dishonest [4]. However, with the massive growth in news media, it has become increasingly challenging and time-consuming to manually verify that the content produced by these outlets follows journalistic standards. Therefore, it is important to develop and explore automated techniques as a way to gain insight and monitor news media.

In English, automated stance detection concerning political leaning has been explored for different topics, such as political elections and candidates, climate change, and abortion rights [5, 6, 7]. However, political stance detection in Estonian is a mostly unexplored topic. Estonian is a small language with around a million speakers [8]. In contrast, English has over 1.4 billion speakers [9]. According to Kadri Vare, the head of EKI's[1] language and speech technology department, the main Estonian language corpora contains approximately 3 billion words, a relatively small amount compared to 800 billion words available for English [10]. This divide can also be noted in the amount of data available for natural language processing tasks, such as training large language models [10, 11].

---

[1] Institute of the Estonian Language (Eesti Keele Instituut)

The rise of online news media in Estonia is notable. Despite a population of around 1.35 million people [12], online news media is growing. Delfi Meedia, one of the largest media companies in Estonia, claims to have over 700 000 readers per month [13], and as of 2023, the platform has accumulated over 100 000 online paid subscribers [14]. Therefore, with this increasing popularity, manually validating news articles is highly challenging and the development of automated approaches for political stance detection is required for Estonian news media.

However, developing automated approaches for low-resource languages can be challenging, as these smaller languages are particularly affected by the non-availability of task and domain-specific data [11]. Furthermore, identifying labeled data requires manual annotation, which is time and cost intensive [11]. Compared to English, the data can be of lower quality, which can lead to poorer results and varied performance [11]. Therefore, identifying political stance in Estonian news media reporting is particularly challenging.

In this thesis, political stance in Estonian news media is analyzed on the topic of immigration. The goal of the thesis is to explore and identify features and techniques that are indicative of political stance in Estonian news media. This is done by an exhaustive analysis of the existing relevant features used for stance detection in English and the proposal of novel features specifically for Estonian news media. Political stance detection is done for Estonian news media, as few related works exist on the topic of automated stance detection in this language. A political stance dataset consisting of 3261 sentences will be analyzed, from which different features are extracted, such Estonian-specific features like diminutives and words in the translative case.

This thesis consists of 7 main chapters. The first chapter introduces the topic and defines the research goals. The second chapter describes the background and motivation and explains the meaning of stance, political stance, and stance detection, followed by an overview of the prior works. The fourth chapter gives an overview of the dataset and describes the preprocessing steps. The fifth chapter describes the methodology and features extracted for the task of automated political stance detection. The sixth chapter contains the exhaustive analysis of extracted features and their usefulness in political stance detection. The seventh chapter summarizes the thesis and mentions ways in which this work could be expanded or improved upon.

## 2.    Background

This section outlines the concept of stance by introducing its definition and detailing the stance detection procedure. Additionally, motivation for conducting this work is provided to explain the novelty and necessity of research in this domain.

### 2.1    Definition

Stance can be thought of as a way of thinking [15], where an individual is either against, neutral, or supportive towards a topic or idea [16, 17].

Du Bois [18] identified that taking a stance involves the three following distinct steps:

1.  Evaluation of an object, usually the topic or idea under discussion;

2.  Positioning oneself, when the stancetaker clarifies their perspective or attitude;

3.  Alignment with others, during which the stance is contextualized, which enables it to be categorized as either in favor (positive), neither (neutral), or against (negative).

Stance detection refers to the process of automatically determining or predicting an individual's position or perspective on a particular topic. By detecting stance from written content, it is being assessed whether the author is against, neutral, or in favor of a particular proposition or target [19, 20].

As a natural language processing task, stance detection can be thought of as a classification problem, where a pair consisting of text and a target is assigned a label that is either positive, negative, or neutral. [19, 20]. For stance detection, the target can encompass any idea or object [18]. Political stance emerges when the target or subject is related to political matters, such as climate change, abortion, or specific politicians. In this thesis, the target is the topic of immigration, and the text is a topic-related sentence. Immigration is a concept that encompasses the international movement of people, usually foreign nationals, who intend to become permanent residents in a different country [21, 22].

Immigration is a suitable target for automated stance detection, as stances towards it are varied and can often veer towards extremes [23, 25]. Immigrants can be viewed as strong and talented workers with great potential or, conversely, burdens on society who take jobs from locals and will not integrate into the local culture [23]. Media coverage of immigration can influence public opinion, especially when it adopts an overly negative stance [23]. These shifts in attitude can potentially translate to negative treatment of immigrants, fueling racism and social division, and the enactment of discriminatory policies [23, 24, 25].

## 2.2    Motivation

Limited research has been conducted on political stance detection in Estonian. The few existing works do not directly address stance detection and are more analytical and context-dependent in their approach rather than automatic and generalizable. The absence of works in this area highlights the importance of exploring the topic of stance detection in Estonian media.

A lack of suitable natural language processing tools is also apparent. In her master's thesis, Natalja Maksimova [26] used machine learning methods to predict and analyze the reasons behind computer science students dropping out early. For this task, essay texts provided by the students were analyzed [26]. However, these essays, some initially written in Estonian, were translated into English to fully utilize all of the features of quanteda, a package in the R programming language used for natural language processing tasks, namely quantitative analysis of textual data [26, 27]. This package does not directly support Estonian natural language processing tasks [28].

# 3.    Related Works

In this chapter, related works are outlined to provide context and insights about existing research and approaches for automated political stance detection. This is done for studies regarding English news media, as well as for works in smaller non-English languages, such as Estonian.

## 3.1    Political Stance Detection in English

Bias refers to the tendency or inclination to show preference or discrimination towards an individual or group based on their inherent or acquired characteristics, potentially leading to unequal treatment or outcomes [29]. Stances, opinions, and biases occur naturally in human discourse, which have discriminatory effects that have been outlined by Mehrabi *et al.* [29], such as content production bias, which arises from structural, lexical, semantic, and syntactic differences in the contents generated by users, or behavioral bias, which is caused by the way users behave on different contexts. Avoiding unfair treatment and bias has become crucial as artificial intelligence applications become more widespread in daily life, especially in sensitive environments, such as those handling medical or immigration data [29].

Political bias in news media has been explored by Spinde *et al.* [30], who revealed a direct correlation between perceived journalist bias and political extremeness. Sen *et al.* [31] analyzed media bias in Indian policy discourse and found that Indian mass media exhibits coverage bias, as it extensively prefers to cover middle-class concerns instead of issues regarding the poor. Political bias was also found by detecting political stance at the statement level, as specific news sources were found to be affiliated with or against the two largest political parties in the country [31].

Stance detection can be approached in many different ways. ALDAyel *et al.* [5] outline possible stance detection techniques and features, such as bias and sentiment lexicons, n-grams, word embeddings, and topic modeling, among others. Spinde *et al.* [32] attempt to identify bias-inducing words from news articles by using linguistic and context-oriented features, such as POS (part-of-speech) tags, syntactic dependencies, named entities, word embeddings, sentiment-related features, and word lexicons. For this task, a dataset of 1700 sentences was

created, where sentences were annotated as biased or non-biased by 784 people of varying backgrounds [32]. An F1-score of 0.43 was achieved to detect biased words, revealing that feature-based approaches can be used for media bias detection [32]. Textual features, such as n-grams and word embeddings were used by Mohammad *et al.* [33] to detect stance in tweets, achieving an average F1-score of 0.69 for an n-gram based SVM (support vector model) classifier.

Guo *et al.* [34] approach media bias detection by fine-tuning BERT models for different media outlets. By fine-tuning a model for each outlet (such as CNN, Fox, Breitbart), it inhibits that specific outlet's innate biases and attitudes, which arise when the model performs prediction tasks, such as prompt-based mask token prediction [34]. The model is asked to finish a sentence, and by doing so, it may reveal hidden opinions or biases that were learned from the articles used for fine-tuning [34]. This makes it possible to generate and analyze attitudinal represent-tations, which can reveal insights about the relative biases present in each outlet [34].

While these approaches have demonstrated applicability in English contexts, their direct transferability to Estonian needs further exploration. Directly applying these approaches could lead to faulty or misleading results, as language-specific nuances might introduce unexpected variables. To better detect the stances expressed in Estonian texts, considering and extracting additional language-specific features is likely beneficial.

## 3.2    Political Stance Detection in Non-English Languages

Political stance detection has been explored in other languages, such as Danish and Swedish [35, 36]. Political stance detection in Danish is attempted in a paper by Lehmann and Derczynski [35], in which word embeddings along with context-based features were used to train three classifiers. An annotated dataset was constructed, consisting of 898 quotes from Danish politicians on the topic of immigration policy. The study revealed that the simpler MLP (multilayer perception) model outperformed the more advanced LSTM (long short-term memory) models. The inclusion of context-based features, such as political party affiliation, significantly enhanced model performance. For the best-performing model, the F-score decreased from 0.72 to 0.58 when context features were excluded.

Yantseva *et al.* [36] conducted political stance classification in Swedish using various machine learning methods, such as support vector machines (SVM), logistic regression, extreme gradient boosting (XGBoost), among others. A dataset of 5701 immigration-related messages sourced from a Swedish online forum were transformed into a vector representation. This data was used to train a model that achieved an F1-score of 0.72.

Stance detection is an area of research relatively unexplored in Estonian, as only a few papers explore this topic. However, some existing works cover similar research areas, such as fairness in news and sentiment analysis of texts.

Mets *et al.* [37] detect political stance from Estonian immigration-related sentences using BERT-like large language models and attempt zero-shot classification using ChatGPT. A dataset of 8000 sentences was extracted from two Estonian media outlets, Uued Uudised and Eesti Ekspress. Est-RoBERTa, a monolingual Estonian language model, performed the best, achieving an average F1-score of 0.66 [37, 38].

In his thesis, Pärt Dolenko [39] conducted a case study on issue and framing bias. Editorials from two large Estonian media providers, Eesti Päevaleht and Postimees, were analyzed during a six month period when the socially conservative government (referred to as EKREIKE[2]) was replaced by a government lead by Kaja Kallas, the leader of the Reform party. Deviations from media neutrality were present for both outlets, which were noted to prefer covering certain policy issues. In addition, both outlets' editorials expressed a negative stance toward the resigning government and a comparatively more positive stance toward the incoming government. However, the overall stance when covering the government or politicians was found to be negative.

Framing analysis has been conducted for Estonian media texts, particularly regarding feminism and women's rights issues. Raili Marling [40] analyzes a textual corpus from the Estonian news outlet Postimees and observes that feminism is framed rather negatively, as adjectives such as *rõve* (obscene), *agressiivne* (aggressive), and *räige* (radical) are used to describe it. Kaidi-Lisa Kivisalu [41] conducts a frame analysis of the #MeToo movement in Estonia, in which 112

---

[2] EKREIKE is a portmanteau of the acronyms that refer to coalition parties in the second cabinet of Jüri Ratas. EKRE (Eesti Konservatiivne Rahvaerakond – Estonian Conservative People's Party), I (Isamaa – Fatherland), KE (Keskerakond – Estonian Centre Party).

opinion pieces from news outlets Postimees, EPL[3], and ERR[4] were found to be rather positive towards the movement, although victim-blaming was found to be the most prevalent frame in the discourse. Elisabeth Kaukonen [42] analyzes the use of gender-marked words (such as compound words either ending with *man* or *woman*) in sports-related news articles. She finds that *man*-suffixed words in news focusing on sports are more prevalent, as they make up 98% of compound word occurrences.

In her bachelor's thesis [43], Grete-Liina Roosve analyzes the interventionist role of the journalist in news articles. She identifies features indicative of the interventionist role, such as using the first person, excessive use of adjectives, and inserting claims or suggestions. A corpus of 2409 articles is analyzed, from which it was found that intervention is present in approximately 40% of news stories, and an opinion is expressed in 13% of articles.

While these approaches provide some useful features and reveal insights about political stance and biases in Estonian news media, they are rather manual and qualitative as they tend to focus on specific issues or contexts in certain time periods. This thesis approaches political stance detection from a computational standpoint and aims to achieve quantitative and generalizable results.

---

[3] Eesti Päevaleht
[4] Eesti Rahvusringhääling

# 4. Dataset

The dataset used in the thesis originates from a paper authored by Mark Mets [37], who acquired the data from two Estonian news providers - Ekspress Grupp (the parent company of Delfi Meedia) and Uued Uudised. The corpus consists of 266 628 articles published between 2015 and 2022.

To obtain a sentence-level dataset covering immigration, Mets *et al.* [37]. constructed a lexicon containing migration-related keywords and word stems. The lexicon was compared to the article corpus, where topical sentences with matching keywords were extracted. As the Estonian language is morphologically complex, regular expressions were also used to capture all possible case forms.

Following this approach, a set of 7345 sentences was extracted [37]. The sentences were annotated by two Estonian-speaking graduate students, who rated each sentence on a 1-5 scale [37]. Ratings 1-2 and 4-5 were indicative of an against or supportive stance towards immigration, respectively, and a rating of 3 deemed the sentence to be neutral [37].
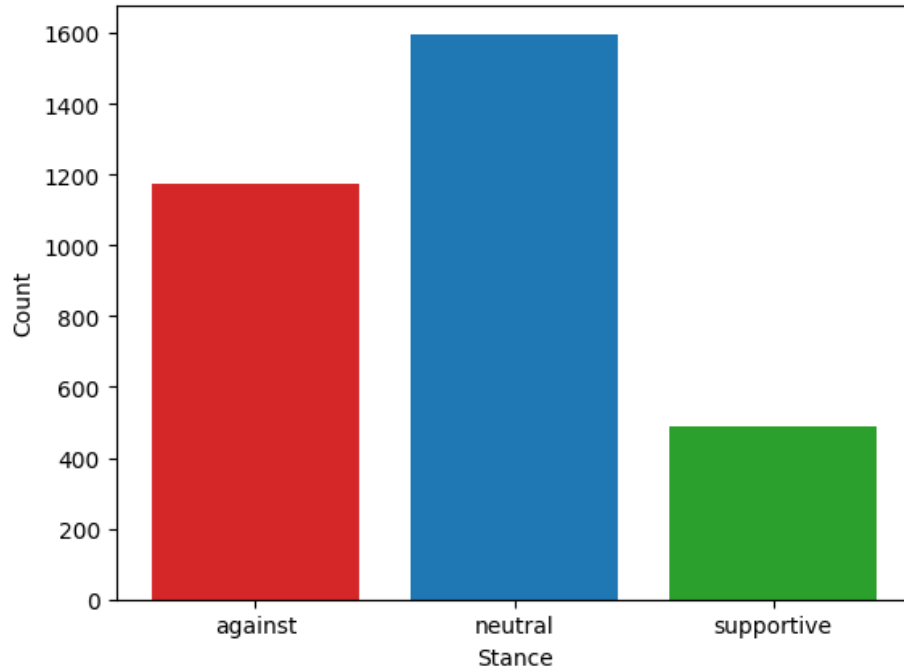


Figure 1. Sentence distribution across the three stance classes.

Additionally, an ambiguous rating was used to exclude sentences that included multiple stances, were unclear, or were not on topic [37]. For this thesis, all sentences annotated as ambiguous were disregarded, reducing the dataset to 3261 sentences, where 1175 were against, 1597 neutral, and 489 supportive towards immigration, as illustrated by Figure 1. The dataset used in the thesis is publicly accessible on GitHub[5] and can be freely downloaded as a CSV file (comma-separated-values file).

## 4.1   Used Tools

The Python programming language is widely used for data science and analysis tasks [44]. In this thesis, Python is used via the interactive graphical browser-based Jupyter notebook environment known as Google Colaboratory [45]. The notebook environment allows users to interlace their code with formatted text, equations, and dynamic visualizations [45]. As Google Colaboratory is a cloud-based environment, it does not require intricate setup from the user and enables the sharing and execution of notebooks by other users [44, 45].  The GitHub repository containing the notebook files for this thesis is linked in Appendix 1. Some additional third-party Python resources used for conducting the work in this thesis are described below.

**Pandas:** Pandas is a library that simplifies data analysis tasks in Python, as it introduces a dynamic tabular data structure known as a DataFrame [46]. In this thesis, the DataFrame rows represent instances of data (sentences), and columns represent values and the presence of different extracted features in that sentence.

**EstNLTK:** EstNLTK is a natural language processing toolkit developed by researchers at the University of Tartu [47]. EstNLTK enables the tokenization and morphological analysis of texts, in part through its *tag_layer* method, which adds different layers of information to the text, such as lemmas, case forms, or part-of-speech tags [47]. This is useful for preprocessing and extracting the Estonian-specific features in this thesis.

**Scikit-learn:** Scikit-learn is a machine-learning library that offers tools for data analysis and implementations of many machine-learning algorithms [48]. Additionally, Scikit-learn provides

---

[5] https://github.com/markmets/immigration-prediction-EST/blob/main/Annotated_Dataset.csv

automated calculation of evaluation metrics that are used in this thesis to judge the performance of sentiment models.

**Matplotlib:** Matplotlib is a Python library for creating visualizations, such as graphs, plots, and charts [49]. In this thesis, it is used to give a visual overview of trends and distributions in data.

**Transformers:** The Transformers is an open-source Python library for sharing pre-trained machine learning models [50]. Through the associated web platform, Hugging Face, users can upload their models for others to access [50]. In this thesis, existing pre-trained models uploaded online are used to conduct sentiment analysis.

## 4.2    Preprocessing

Before further analysis was conducted, the data was preprocessed to remove noise and aid the performance and accuracy of automated tools to ensure effective identification of features.

Mets *et al.* [37] posed the sentence as the unit of analysis instead of paragraphs or articles, which have varying lengths between instances and may contain more than one stance. Therefore, it was assumed that each row contained one sentence. However, EstNLTK's morphological analysis revealed that several rows in the dataset contained more than one sentence. This inconsistency must be accounted for, as automated tools may produce unintended or faulty results, given that some features, like dependency tree height, only work under the assumption of a single-sentence structure. Taking this into consideration, the total count of a frequency-based feature was divided by the respective number of sentences in a given row.

The Estonian language has two letters with diacritics (š, ž). Sentences in the dataset containing words with these letters were replaced by question marks or other nonsensical symbols, leading the morphological analysis tool to interpret them as separate sentences. This was fixed using EstNLTK's SpellCheckRetagger[6], a tool that identifies misspellings and adds corrected forms [47]. The broken symbol within the word was replaced by either *š* or *ž*, checked by the tagger,

---

6 https://github.com/estnltk/estnltk/blob/5c5ce3a810f7a5e41602156b7044edb63e83532d/tutorials/nlp_pipeline/B_03_segmentation_words_spelling_normalization.ipynb

and subsequently corrected if it was identified as the correct spelling. This approach did not correct spelling for words (mostly named entities) that contained letters with diacritics other than *š* and *ž* since the Estonian spellcheck tagger could not provide analysis for named entities, such as foreign place- or surnames.

Spacing issues, repeated symbols, or missing punctuation at the beginnings and ends of sentences caused additional issues, due to which some single-sentence rows were incorrectly interpreted as having multiple sentences, or some multi-sentence rows were considered to have only one sentence. Redundant symbols were removed, and additional spaces and punctuation were added to fix these problems.

In large numbers, spaces are inserted to separate the thousands to enhance readability (200000 → 200 000). However, these were removed due to causing issues with dependency tree parsing.

# 5.    Methodology

A news article is a non-fictional piece of writing that informs readers of recent events [51, 52]. In this section, different features that can help in identifying the political stance of Estonian news media are studied. These are segregated into lexical features, features specific to the Estonian language, framing-related features, and sentiment features, which will be discussed in the following subchapters.

Features are extracted solely from the text of the news article. While additional meta-features, such as the title, author, publication date, and publisher, are available, they are not currently considered as they require prior outside knowledge about a media outlet or author and their stance on specific issues.

## 5.1    Lexical Features

Lexical features are related to the grammar and construction of words [53]. These features are extracted from the article's contents using natural language processing techniques. The following lexical features are covered in this thesis: word count, dependency tree height, Flesch Reading Ease Score, named entities, noun phrases, adjectives, quotes, and quoted phrases.

### 5.1.1  Sentence Complexity

The structure of a sentence influences its readability. A short and simple sentence can be read easily and quickly, whereas a lengthy and intricate one demands more attention and focus. By analyzing features related to the complexity of a sentence, it can be seen how a particular stance is expressed. Anti- or pro-immigration sentences might be deliberately complex to confuse or demand more attention from the reader, or conversely, easy and quick to understand so that the author could convey and spread their stance as clearly and to as many people as possible.

In this subchapter, three complexity-related features are described: word count (the number of words in a sentence), dependency tree height, and Flesch Reading Ease Score. The latter two are discussed in detail.

**Dependency Analysis:** Dependency trees represent the grammatical structure of a sentence. Words are connected by labeled edges that indicate syntactic relationships between subjects, words, and objects, among others [54, 55]. The height of the dependency tree corresponds to the distance from the root node, typically a verb, to the farthest branch [55]. A taller tree indicates longer dependency paths that allude to a more complex sentence structure [55].

For Estonian dependency tree parsing, EstNLTK provides models trained on the Estonian Dependency Treebank[7], a dataset containing 30 000 syntactically annotated sentences.

**Readability Analysis:** Readability scores are used to quantify and measure a text's clarity and comprehensibility [56]. A standard score used to evaluate readability is referred to as Flesch Reading Ease Score (FRES), which is defined by the formula

$$\text{FRES} = 206.835 - (1.015 \cdot \text{ASL}) - 84.6 \cdot \text{ASW} \ , \tag{1}$$

where ASL represents the average number of words per sentence, and ASW represents the average number of syllables per word [56].

The formula yields a score on a scale of 0-100, where a lower score indicates a more complex text [56]. Texts with scores between 60-100 are considered easy to read for individuals with basic education, as opposed to texts with scores between 0-30, which may be challenging even for highly educated individuals [56]. For calculating FRES, EstNLTK can be used, as the toolkit supports syllabifying Estonian texts and provides a tagger[8] for calculating the complexity score.

### 5.1.2 Named Entities

In text, named entities are phrases that typically encompass people, geographic locations, organizations, or any other entities that can be referred to with a proper name [57]. By detecting and analyzing named entities, it can be seen whether some individuals or organizations are central to a certain stance.

---

For named entity recognition, EstNLTK's NerTagger[9] is used. The tagger was developed by Tkachenko *et al.* [58] and further improved by Rasmus Maide in their bachelor's thesis [59]. A machine learning algorithm known as conditional random fields (CRF) was used for training on a labelled dataset containing orthographical, morphological, WordNet, and dictionary-based features [58]. Despite Estonian being a morphologically complex language, the model achieved an F1-score of 0.87 [58].

### 5.1.3 Noun Phrases

A noun phrase is a phrase that consists of a main (head) noun and its modifiers. In the Estonian language, modifiers can be adjectives (*suur maja – big house*), possessive nouns (*minu maja – my house*), and prepositional phrases (*maja kõrval – next to the house*) [60]. For extracting noun phrases, EstNLTK's experimental noun phrase chunker[10] is used. It detects non-overlapping noun phrases from text by combining consecutive words that share a dependency relation, therefore being syntactically connected.

### 5.1.4 Adjectives

Adjectives are words that modify and describe the qualities or attributes of nouns [61]. Adjectives provide valuable insight about the stance expressed within sentences. A high concentration of adjectives in a sentence may suggest that the writer is emotional and might have a stance on a specific topic. Adjectives can also carry positive or negative sentiments, the use of which may suggest a favorable or unfavorable stance on immigration. Adjectives were extracted using EstNLTK's morphological analysis, which included part-of-speech (POS) tags for words [47].

### 5.1.5 Quotes

Quotes are punctuation marks typically used to mark direct speech, citations, or titles. However, quotes can also be used to convey irony and sarcasm or detachment and skepticism [62, 63].

---

[9] https://github.com/estnltk/estnltk/blob/main/tutorials/nlp_pipeline/D_information_extraction/02_named_entities.ipynb
[10] https://github.com/estnltk/estnltk/blob/main/estnltk/estnltk/taggers/miscellaneous/np_chunker_v1_4_1.py

Therefore, in certain instances, the use of quotes might signal a particular stance, with the author saying one thing while intending to convey something else [62, 64]. This is illustrated by the sentences in Example 1.

| (a) She did a good job. | (b) She did a "good" job. |
|---|---|

Example 1. Alternate meaning evoked by quotation marks.

When the word *good* is not between quotes, the speaker genuinely implies that the job was well done. However, in sentence B, the use of quotes suggests skepticism or doubt about the quality of the job. The usage of quotes can imply a contrary or ironic meaning of the literal expression [62-65]. A similar usage of quotes is also apparent in some of the sentences in the dataset, as shown in Example 2.

| [---] aastaga tuli Eestisse 22 000 **"ajutist"** töölist [---] |
|---|
| [---] in a year, 22 000 **"temporary"** workers came to Estonia [---] |

Example 2. A sentence from the dataset that uses quotation marks to express doubt.

In Example 2, the presence of quotes for the word *temporary* questions the duration of the employment status, suggesting that the workers are staying for longer or seeking residency. In contrast, a lack of quotes would imply that the workers are genuinely employed on a temporary basis.

| [---] võttes vastu inimsmugeldajate **"ohvreid"** aafriklaste ja teiste migrantide näol. |
|---|
| [---] by accepting **"victims"** of human traffickers in the form of Africans and other migrants. |

Example 3. A sentence from the dataset that uses quotation marks to express doubt.

In Example 3, removing the quotes around the word *victims* would convey that the individuals are indeed victims of human trafficking without casting doubt on the validity of their victim status.

> **"Sallivuslased"** aitavad neil oma soovituste ja muu **"abiga"** ennast hädalistena tutvustada ja mõrvarid seavad ennast **"pagulastena"** Euroopas sisse.

> **"Tolerance advocates"** help them with suggestions and other **"assistance"** so they could present themselves as sufferers, as murderers establish themselves as **"refugees"** in Europe.

Example 4. A sentence from the dataset that uses quotation marks to express irony.

Example 4 showcases irony. The use of quotation marks around the given terms implies skepticism and disbelief. The author is not sincere and questions the motives and identities of those described as *tolerance advocates* or *refugees*.

Quoted phrases and words were detected from the dataset by using regular expressions to split the sentence based on multiple variations of quotation marks (" " ",, ', «»). To account for titles and citations, quoted phrases starting with an uppercase letter and ending with a punctuation mark were skipped. In addition, phrases longer than 5 words were skipped, as they were likely to be direct speech.

## 5.2    Estonian-specific Features

This subchapter describes features that could be indicative of stance and are specific to the Estonian language. Extracting these features from English texts can be difficult due to the differences between the two languages.

The Estonian language is morphologically complex [37, 66, 67]. This complexity is partly exemplified by the abundance of verb conjugation forms and grammatical cases for nouns and adjectives [66, 67]. These characteristics can make analyzing Estonian texts challenging but, conversely, provide helpful features for stance detection tasks, as specific forms of words could be more prevalent in a certain stance.

The following Estonian-specific features are discussed in this subchapter: diminutives, superlatives, conditional form, translative case. indirect speech.

### 5.2.1 Diminutives

In Estonian, diminutive words are formed by adding the suffix *-ke* or *-kene*. The usage of these suffixes changes the emotional tone of the word, which can either shift to being positive (affectionate) or negative (derogatory and belittling) [68]. The diminutive suffix can be added to both nouns and adjectives. In the case of adjectives, the suffix diminishes the characteristic indicated by the adjective (*loll* – stupid, *lollike* – dummy) [68, 69]. English does not have a consistent suffix for diminutive words, unlike Estonian where forming the diminutive is mostly uniform across nouns and adjectives.

The use of diminutives in Estonian spoken language is analyzed by Mirjam Liivak [68] in her master's thesis. Out of 143 instances of diminutives, 43 were used to express a positive sentiment, and 27 were used to express a negative sentiment. To better illustrate the use of the diminutive form, consider Example 5.

| |
|---|
| Aga selleks ju migrandipaadid **kehvakesed** ongi, ja ilmselt lastakse need mõnda laeva märgates meelega vett täis. |
| But that's exactly why migrant boats are so **flimsy**, and presumably they are intentionally filled with water when spotted by a ship. |

Example 5. Use of diminutives in a sentence from the dataset to express a stance.

There is a subtle shift in stance from the word *kehv* (poor, bad) to *kehvakene* (flimsy). The usage of the diminutive form suggests that the author does not consider the poor construction quality of the boat to be a cause of concern and suggests that it is deliberate rather than resulting from a lack of adequate resources.

Diminutives were extracted by checking whether the suffix was present in the word's root form, as provided by EstNLTK. The extraction was more difficult for adjectives since the root form did not differentiate the suffix. Additionally, some adjectives, such as *raske* (difficult) and *väike* (small) with *-ke*, but are not diminutive forms.

### 5.2.2 Superlative Form

The superlative is a grammatical form for adjectives that indicates that the subject under discussion possesses some characteristic to a greater extent than any other subject of the same category [63]. In Estonian, the superlative is usually denoted by the suffix *-im* (*suurim* – biggest) or by the word *kõige* preceding the comparative form (*kõige kiirem* – fastest) [63].

| |
|---|
| **Suurim** probleem ongi see, et kogu Euroopa on sunnitud migrantidega tegelema [---] |
| The **biggest** problem is that the whole of Europe is forced to deal with migrants [---] |

Example 6. Use of the superlative form in a sentence from the dataset to express a stance.

The use of the superlative form can convey extreme opinions or positions, which can indicate stance. This is illustrated in Example 6, where the author expresses an anti-immigration stance by using the superlative form, suggesting that immigration surpasses all other problems in magnitude or significance.

### 5.2.3 Conditional Form

In Estonian, verbs in the conditional form end with the suffix *-ks*. The use of the conditional form can imply implausibility or that the situation being described is unrealistic [63].

| |
|---|
| Massimigratsiooni mahitajad aga **ujutaksid** kontinendid pigem migrantidega üle ja **segaksid** ära kogu maailma rahvastiku. |
| The proponents of mass migration, however, would rather **flood** continents with migrants and **mix up** the entire world population. |
| Pigem läheks vastupidi, nooremad **tõttaks** multikultiriiki, kus Hiina rahadega **loodaks** näiline jõukus (kuni migrandimassid selle ära söövad). |
| On the contrary, younger individuals **would rather flock to** multicultural countries, where apparent prosperity **is hoped to be created** with Chinese money (until migrant masses consume it). |

Example 7. Use of the conditional form in sentences from the dataset to express a stance.

Example 7 showcases sentences from the dataset where an against stance is expressed. In both sentences from Example 7, the authors describe unrealistic and hyperbolic events by using the conditional form.

**Translative case:** Similarly, nouns in the translative case end also end with the suffix *-ks*. The translative case is often used to express change, as a state or condition into which one transitions (*Ma muutusin õnnelikuks*, I became happy) [63]. This is further explored by Kristina Pai [70] in her master's thesis, where a form of the translative case, known as the translative predicative adverbial is examined. Pai found that depending on the context, words in the translative case can suggest peculiarities or express attitudes (*Pean teda lolliks*, I think he's stupid) [70].

### 5.2.4  Indirect Speech

In Estonian, indirect speech expresses a statement heard from someone else rather than directly from the speaker [63]. It can be recognized by the use of verbs that end with the suffix *-vat*. According to Lee *et al.* [71], indirect speech enables plausible deniability and is often used to introduce uncertainty, deflect responsibility, and distance individuals from the information or opinion they are expressing.

| |
|---|
| [---] kõik nad on viimasel ajal kukkunud aktiivselt ladistama sellest, et Eesti **polevat** migratsiooni sihtriik, meil **polevat** mingit massiimmigratsiooni ega ka **tulevat** seda, ning rändepakt **olevat** vaid õnnistus. |
| [---] all of them have recently been chatting about how Estonia is **supposedly not** a destination country for migration, there is **supposedly no** mass immigration here, **nor is there going to be** any, and the migration pact **is said to be** nothing but a blessing. |

Example 8. Use of indirect speech in a sentence from the dataset to express a stance.

Example 8 shows how an anti-immigration stance is expressed by using indirect speech to convey statements heard from others. Phrases like *polevat* (supposedly) and *olevat* (said to be) express uncertainty, as the author is suggesting that the actual situation regarding immigration is different from what is being claimed.

## 5.3    Framing Analysis

This subchapter describes framing-related features. The way authors frame and present certain concepts can be indicative of stance [5, 40-42]. The following concepts and features are related to framing discussed: black-and-white thinking, bigram analysis, adjective-based framing.

### 5.3.1  Black-and-white Thinking

Black-and-white thinking is a logical fallacy in which a complex situation is simplified into two extremes [72]. When authors use extreme or polarizing language, they often eliminate or do not consider alternate perspectives or possibilities [72]. In this thesis, black-and-white thinking is detected by word choice.  Table 1 contains a list of hyperbolic words that could be considered polarizing. By detecting words from this list, it can be assessed whether a particular stance is being portrayed in a binary matter and lacks a middle ground.

Table 1. List of words to detect black-and-white thinking. English translation in parenthesis.

| |
| --- |
| kõik, kõige (all), kunagi, eales (ever), iial (never), alati (always), igavesti (forever), tervenisti, täiesti, üleni (entirely), täitsa, täielikult (completely), üdini, läbinisti (thoroughly), läbini (through and through), absoluutne (absolute), absoluutselt (absolutely), totaalne (total), totaalselt (totally), ainult (only), ainus (sole), kogu (whole) |

The basis of this list was constructed by hand and further expanded using the Estonian Synonym Dictionary. To detect black-and-white thinking, the words from this list were matched to the sentences in the dataset by counting their occurrences. A sentence with black-and-white thinking is showcased in Example 9.

| |
| --- |
| Vahemere paadipõgenike ümber toimuv jätab üha enam mulje, et rändekriis hakkab **kõigile** närvidele käima, välja arvatud inimõiguslased ja teised sallivuslased, kes ei muutu **kunagi**. |
| The events around Mediterranean boat refugees increasingly give the impression that the migration crisis is getting on **everyone's** nerves, except for human rights activists and other tolerant individuals who **never** change. |

Example 9. Black-and-white thinking in a sentence from the dataset that expresses a stance.

The words *kõigile* (everyone's) and *kunagi* (never) convey black-and-white thinking as they help express a polarized perspective. By using the word *everyone*, the author implicitly expresses that this belief is universal and adopted by all people, whereas the word *never* implies an absolute and non-changing stance. This binary word use disregards nuances and variations and conveys a complex situation rather plainly.

### 5.3.2 Bigram Analysis

An *n*-gram is a sequence of $n$ words [73]. For the case $n = 2$, the two-word sequence is referred to as a bigram. In a sentence with $l$ words, $n$ bigrams can be formed where $n = l - 1$, as the sentence contains $l - 1$ consecutive pairs of words [73]. Example 10 illustrates the bigram extraction process.

| Example sentence: | This is an example sentence. |
|---|---|
| Bigrams: | (this, is), (is, an), (an, example), (example, sentence) |

Example 10. Bigrams formed from an example sentence.

For bigram analysis, the dataset was lemmatized and cleared of stopwords, after which the most common bigrams were examined. This analysis aimed to identify any specific word pairs associated with a negative or positive stance.

### 5.3.3 Adjective-Based Framing

Framing bias refers to the way certain concepts are presented (framed) in text [74]. By reinforcing or emphasizing different aspects, the same concept can be framed differently [74]. In this thesis, word pairs consisting of an adjective and a noun were detected. It could be hypothesized that the concept of immigration is referred to as illegal, uncontrollable or unlawful in sentences with the against stance, whereas adjectives like lawful or controlled are used to frame immigration in the supportive stance.

## 5.4 Sentiment Analysis

Sentiment analysis is a task that involves using computational methods to determine the opinions, attitudes, and emotions expressed in text [75]. Based on its contents, a piece of text is characterized as either positive, negative, or neutral [20, 75]. Stance detection and sentiment analysis are closely related but differ slightly in their approach [20, 33, 76]. Sentiment analysis focuses on the general polarity of the text, whereas stance detection focuses on the viewpoint expressed towards a specific subject or target [20, 33, 76]. The following subchapters describe tools and resources available for sentiment analysis in Estonian texts and their use in this thesis.

### 5.4.1 Lexicon-based Approaches

Sentiment can be analyzed by using a list of annotated words. In her thesis, Regita Luukas [77] analyzed sentiment in-course feedback from the Study Information System using a lexicon of sentiment-annotated words. The lexicon was provided by the Institute of the Estonian Language (EKI) as a part of their valence corpus. Words that expressed positive sentiment were annotated with a score of 1, and negative words were annotated with -1. For each feedback entry, a base score of 0 was assigned. Entries were compared to the lexicon, and the base score was incremented or decremented accordingly when a match was found. Words with a negation had their score reversed. Following this, the final score for each entry was used to determine sentiment. Positive scores indicated positive sentiment and negative scores indicated negative sentiment. A score of zero either meant that the sentiment was ambiguous when there was an equal amount of positive or negative words or that the sentiment was neutral when no positive or negative words were found.

A similar approach to sentiment analysis was also used in this thesis to analyze the stance dataset. Before conducting the analysis, EKI's lexicon[11] was lemmatized and cleared of duplicates, as it contained all case forms for each word. After preprocessing, the lexicon consisted of 2454 words, of which 987 had positive sentiment and 1467 had negative sentiment.

---

[11] https://github.com/EKT1/valence/blob/master/valence/sqnad.csv

The lemmatized lexicon was compared to lemmatized sentences to calculate the final sentiment score for each row.

A second lexicon provided by Mohammad *et al.* [78] was also used. This lexicon, referred to as EmoLex, was initially constructed in English. However, a machine-translated version for Estonian was also provided. Google Translate was used to obtain the Estonian lexicon.

Words in EmoLex were annotated for ten emotions: anger, anticipation, disgust, fear, joy, negative, positive, sadness, surprise, trust [78]. These were reduced to two classes. Anticipation, joy, surprise, and trust were grouped with positive, and anger, disgust, fear, and sadness were grouped with negative. Words with both positive and negative scores greater than zero were considered ambiguous and, therefore, removed from the dataset, alongside neutral words with a zero score for both sentiments.

Similarly to EKI's lexicon, duplicate words were removed, and scores were unified to 1 and -1. Following this, a lexicon of 3693 sentiment-annotated words was obtained, of which 1785 had positive sentiment and 1908 had negative sentiment.

### 5.4.2 Emotsioonidetektor

The Institute of Estonian Language (EKI – Eesti Keele Instituut) is a national institution dedicated to the long-term survival of the Estonian language [79]. As part of a project led by Hille Pajupuu, the institute developed an Estonian language emotion detection tool [80]. The tool, named Emotsioonidetektor (Emotion Detector) is freely available and can be used via a web interface. It takes text as an input and classifies it as negative, neutral, or positive.

Emotsioonidetektor is based on a Naive Bayes classifier and trained on the Estonian Valence Corpus. This language resource consists of sentiment-annotated words and paragraphs [80]. Emotsioonidetektor differs from lexicon-based approaches since it also considers context, such as cases where a positive or negative word was negated, obtaining the opposite sentiment [80].

### 5.4.3  BERT

BERT (Bidirectional Encoder Representations from Transformers) is a pre-trained natural language processing model introduced by researchers at Google in 2018 [81]. BERT models can be used for various tasks, such as POS (part-of-speech) tagging, named entity recognition, dependency parsing, and sentiment analysis [81]. BERT's bidirectional training approach enables the model to be fine-tuned with a singular output layer, enabling its use for a wide variety of tasks without further architectural modifications [81].

EstBERT, an Estonian language-specific model, was trained in 2021 [82]. The model was trained on the Estonian National Corpus, the largest language resource available for Estonian consisting of around 1.34 billion words [82]. The relatively large amount of training data enabled it to achieve better results on some tasks than previously available multi-language BERT models [82]. For sentiment analysis in this thesis, both a fine-tuned EstBERT model and a multilingual XLM-RoBERTa model are used. The model *EstBERT128_Sentiment* was fine-tuned on the Estonian Valence Corpus and achieved an accuracy score of 0.74 [82]. However, it was noted that the multilingual XLM-RoBERTa model achieved a slightly better score of 0.76 [82].

## 5.5  Summary of Features

Table 2 provides an overview of the features discussed in this thesis. The total number of features is 32, out of which 10 are novel Estonian-specific features. For some lexical and Estonian-specific features, both the content (as *adjectives*) and frequency (as *adjectives_count*) are considered.

Table 2. Summary of features.

| Feature name | | Description |
|---|---|---|
| **LEXICAL** | *word_count* | The number of words in a sentence. |
| | *dependency_tree_height* | The height of a dependency tree based on automatic syntactic analysis performed by EstNLTK's Maltparser model. |
| | *flesch_score* | The Flesch Reading Ease Score as calculated by EstNLTK's *SentenceFleschScoreRetagger*. |
| | *named_entities* | A list of named entities extracted by EstNLTK's named entity tagger. |
| | *named_entites_count* | Number of named entities in a sentence. |
| | *noun_phrases* | A list of noun phrases extracted by EstNLTK's experimental noun phrase chunker. |
| | *noun_phrases_count* | Number of noun phrases in a sentence. |
| | *adjectives* | Lemmas of adjectives used in a sentence. |
| | *adjectives_count* | Number of adjectives used in a sentence. |
| | *quotes_count* | Number of quotes in a sentence. |
| | *quoted_words* | A list of words and short phrases that are between quotes in a sentence. |
| | *quoted_words_count* | Number of quoted words and short phrases. |
| **ESTONIAN-SPECIFIC** | *diminutives* | A list of words that are in the diminutive form, noted by the ending *-ke* or *-kene*. |
| | *diminutives_count* | Number of words in the diminutive form. |
| | *superlatives* | A list of adjectives in the superlative form. |
| | *superlatives_count* | The number of adjectives in the superlative form. |
| | *conditionals* | A list of verbs that are in the conditional form, noted by the suffix *-ks*. |
| | *conditionals_count* | Number of words in the conditional form. |
| | *translatives* | A list of nouns that are in the translative case, noted by the suffix *-ks*. |
| | *translatives_count* | Number of words in the translative case. |
| | *indirects* | A list of verbs that are indirect, noted by the suffix *-vat*. |
| | *indirects_count* | Number of indirect words. |
| **FRAMING** | *bw_count* | Number of words that insinuate black and white thinking. |
| | *has_against_bigram* | A categorical variable indicating whether an against or supportive stance bigram or adjective used for framing was present in the sentence or not. |
| | *has_support_bigram* | |
| | *framing_against* | |
| | *framing_supportive* | |
| **SENTI-MENT** | *ekilex_sentiment* | A sentiment classification of either negative, neutral, positive, as determined by the respective model. |
| | *emolex_sentiment* | |
| | *eki_emotion* | |
| | *estbert_sentiment* | |
| | *xlmroberta_sentiment* | |

By analyzing these features and their frequencies across sentences, their relevance in identifying political stance in Estonian news media can be determined.

# 6.    Results

This chapter provides an overview of the achieved results. Each feature was analyzed to determine its relevance for political stance detection.

## 6.1    Lexical Features

This chapter explores the results of the lexical features outlined in methodology. For each feature, some statistics and insights are provided, as well as their applicability in stance detection for Estonian texts.

### 6.1.1   Sentence Complexity

Three features related to the complexity of sentences were analyzed: word count, dependency tree height, and Flesch Reading Ease Score.

Table 3. Summary of statistics for feature *word_count*.

| stance | count | mean | std | min | 25% | 50% | 75% | max |
|---|---|---|---|---|---|---|---|---|
| against | 1175 | 22.32 | 10.19 | 3 | 15 | 20.33 | 28 | 94 |
| neutral | 1597 | 18.31 | 8.04 | 3 | 13 | 17 | 22.5 | 60 |
| supportive | 489 | 18.95 | 7.95 | 4 | 13 | 18 | 23 | 52 |
| combined | 3261 | 19.85 | 9.06 | 3 | 14 | 18 | 24 | 94 |

Table 3 outlines the statistics related to word count. The average number of words per sentence in the dataset is 19.85. With a mean of 22.32 words per sentence, it can also be observed that anti-immigration sentences have more words on average than sentences expressing a supportive or neutral stance.

Table 4. Summary of statistics for feature *dependency_tree_height*.

| stance | count | mean | std | min | 25% | 50% | 75% | max |
|---|---|---|---|---|---|---|---|---|
| against | 1175 | 6.33 | 1.94 | 2 | 5 | 6 | 7 | 18 |
| neutral | 1597 | 5.72 | 1.71 | 2 | 5 | 5.33 | 7 | 17 |
| supportive | 489 | 6.06 | 1.73 | 2 | 5 | 6 | 7 | 13 |
| combined | 3261 | 5.99 | 1.82 | 2 | 5 | 6 | 7 | 18 |

Table 4 gives an overview of the statistics for dependency tree height. The average height of the dependency tree across the dataset is 5.99. It can be noted that anti-immigration sentences have a slightly more complex structure, as the mean height of 6.33 is higher than for both neutral and supportive sentences.

Table 5. Summary of statistics for feature *flesch_score*.

| stance | count | mean | std | min | 25% | 50% | 75% | max |
|---|---|---|---|---|---|---|---|---|
| against | 1175 | 49.46 | 25.37 | -91.73 | 35.29 | 51.13 | 66.30 | 123.93 |
| neutral | 1597 | 56.76 | 22.74 | -36.52 | 42.73 | 57.50 | 72.38 | 134.12 |
| supportive | 489 | 54.56 | 23.11 | -29.21 | 40.53 | 53.76 | 70.30 | 129.57 |
| combined | 3621 | 53.80 | 24.00 | -91.73 | 39.49 | 54.96 | 69.79 | 134.12 |

Table 5 gives an overview of the statistics for the Flesch Reading Ease Score (FRES), a metric used to gauge the complexity of texts. FRES was calculated to be 53.8 in the dataset. As with word count and dependency tree height, it also followed that anti-immigration sentences were slightly more complex than neutral and supportive sentences. If these results on Estonian texts are interpreted in a similar manner as with English, then the average score of 49.46 for the against stance would indicate that the texts are difficult to read and are mostly appropriate for college-grade students. A score above 50 was present for both neutral and supportive stances, meaning that the reading level for these texts is generally suitable for high school students.

However, as illustrated by Table 5, directly applying this formula to Estonian texts may lead to unintended results, as some of the scores for some sentences fall outside of the predefined 0-100 range. There were 73 rows in the dataset with a negative score and 76 rows with a score above 100. This is likely because the structure of the Estonian language is more complex than English [83]. Given that the average number of words per sentence is 19.85, as per Table 3, the average number of syllables per word cannot surpass 2.20, in which case the score becomes negative. However, Estonian words likely contain more syllables on average than English words. For instance, the stopword list used in this thesis revealed an average syllable count of 3.47 for 5116 total words.

Following this, it can be concluded that the Flesch Reading Ease Score can be used to analyze trends or to get a general idea of the complexity of Estonian texts, but directly equating the

scores to grade levels as outlined by Zamanian and Heydari [56] is not appropriate, as education systems and grade levels vary by country, in addition to differing syllable and word counts across languages not being accounted for by the formula. By adjusting some of the values or scaling the results, a formula and result applicable to Estonian could be achieved. However, this requires further analysis of Estonian texts and the education system to produce a meaningful mapping from score to grade level.

In conclusion, all of the examined features related to text complexity reflected that sentences with an against stance towards immigration are slightly more complex than neutral or supportive sentences. However, the standard deviation in the against stance was also higher for all three features, suggesting varied complexity for anti-immigration sentences.

### 6.1.2 Named Entities

Table 6. Summary of statistics for feature *named_entities_count*.

| *named_entities_count* across all sentences. | | | | | | | | |
|---|---|---|---|---|---|---|---|---|
| stance | count | mean | std | min | 25% | 50% | 75% | max |
| against | 1175 | 1.54 | 1.43 | 0 | 0 | 1 | 2 | 10 |
| neutral | 1597 | 1.60 | 1.56 | 0 | 0 | 1 | 2 | 11 |
| supportive | 489 | 1.55 | 1.60 | 0 | 0 | 1 | 2 | 11 |
| combined | 3621 | 1.57 | 1.53 | 0 | 0 | 1 | 2 | 11 |
| Sentences with feature *named_entities_count* not equaling 0. | | | | | | | | |
| stance | count | mean | std | min | 25% | 50% | 75% | max |
| against | 897 | 2.01 | 1.33 | 0.33 | 1 | 2 | 3 | 10 |
| neutral | 1185 | 2.16 | 1.45 | 0.25 | 1 | 2 | 3 | 11 |
| supportive | 359 | 2.11 | 1.52 | 0.50 | 1 | 2 | 3 | 11 |
| combined | 2441 | 2.10 | 1.42 | 0.25 | 1 | 2 | 3 | 11 |

Table 6 gives an overview of the statistics related to the frequency of named entities. Around 70% of sentences for each stance contained at least one named entity. However, the mere presence of a named entity does not indicate that the sentence has a stance ($\chi^2 = 2.28, p = 0.32$). Furthermore, the variance for means across stances is minimal, indicating that counting named entities is not relevant for stance detection.

Table 7. Top 10 most common named entities per stance, skipping the top 2 for each stance, which were *eesti* (Estonia) and *euroopa* (Europe).

|     | against      | count | neutral      | count | supportive  | count |
|-----|--------------|-------|--------------|-------|-------------|-------|
| 3.  | rootsi       | 51    | kreeka       | 75    | kreeka      | 22    |
| 4.  | saksamaa     | 51    | euroopa liit | 67    | türgi       | 22    |
| 5.  | EKRE         | 47    | türgi        | 62    | saksamaa    | 20    |
| 6.  | euroopa liit | 44    | saksamaa     | 53    | soome       | 20    |
| 7.  | ungari       | 42    | rootsi       | 45    | euroopa liit| 19    |
| 8.  | itaalia      | 31    | itaalia      | 45    | süüria      | 16    |
| 9.  | prantsusmaa  | 29    | süüria       | 41    | prantsusmaa | 12    |
| 10. | vahemeri     | 28    | valgevene    | 34    | rootsi      | 12    |
| 11. | kreeka       | 28    | aafrika      | 34    | vahemeri    | 12    |
| 12. | helme        | 25    | ungari       | 34    | aafrika     | 10    |

The most common named entity is the country of Estonia (*eesti*), which occurred 751 times, followed by Europe (*euroopa*) with 619 occurrences. As shown in Table 7, Greece (*kreeka*), Sweden (*rootsi*) and Germany (*saksamaa*) are also mentioned, likely due to their role in the European migrant crisis in 2015.



Figure 2. Wordclouds of named entities for anti- and pro-immigration sentences.

Anti-immigration sentences mention the right-wing Estonian Conservative People's Party (EKRE) and the surname Helme, alluding to the leader of the party and his father, who are prominent figures in Estonian politics (Figure 2). However, the high number of mentions is likely because the dataset uses data from Uued Uudised, the media outlet of the party. The anti-immigration sentences also feature Hungary (*ungari*) and its prime minister, Viktor Orban, who is known for his anti-immigration stance [84].

These findings suggest that while named entities can offer insights about stance, the stance is primarily influenced by the political views of mentioned figures or governments, which may require external knowledge and context. In this dataset, frequently occurring named entities appear across all three classes, suggesting that they are not strongly indicative of stance.

### 6.1.3 Noun Phrases

Table 8. Summary of statistics for feature *noun_phrases_count*.

| stance | count | mean | std | min | 25% | 50% | 75% | max |
|---|---|---|---|---|---|---|---|---|
| *noun_phrases_count* across all sentences. | | | | | | | | |
| against | 1175 | 3.11 | 1.87 | 0 | 2 | 3 | 4 | 17 |
| neutral | 1597 | 2.83 | 1.69 | 0 | 2 | 3 | 4 | 15 |
| supportive | 489 | 2.82 | 1.68 | 0 | 2 | 3 | 4 | 11 |
| combined | 3621 | 2.93 | 1.76 | 0 | 2 | 3 | 4 | 17 |
| Sentences with feature *noun_phrases_count* not equaling 0. | | | | | | | | |
| stance | count | mean | std | min | 25% | 50% | 75% | max |
| against | 1132 | 3.23 | 1.80 | 0.50 | 2 | 3 | 4 | 17 |
| neutral | 1529 | 2.95 | 1.62 | 0.50 | 2 | 3 | 4 | 15 |
| supportive | 456 | 3.02 | 1.56 | 1 | 2 | 3 | 4 | 11 |
| combined | 3117 | 3.06 | 1.67 | 0.50 | 2 | 3 | 4 | 17 |

Table 8 gives an overview of the statistics related to the frequency of noun phrases. Over 90% of sentences in each class contained at least one noun phrase. Anti-immigration sentences have slightly more noun phrases than neutral or supportive sentences. Noun phrases were not significantly associated with stance in this dataset $(\chi^2 = 0.12, p = 0.73)$.

### 6.1.4 Adjectives

Sentences against immigration were found to have a higher number of adjectives compared to neutral or supportive stances. As seen in Table 9, sentences with an against stance have around 2.18 adjectives per sentence, more than in neutral and supportive sentences. Of 1175 sentences with an anti-immigration stance, 84% contained at least one adjective.

Table 9. Summary of statistics for feature *adjectives_count*.

| *adjectives_count* across all sentences. | | | | | | | | |
|---|---|---|---|---|---|---|---|---|
| **stance** | **count** | **mean** | **std** | **min** | **25%** | **50%** | **75%** | **max** |
| **against** | 1175 | 2.18 | 1.81 | 0 | 1 | 2 | 3 | 11 |
| **neutral** | 1597 | 1.50 | 1.47 | 0 | 0 | 1 | 2 | 13 |
| **supportive** | 489 | 1.64 | 1.60 | 0 | 1 | 1 | 2 | 12 |
| **combined** | 3621 | 1.77 | 1.66 | 0 | 1 | 1 | 3 | 13 |
| Sentences with feature *adjectives_count* not equaling 0. | | | | | | | | |
| **stance** | **count** | **mean** | **std** | **min** | **25%** | **50%** | **75%** | **max** |
| **against** | 1132 | 3.23 | 1.80 | 0.50 | 2 | 3 | 4 | 17 |
| **neutral** | 1529 | 2.95 | 1.62 | 0.50 | 2 | 3 | 4 | 15 |
| **supportive** | 456 | 3.02 | 1.56 | 1 | 2 | 3 | 4 | 11 |
| **combined** | 3117 | 3.06 | 1.67 | 0.50 | 2 | 3 | 4 | 17 |

By conducting a chi-square test, it was found that there is a significant association between stance and the presence of adjectives $(\chi^2 = 54.11, p < 0.01)$. These results confirm that adjectives can be used to detect stance, and suggest that authors against immigration tend to use more adjectives in their writing, likely as a way to express stronger emotion or judgment, which can influence the reader.



Figure 3. Wordclouds of adjectives for anti- and pro-immigration sentences.

Figure 3 shows the different adjectives used in against and supportive stances. Big (*suur*) is the most common adjective in anti-immigration sentences, followed by entire (*kogu*), and new (*uus*). These three adjectives are also present in the supportive stance. The word *new* is likely because of Uued Uudised (New News), the name of a media outlet where the data origi-nates from.

However, a difference can be noted when comparing rather negative adjectives such as illegal (*illegaalne*), massive (*massiivne*), and foreign (*võõras*), in the against stance, with hopeful and positive ones, such as international (*rahvusvaheline*), local (*kohalik*), and possible (*võimalik*), in the supportive stance. Therefore, the frequency and content of adjectives can aid in detection of political stance.

## 6.1.5 Quotes

Quotes were less common in the dataset, as around 20% of sentences contained them. Quotes were proportionally most common in the supportive stance, however it must be noted that quotes used for direct speech or titles are included.

Table 10. Summary of statistics for feature *quoted_words_count*.

| Sentences with feature *quoted_words_count* not equaling 0. | | | | | | | | |
|---|---|---|---|---|---|---|---|---|
| stance | count | mean | std | min | 25% | 50% | 75% | max |
| against | 104 | 1.05 | 0.60 | 0.25 | 0.63 | 1 | 1 | 4 |
| neutral | 53 | 1.06 | 0.55 | 0.25 | 1 | 1 | 1 | 3 |
| supportive | 27 | 1.18 | 0.67 | 0.25 | 1 | 1 | 1 | 3 |
| combined | 184 | 1.07 | 0.60 | 0.25 | 1 | 1 | 1 | 4 |

When excluding quotes used for their intended purposes, the against stance has the most words between quotes (and short phrases no greater than 5 words). Table 10 provides an overview and shows that 104 anti-immigration sentences contained a quoted word or phrase, which is higher than in both neutral and supportive stances combined.

Among the anti-immigration sentences, phrases such as European values, victims, and hate speech were between quotes, likely suggesting that these concepts were not taken seriously or mentioned in an ironic context, thus conveying an anti-immigration stance.

## 6.2 Estonian-specific Features

In this subchapter, the results and suitability for political stance detection of Estonian-specific features are analyzed.

### 6.2.1 Diminutives

Diminutives were infrequent in the dataset. Out of 3261 sentences, only 5 contained a word in the diminutive form. However, three of these sentences were annotated as against, and the remaining two were neutral. Although the diminutive is typically used to sound more gentle and positive, no pro-immigration sentences were found with this feature.

| |
|---|
| Kui Ameerikas tuli võimule Trump, lubasid paljud Hollywoodi näitlejad samuti emigreeruda ja **lumehelbekesed** akendest välja viskuda, aga jäid siiski kohapeale ussitama – BLM-i suitsulõhnalised meeleavaldused lubasid ennast vabalt maha maandada, Portlandis loodi koguni oma anarhistlik "autonoomia". [---] |
| When Trump came to power in America, many Hollywood actors promised to emigrate and **snowflakes** [promised] to throw themselves out of windows, but they still stuck around to nag - the BLM smoke-smelling protests allowed them to calm down, and in Portland, an anarchist "autonomy" was created. [---] |
| Hiljuti lõi "progressiivses maailmas" laineid Rootsi **lumehelbeke**, kes olevat justkui väljasaadetud afgaani elu päästnud – tegu oli paraku Rootsis juba tuntud kriminaaliga. |
| Recently, a Swedish **snowflake** made waves in the "progressive world" for supposedly saving the life of a deported Afghan – who was already a known criminal in Sweden, however. |

Example 10. Anti-immigration sentences using the diminutive form.

Two anti-immigration sentences shown in Example 10 contained the word *lumehelbeke* (snowflake), a derogatory term used to mock sensitive and delicate young adults who easily take offense and cannot tolerate conflict or criticism [85]. Both sentences also included quoted words, insinuating doubt and judgment.

These findings suggest that detecting and analyzing diminutives can be useful in political stance detection. However, as this feature is uncommon in this dataset and in news media texts overall, no significant conclusions can be made.

### 6.2.2 Superlative Form

Adjectives in the superlative form were uncommon among the sentences, as only 87 sentences contained them. Proportionally, the superlative form was slightly more common in the anti-immigration stance, but likely due to a lack of data, no significant association between stance and superlatives was found.

The most common superlative adjective was *suurim* (biggest), with 24 occurrences, followed by *parim* (best) and *kõige olulisem* (most important). However, no specific superlative adjective was typical for any stance. Similarly to diminutives, superlatives can possibly give insight into political stance detection but are not statistically significant ($\chi^2 < 0, p = 0.98$). As it follows, drawing significant conclusions is not possible due to the lack of data.

### 6.2.3 Conditional Form

Table 11. Summary of statistics for feature *conditionals_count*.

| Sentences with feature *conditionals_count* not equaling 0. | | | | | | | | |
|---|---|---|---|---|---|---|---|---|
| stance | count | mean | std | min | 25% | 50% | 75% | max |
| against | 167 | 1.20 | 0.72 | 0.33 | 1 | 1 | 1 | 4 |
| neutral | 140 | 1.15 | 0.55 | 0.50 | 1 | 1 | 1 | 5 |
| supportive | 73 | 1.14 | 0.49 | 0.33 | 1 | 1 | 1 | 3 |
| combined | 380 | 1.17 | 0.62 | 0.33 | 1 | 1 | 1 | 5 |

The conditional form was present in 380 sentences, as seen in Table 11. Proportionally, for both anti- and pro-immigration stances, around 14% of sentences contain a conditional form, in contrast to neutral sentences, where only 8% use conditionals.

Table 12. Summary of statistics for feature *translatives_count*.

| Sentences with feature *translatives_count* not equaling 0. | | | | | | | | |
|---|---|---|---|---|---|---|---|---|
| stance | count | mean | std | min | 25% | 50% | 75% | max |
| against | 326 | 1.27 | 0.74 | 0.33 | 1 | 1 | 1.88 | 6 |
| neutral | 310 | 1.23 | 0.60 | 0.33 | 1 | 1 | 1 | 4 |
| supportive | 127 | 1.22 | 0.60 | 0.33 | 1 | 1 | 1 | 5 |
| combined | 763 | 1.25 | 0.66 | 0.33 | 1 | 1 | 1 | 6 |

A similar trend can be noted also for the translative case, where 28% of against and 26% of supportive sentences use it, compared to 19% in the neutral sentences. The translative case was present in 763 sentences, as seen in Table 12.

Therefore, it can be concluded that these features can aid in political stance detection and that there is a statistically significant association between stance and frequency of both conditional form $(\chi^2 = 24.78, p < 0.01)$ and translative case $(\chi^2 = 27.31, p < 0.01)$. However, the content of the words in both conditional form and translative case do not reveal much insight about stance. The most common word across stances in the translative case is *näiteks* (for example). Similarly, the two most common words in the translative case are *oleks* (would be) and *peaks* (should be). These words are not indicative of stance solely on their own.

### 6.2.4 Indirect Speech

Indirect speech was also uncommon, as it was present only in 38 sentences. However, it was most prevalent in the anti-immigration sentences, with 22 occurrences, which was more than in the neutral and supportive stances combined. By far, the most common word in indirect speech was *olevat* (is said to be), with 24 occurrences, which was 50% of all the indirect words. This suggests that for political stance detection, the frequency of indirect verbs within sentences is more relevant $(\chi^2 = 5.38, p = 0.02)$, and the specific meaning of the word does not offer much insight.

## 6.3 Framing Analysis

This subchapter gives an overview of the results of framing-related features explored in the thesis and their suitability for political stance detection in Estonian news media.

### 6.3.1 Black-and-white Thinking

Polarizing words indicative of black-and-white thinking were present in 420 sentences. Anti-immigration sentences have the highest proportion, as 19% contain at least one of these words. There is a significant association between stance and the presence of polarizing words $(\chi^2 = 38.31, p < 0.01)$.

### 6.3.2 Bigram Analysis

From bigram analysis, the most common bigrams were extracted for an against and supportive stance. Table 13 provides an overview of the top 10 most common bigrams per stance.

Table 13. 10 most common bigrams in the against and supportive stances with English translations. Bigrams of interest are noted by a colored cell.

| | Against | Count | Supportive | Count |
|---|---|---|---|---|
| 1. | (euroopa, liit) (european, union) | 38 | (euroopa, liit) (european, union) | 29 |
| 2. | (mart, helme) (mart, helme) | 18 | (eesti, keel) (estonian, language) | 11 |
| 3. | (eesti, keel) (estonian, language) | 18 | (euroopa, komisjon) (european, commission) | 8 |
| 4. | (konservatiivne, rahvaerakond) (conservative, peoples party) | 14 | (miljon, euro) (million, euro) | 7 |
| 5. | (illegaalne, immigrant) (illegal, immigrant) | 13 | (välismaalane, seadus) (foreigner, law) | 6 |
| 6. | (kogu, euroopa) (whole, [of] europe) | 12 | (süüria, põgenik) (syrian, refugee) | 5 |
| 7. | (martin, helme) (martin, helme) | 11 | (eesti, pagulasabi) (estonian, refugee aid) | 5 |
| 8. | (tooma, kaasa) (bring, along) | 11 | (aafrika, päritolu) (african, origin) | 5 |
| 9. | (eesti, konservatiivne) (estonian, conservative) | 10 | (sisseränne, piirarv) (immigration, limit) | 5 |
| 10. | (neeger, araablane) (negro, arab) | 10 | (globaalne, ränderaamistik) (global, migration framework) | 5 |

In the against stance, terms like "illegal immigrant" are used, whereas the term "refugee" is more common in the supportive bigrams. Specific political figures from the EKRE party are also mentioned. However, this can be attributed to the sentences originating from the party's media outlet, Uued Uudised. Supportive stance bigrams feature terms related to humanitarian efforts and international collaboration. It can also be noted that the term *neeger*, which can be considered offensive, appears in the against stance, whereas the term *African origin* is used instead in the supportive stance.

Bigrams that were present in the neutral list were filtered out of both the against and supportive bigram lists, which resulted in 38 negative bigrams and 4 positive bigrams. The bigrams in these lists were subsequently matched with the lemmatized sentences. Following this, 70% of sentences that contained an against bigram were correctly identified to have an against stance. For supportive bigrams, the result was 61%.

### 6.3.3 Adjective-based Framing

Table 14. 10 most common adjective-noun pairs in the against and supportive stances with English translations. Pairs of interest are noted by a colored cell.

| | Against | Count | Supportive | Count |
|---|---|---|---|---|
| 1. | (illegaalsete, immigrantide) (illegal, immigrants) | 8 | (ebaseadusliku, rände) *(unlawful, migration)* | 3 |
| 2. | (odava, tööjõu) (cheap, labour) | 8 | (rahvusvahelist, kaitset) *(international, defense)* | 3 |
| 3. | (konservatiivne, rahvaerakond) (conservative, peoples party) | 7 | (rahvusvahelise, rändekava) *(international, migration plan)* | 2 |
| 4. | (massilise, sisserände) (massive, immigration) | 4 | (soolise, võrdõiguslikkuse) *(gender, equality)* | 2 |
| 5. | (uute, uudiste) (new, news) | 4 | (avatud, algus) *(open, beginning)* | 2 |
| 6. | (illegaalse, immigratsiooni) (illegal, immigration) | 4 | (salliva, õpikeskkonna) *(tolerant, learning environment)* | 2 |
| 7. | (uus, valitsus) (new, government) | 3 | (kogu, maailmas) *([in the] entire, world)* | 2 |
| 8. | (uued, uudised) (new, news) | 3 | (suure, panuse) *(big, contribution)* | 2 |
| 9. | (illegaalseid, immigrante) (illegal, immigrants) | 3 | (globaalses, ränderaamistikus) *(global, migration framework)* | 2 |
| 10. | (suur, probleem) (big, problem) | 3 | (rassilise, diskrimineerimise) *(racial, discrimination)* | 2 |

In the against stance, terms such as *illegal immigrants*, *cheap labor*, and *massive immigration* are present as seen in Table 14, thus focusing more on the legality, scale, and economic impact of immigration. In the positive stance, terms like *international defense*, *gender equality*, and

*big contribution* are highlighted, which frame immigration in a more humane, tolerant, and progressive manner.

Table 15 shows how different adjectives were used to frame the same concept, namely immigration. Adjectives preceding word stems *immigra* and *rän*[12] were captured. Duplicate and ambiguous terms that also appeared in the neutral stance were removed.

Table 15. Lemmatized list of unique adjectives in Estonian used to frame immigration, that preceded the stems *immigra* and *rän*.

| Against | Terms in both | Supportive |
|---|---|---|
| agressiivne, allaheitlik, avantüristlik, efektiivne, elama, isiklik, islamiusuline, jahtiv, järgmine, jätkuv, kahjulik, kogu, konservatiivne, kriminaalne, kuritahtlik, käiv, kõrge, lõtv, ohtlik, paarituhandeline, potentsiaalne, range, rekordkõrge, riiklik, salakaval, sarnane, sealne, senine, seotud, suunduv, suvaline, tark, teisene, toimuv, tugevnev, tuntud, tülikas, valimatu, ähvardav, ühine, üksik, üleeuroopaline | lähtuv, piiramatu, kasvav | esitatud, hiiglaslik, inimlik, laiahaardeline, lubatud, noor, oluline, seaduslik, tõstatatud, vaba, väärikas, üleilmne |

The against stance uses a more varied set of words to frame immigration. Words such as *criminal*, *record high*, and *threatening* are used. In contrast, the supportive stance uses terms like *humane*, *legal*, and *allowed*.

Adjectives in the against and supportive lists were matched to lemmatized sentences in the dataset. For the against stance, 53% of sentences that mentioned immigration and contained an adjective from the corresponding list were correctly identified as being anti-immigration. This approach did not work for the positive stance, as the result was only 24%, likely resulting from the smaller and less varied list of adjectives.

---

[12] *rän* comes from the word *ränne*, which is a word often used to refer to *migration* or *travel*.

## 6.4    Sentiment Analysis

This chapter explores the results of the sentiment analysis tools. For interpreting the results, the against, supportive, and neutral stances were mapped to negative, positive, and neutral sentiment, respectively.

Precision, recall, and F1-score are metrics that are used to evaluate the results of classification tasks [86]. These measures are calculated based on the classification model's predictions and the actual true labels in the dataset. Confusion matrices, similar to Figure 4, are used to give an overview of a model's classification performance.

|  |  | **Predicted label** annotated by the model | |
|---|---|---|---|
|  |  | Positive | Negative |
| **Actual (true) label** annotated by a human | Positive | TP (true positive) | FN (false negative) |
|  | Negative | FP (false positive) | TN (true negative) |

Figure 4. Confusion matrix.

**Precision (P):** Precision highlights the model's ability to make accurate positive predictions. It measures the proportion of true positive predictions among all positive predictions, and is defined by the formula (2)

$$P = \frac{TP}{TP + FP} \tag{2}$$

**Recall (R):** Recall highlights the model's ability to capture all positive instances. It measures the proportion of true positive predictions among all actual positive instances in the dataset, and is defined by the formula (3)

$$R = \frac{TP}{TP + FN} \tag{3}$$

**F1 score:** The F1 score is the weighted average of the precision and recall measures. It accounts for imbalanced class sizes, in the case of which precision and recall can be misleading. The F1 score is defined by the formula (4)

$$\text{F1} = \frac{TP}{TP + FN} \qquad (4)$$

For multiclass classification, these metrics are calculated separately for each class by using the respective true positive, false positive, and false negative counts.

### 6.4.1  Lexicon-based Approaches

Lexicon-based approaches performed poorly on the dataset. Using EKI's lexicon, an accuracy of 0.45 was achieved (Figure 5).
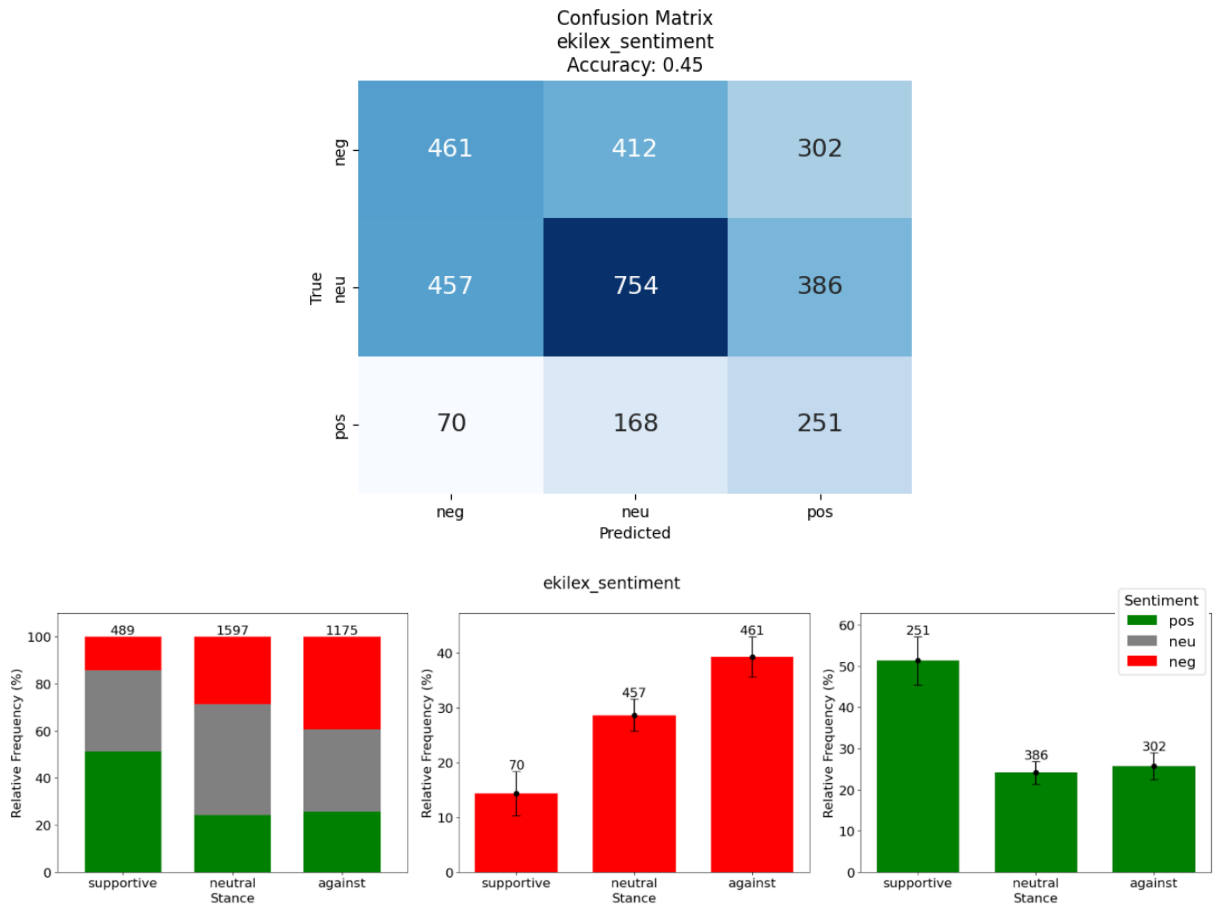


Figure 5. Confusion matrix and relative frequency graphs for sentiment predictions using the lexicon provided by the Institute of the Estonian Language (EKI)

For EKI's lexicon, a trend can be noted, as seen in Figure 5. The proportion of negative sentiment increased consistently, as the supportive stance had the least negative sentiment (15%), whereas the against stance had the most (48%). However, this trend did not follow for

positive sentiment, where the neutral and against classes have overlapping confidence intervals and roughly an equal proportion of positive sentences, 24% and 26%, respectively.
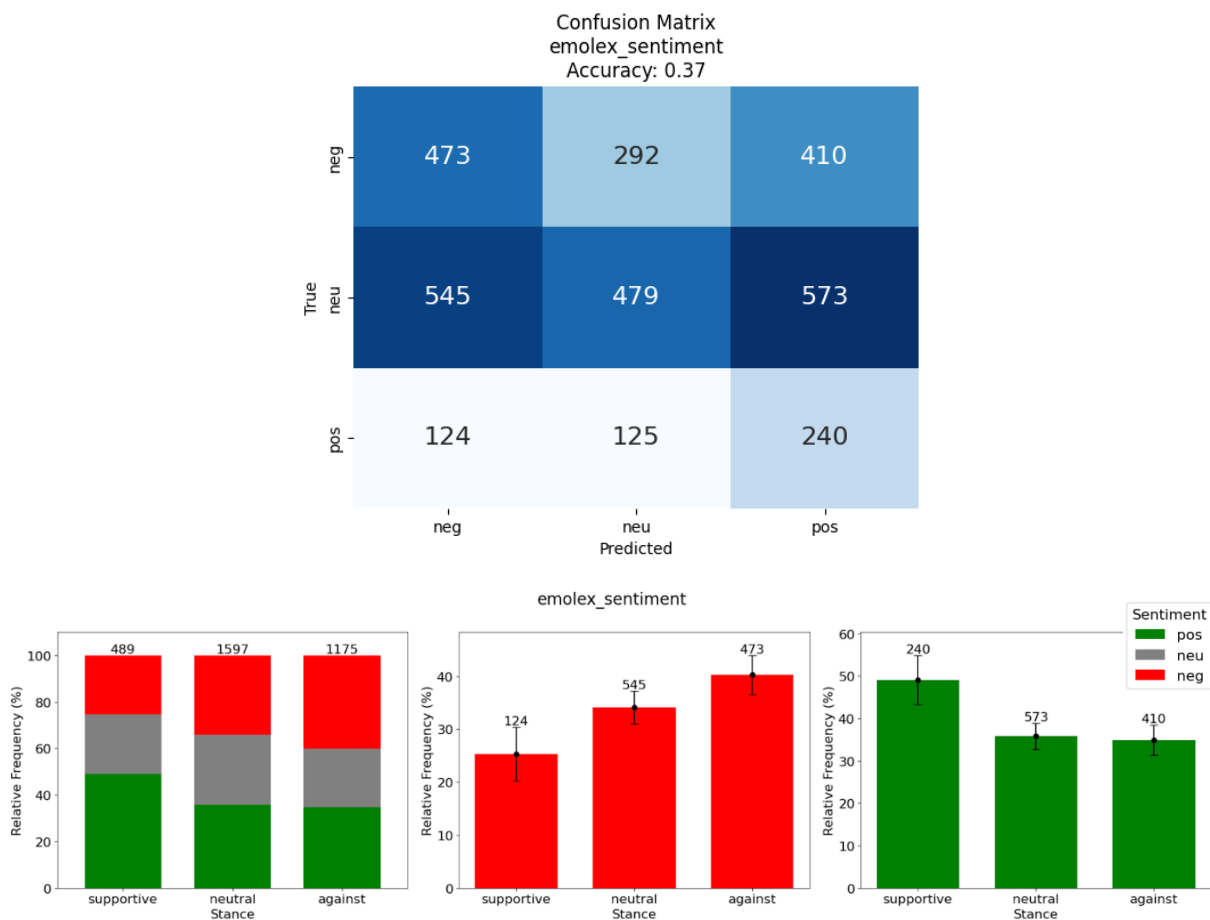


Figure 6. Confusion matrix and relative frequency graphs for sentiment predictions using the lexicon provided by Mohammad *et al.* [77].

The EmoLex lexicon provided by Mohammad *et al.* [77] achieved an accuracy of 0.35, as it largely over-favored the positive class. As pro-immigration sentences had the lowest representation in the dataset, predicting positive sentiment increased the amount of false negatives and subsequently resulted in low precision (0.20). Therefore, the results are not confident for both positive and negative sentiment, as shown by Figure 6.

### 6.4.2 Emotsioonidetektor

Emotsioonidetektor performed poorly on the stance dataset. As shown by Figure 7, Emotsioonidetektor achieved an accuracy score of 0.33, the same as random guessing with three classes. In terms of sentiment, the tool was polarizing, as it tended to predict the sentence to either be negative or positive, thus largely ignoring the neutral class. Only 173 sentences were predicted to be neutral, out of which 116 were predicted correctly. However, the dataset contained 1597 neutral sentences in total, the highest proportion (48%) out of all the three classes.
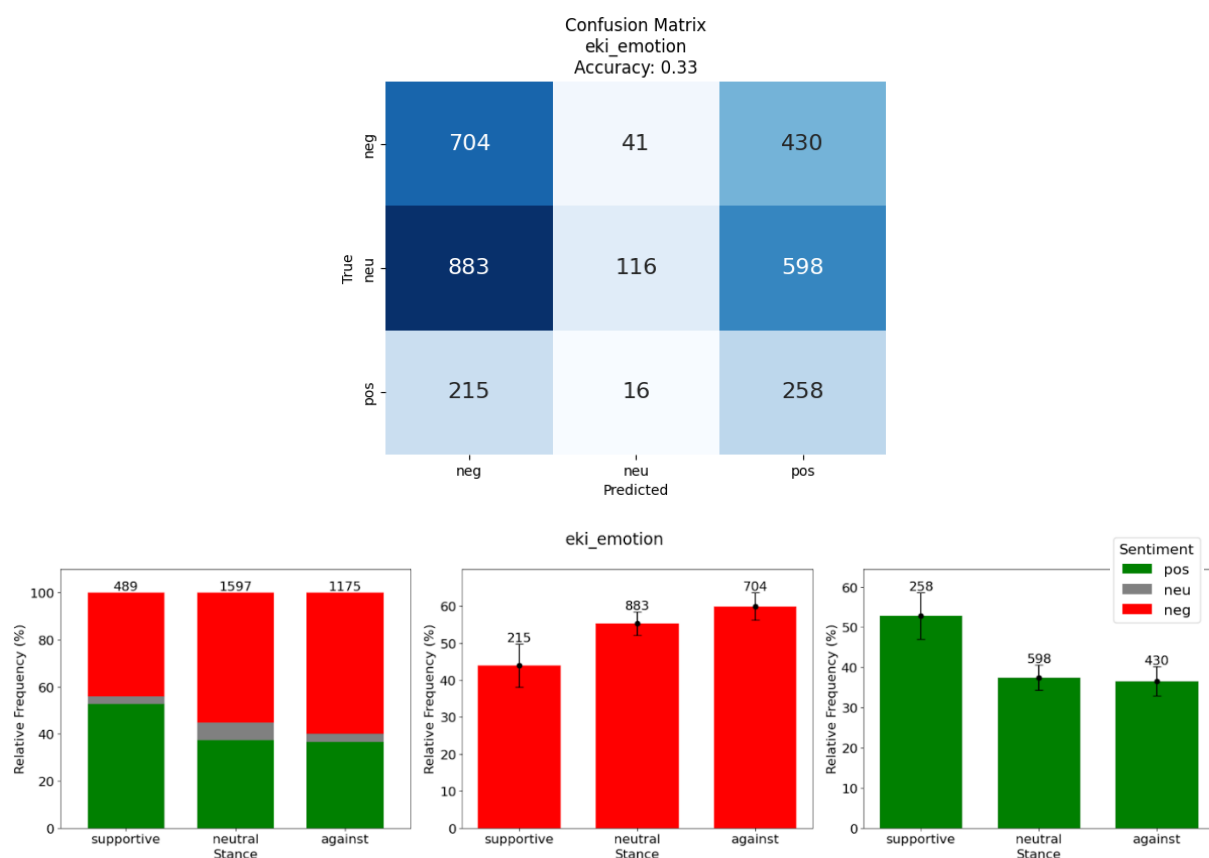


Figure 7. Confusion matrix and relative frequency graphs displaying the results of Emotsioonidetektor predictions.

The poor results can be partly attributed to the unit of analysis being a sentence instead of a longer paragraph, which was used for training the tool. In addition, the main objective of Pajupuu *et al.* was to predict the possible effect of a written text on the reader instead of the general sentiments expressed in written text by its author [79].

48

### 6.4.3 BERT

The EstBERT sentiment model achieved an accuracy score of 0.49, while XLM-RoBERTa achieved a score of 0.58, as per Figure 8.
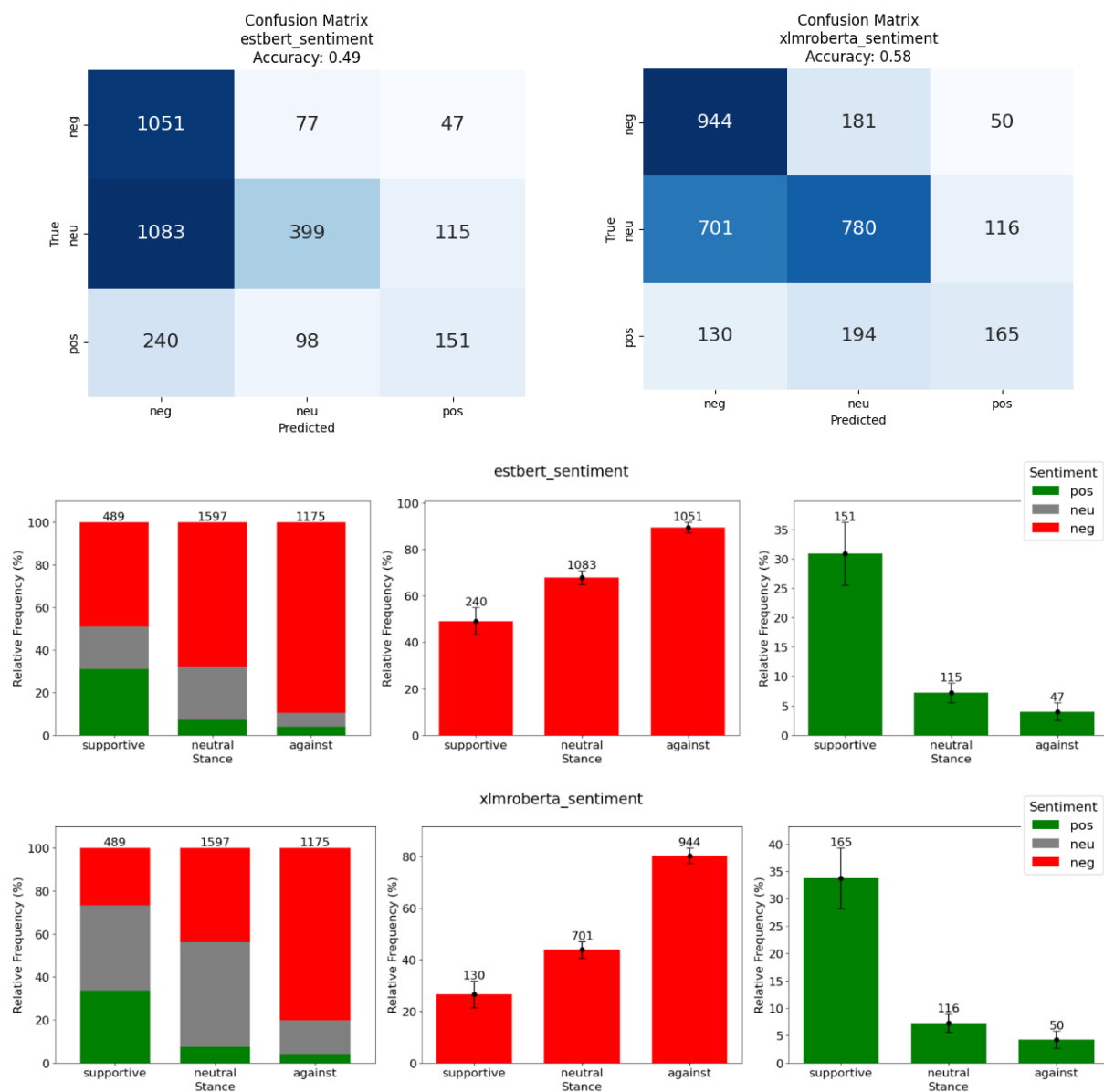


Figure 8. Confusion matrices and relative frequency graphs displaying the results of BERT sentiment model predictions.

The results of the multilingual model are better since it predicted a higher proportion of sentences to be neutral. It can also be seen that both models tended to predict sentences to be

negative, thus struggling more with the neutral and positive classes. Neutral sentences were labeled negative, and positive sentences were labeled neutral. Anti-immigration sentences were predicted to have negative sentiment, resulting in a high recall score for both models, 0.89 for EstBERT and 0.8 for XLM-RoBERTa.

### 6.4.4  Sentiment Analysis Results

Table 16 showcases the precision, recall and F1-score per stance for each of the sentiment analysis tools.

Table 16. Evaluation metrics for each stance class. P – precision, R – recall, F1 – F1-score.

| Model | against | | | neutral | | | supportive | | |
|---|---|---|---|---|---|---|---|---|---|
| | P | R | F1 | P | R | F1 | P | R | F1 |
| EKI | 0.47 | 0.39 | 0.43 | 0.57 | 0.47 | 0.51 | 0.27 | 0.51 | 0.35 |
| EmoLex | 0.41 | 0.40 | 0.41 | 0.53 | 0.30 | 0.38 | 0.20 | 0.49 | 0.28 |
| Emotsioonidetektor | 0.39 | 0.60 | 0.47 | 0.67 | 0.07 | 0.13 | 0.20 | **0.53** | 0.29 |
| EstBERT | 0.44 | **0.89** | 0.59 | **0.70** | 0.25 | 0.37 | 0.48 | 0.31 | 0.38 |
| XLM-RoBERTa | **0.53** | 0.80 | **0.64** | 0.68 | **0.49** | **0.57** | **0.50** | 0.34 | **0.40** |

In comparison to all other models, the multilingual XLM-RoBERTa sentiment model achieved better F1-scores in all three classes. The F1-scores range from 7.8%-35.9% in the negative class, 10.5%-77.2% in the neutral class, and 5%-30% in the positive class.

## 6.5    Summary of Insights

In this thesis, 32 features were extracted from immigration-related sentences, and their suitability for stance detection was analyzed. Table 17 summarizes which features were considered helpful in stance detection.

Table 17. Classification of features based on their usefulness in political stance detection.

| Useful | Not Useful |
|---|---|
| *adjectives, adjectives_count, quotes_count, quoted_words, quoted_words_count, diminutives, conditionals_count, translatives_count, indirects_count, bw_count, has_against_bigram, has_support_bigram, framing_against, estbert_sentiment, xlmroberta_sentiment* | *word_count, dependency_tree_height, flesch_score, named_entities, named_entities_count, noun_phrases, noun_phrases_count, diminutives_count, superlatives, superlatives_count, conditionals, translatives, indirects, framing_supportive, ekilex_sentiment, emolex_sentiment, eki_emotion* |

In total, out of 32 features, 15 features were deemed useful and 17 were not useful. The useful features can be used for automated political stance detection tasks, as a significant association could be noted between the frequency of some of these features and an anti- or pro-immigration stance present in a sentence. The contents of these features (adjectives, quoted words) also show a divide between against, neutral and supportive sentences.

Unhelpful features were either unsuitable for Estonian texts (*flesch_score*), lacked clear intuition (*noun_phrases*), or were too uncommon (*diminutives_count*). Uncommon features were most notable in the Estonian-specific features, which seemed promising and had the possibility reveal insights. However, due to the small dataset and a lack of these features, their usefulness in stance detection could not be determined.

Additionally, some features, such like the counts of punctuation marks or stopwords were also considered. However, these did not reveal much insight, as they are fairly generic and lacked intuition in stance detection.

# 7.    Conclusion

The goal of this thesis was to explore and identify features that are indicative of political stance in Estonian news media. Stance detection is the classification of a text as either against, neutral, or supportive of a certain topic. In this thesis, the stance was detected regarding the politically divisive topic of immigration, a concept related to the international movement of people.

Automated political stance detection in Estonian is challenging as highlighted by the lack of data, language processing tools, and related works in this field. The research in this thesis included outlining different biases and stance detection techniques in existing works, extrapolating existing English approaches to Estonian texts, proposing novel features specific to Estonian news media, and analyzing the framing and sentiment of 3621 sentences in an immigration-related dataset.

In this thesis, 32 total features were identified, which were split into four main categories: lexical features, Estonian-specific features, framing-related features, and sentiment-related features. These features were exhaustively analyzed to determine their suitability for political stance detection in Estonian news media. The goal of the thesis was achieved, as 15 features were shown to be helpful in political stance detection. Out of 10 novel Estonian-specific features, 4 were identified to be useful in political stance detection: diminutives, and the frequency of words in the conditional form, translative case, and indirect speech.

As online news media continues to grow, ensuring the integrity and fairness of news reporting becomes increasingly crucial. Being aware of unethical practices can aid the public in making informed decisions, whether in voting for political representation or treating others with respect.

Improvements to this work could be made by introducing and detecting words from a lexicon that is biased or derogatory towards immigrants. Additionally, the rich morphology of Estonian could be analyzed further by conducting a full frequency analysis of all cases and conjugation forms, which could reveal more features and further insights about stance or sentiment. A more extensive political stance dataset would also enable a model to be trained or fine-tuned that utilizes these features.

# References

[1]     Bhuller M, Havnes T, McCauley J, Mogstad M. How the Internet Changed the Market for Print Media. *American Economic Journal: Applied Economics.* 2024, vol. 16, p. 318–358. https://doi.org/10.1257/app.20210689

[2]     Bucy E, Gantz W, Wang Z. Media Technology and the 24-Hour News Cycle. *Communication technology and social change.* 2007, p. 143–164. https://www.researchgate.net/publication/313055194 (15.05.2024)

[3]     Spinde T. An Interdisciplinary Approach for the Automated Detection and Visualization of Media Bias in News Articles. *2021 IEEE International Conference on Data Mining Workshops.* 2021, p. 1096-1103. https://doi.org/10.48550/arXiv.2112.13352

[4]     Lei Y, Huang R, Wang L, Beauchamp N. Sentence-level Media Bias Analysis Informed by Discourse Structures. *Proceedings of the 2022 Conference on Empirical Methods in Natural Language Processing.* 2022, p. 10040–10050. https://doi.org/10.18653/v1/2022.emnlp-main.682

[5]     ALDayel A, Magdy W. Stance detection on social media: State of the art and trends. *Information Processing & Management.* 2021, vol. 58, p. 102597. https://doi.org/10.1016/j.ipm.2021.102597

[6]     Farsi S, Eusha A, Arefin M. CUET_Binary_Hackers at ClimateActivism 2024: A Comprehensive Evaluation and Superior Performance of Transformer-based Models in Hate Speech Event Detection and Stance Classification for Climate Activism. *Proceedings of the 7th Workshop on Challenges and Applications of Automated Extraction of Socio-political Events from Text (CASE 2024).* 2024, p. 145-155. https://aclanthology.org/2024.case-1.20/ (15.05.2024)

[7]     Mohammad S, Kiritchenko S, Sobhani P, Zhu X, Cherry C. A Dataset for Detecting Stance in Tweets. *Proceedings of the Tenth International Conference on Language Resources and Evaluation (LREC'16).* 2016, p. 3945-3952. https://aclanthology.org/L16-1623 (15.05.2024)

[8]     Statistics Estonia. Population by ethnic nationality. https://stat.ee/en/find-statistics/statistics-theme/population/population-figure#:~:text=Population%20by%20ethnic%20nationality (15.05.2024)

[9]     Dyvik E. The most spoken languages worldwide 2023. https://www.statista.com/statistics/266808/the-most-spoken-languages-worldwide (15.05.2024)

[10] Piir R. Finland's ChatGPT equivalent begins to think in Estonian as well. *ERR*, 2023. https://news.err.ee/1609120697/finland-s-chatgpt-equivalent-begins-to-think-in-estonian-as-well (15.05.2024)

[11] Hedderich M, Lange L, Adel H, Strötgen J, Klakow D. A Survey on Recent Approaches for Natural Language Processing in Low-Resource Scenarios. *arXiv preprint arXiv:2010.12309.* 2020. https://arxiv.org/pdf/2010.12309 (15.05.2024)

[12] Statistics Estonia. Population. https://stat.ee/en/find-statistics/statistics-theme/population/population-figure#:~:text=Tabel%20RV021 (15.05.2024)

[13] Delfi Meedia. https://delfimeedia.ee/ (15.05.2024)

[14] Eesti Meediaettevõtete Liit. Tasulised Digitellimused 2023. https://meedialiit.ee/statistika/statistika-2023/ (15.05.2024)

[15] Cambridge Dictionary Stance. Cambridge University Press. https://dictionary.cambridge.org/dictionary/english/stance

[16] Biber D, Finegan E. Adverbial stance types in English. *Discourse Processes.* 1988, vol. 11, p. 1–34. https://doi.org/10.1080/01638538809544689

[17] Kiesling S. Stance and Stancetaking. *Annual Review of Linguistics.* 2022, vol. 8, p. 409–426. https://doi.org/10.1146/annurev-linguistics-031120-121256

[18] Du Bois J. The stance triangle. *Pragmatics & Beyond New Series*. Amsterdam: John Benjamins Publishing Company. 2007, p. 139–182. http://dx.doi.org/10.1075/pbns.164.07du

[19] Mohammad S, Sobhani P, Kiritchenko S. Stance and Sentiment in Tweets. *ACM Transactions on Internet Technology (TOIT)*. 2017, vol. 17, p. 1–23. https://doi.org/10.1145/3003433

[20] Küçük D, Can F. Stance Detection: A Survey. *ACM Computing Surveys (CSUR).* 2020, vol. 53, p. 1-37. https://doi.org/10.1145/3369026

[21] Merriam-Webster Article. An Encyclopaedia Britannica Company. https://www.merriam-webster.com/dictionary/immigrate

[22] International Organization for Migration. Glossary on Migration. 2019. https://publications.iom.int/system/files/pdf/iml_34_glossary.pdf (15.05.2024)

[23] Kosho J. Media Influence On Public Opinion Attitudes Toward The Migration Crisis. *International Journal of Scientific & Technology Research.* 2016, vol. 5, p. 86–91. https://www.ijstr.org/final-print/may2016/Media-Influence-On-Public-Opinion-Attitudes-Toward-The-Migration-Crisis.pdf (15.05.2024)

[24] Vetik R. Eesti elanike hoiakud poliitilise integratsiooniga seoses. 2000. https://integratsioon.ee/sites/default/files/Mon2000_5Vetik.pdf (15.05.2024)

[25] Päll R. SISSERÄNDEVASTASTE HOIAKUTE MÕJU POLIITILISELE USALDUSELE EUROOPA RÄNDEKRIISI VALGUSES: UNGARI JA POOLA NÄITEL. Bachelor's thesis, University of Tartu, Johan Skytte Institute of Political Studies. 2021. https://dspace.ut.ee/server/api/core/bitstreams/2726d87c-bb9b-42e7-9208-6e19eea6ddef/content (15.05.2024)

[26] Maksimova N. IKT eriala üliõpilaste varajase väljalangemise ennustamine masinõppe meetodite abil. Master's thesis, Tallinn University of Technology, School of Information Technologies. 2022. https://digikogu.taltech.ee/en/Item/380c10c6-dc43-450c-8051-4dd9d3f88d97 (15.05.2024)

[27] Benoit K, Watanabe K, Wang H, Nulty P, Obeng A, Müller S, Matsuo A. quanteda: An R package for the quantitative analysis of textual data. *Journal of Open Source Software.* 2018, vol. 3, p. 774. https://doi.org/10.21105/joss.00774

[28] Watanabe K. Different Languages. *quanteda tutorials.* https://tutorials.quanteda.io/multilingual/overview/ (15.05.2024)

[29] Mehrabi N, Morstatter F, Saxena N, Lerman K, Galstyan A. A Survey on Bias and Fairness in Machine Learning. *ACM computing surveys (CSUR)*. 2021, vol. 54, p. 1–35. https://doi.org/10.1145/3457607

[30] Spinde T, Hamborg F, Donnay K, Becerra A, Gipp B. Enabling News Consumers to View and Understand Biased News Coverage: A Study on the Perception and Visualization of Media Bias. *Proceedings of the ACM/IEEE Joint Conference on Digital Libraries in 2020*. 2020, p. 389-392. https://doi.org/10.1145/3383583.3398619

[31] Sen A, Ghatak D, Khanuja G, Rekha K, Gupta M, Dhakate S, Sharma K, Seth A.. Analysis of Media Bias in Policy Discourse in India. *ACM SIGCAS/SIGCHI Conference on Computing and Sustainable Societies (COMPASS)*. 2022, pp. 57–77. https://doi.org/10.1145/3530190.3534798

[32] Spinde T, Rudnitckaia L, Mitrović J, Hamborg F, Granitzer M, Gipp B, Donnay K. Automated identification of bias inducing words in news articles using linguistic and context-oriented features. *Information Processing & Management.* 2021, vol. 58, p. 102505. https://doi.org/10.1016/j.ipm.2021.102505

[33] Mohammad S, Kiritchenko S, Sobhani P, Zhu X, Cherry C. Semeval-2016 Task 6: Detecting Stance in Tweets. *Proceedings of the 10th international workshop on semantic evaluation (SemEval-2016).* 2016, p. 31–41. https://aclanthology.org/S16-1003.pdf (15.05.2024)

[34] Guo X, Ma W, Vosoughi S. Measuring Media Bias via Masked Language Modeling. *Proceedings of the International AAAI Conference on Web and Social Media*. 2022, vol. 16, p. 1404–1408. https://doi.org/10.1609/icwsm.v16i1.19396

[35] Lehmann R, Derczynski L. Political Stance in Danish. *Proceedings of the 22nd Nordic Conference on Computational Linguistics*. 2019, p. 197–207. https://aclanthology.org/W19-6121 (15.05.2024)

[36] Yantseva V, Kucher K. Machine Learning for Social Sciences: Stance Classification of User Messages on a Migrant-Critical Discussion Forum. *2021 Swedish Workshop on Data Science (SweDS)*. Växjö: IEEE. 2021, p. 1–8. https://doi.org/10.1109/SweDS53855.2021.9637718

[37] Mets M, Karjus A, Ibrus I, Schich M. Automated stance detection in complex topics and small languages: the challenging case of immigration in polarizing news media. *arXiv e-prints*. 2023. https://doi.org/10.48550/arXiv.2305.13047

[38] Ulčar M, Robnik-Šikonja M. Training Dataset and Dictionary Sizes Matter in BERT Models: The Case of Baltic Languages. *Analysis of Images, Social Networks and Texts.* Cham: Springer International Publishing. 2021, vol. 13217, p. 162–172. https://doi.org/10.1007/978-3-031-16500-9_14

[39] Dolenko P. POLITICAL NEUTRALITY OR EDITORIAL SLANT? COMPARING COVERAGE AND CONTENT OF ESTONIA'S LARGEST NATIONWIDE DAILY NEWSPAPERS. Master's thesis, University of Tartu, Johan Skytte Institute of Political Studies. 2022. https://dspace.ut.ee/server/api/core/bitstreams/229c93b5-a59f-475c-bd08-cdb6e46ddd0e/content (15.05.2024)

[40] Marling R. The Intimidating Other: Feminist Critical Discourse Analysis of the Representation of Feminism in Estonian Print Media. *NORA—Nordic Journal of Feminist and Gender Research.* 2010, vol. 18, p. 7–19. https://www.doi.org/10.1080/08038741003626767

[41] Kivisalu K. #METOO MOVEMENT IN ESTONIA: A FRAME ANALYSIS. Master's thesis, University of Tartu, Johan Skytte Institute of Political Studies. 2019. https://dspace.ut.ee/server/api/core/bitstreams/b09be81f-f349-4522-baa4-870c7a3435bc/content (15.05.2024)

[42] Kaukonen E. Sooliselt markeeritud sõnad eesti spordiuudistes. *Keel ja Kirjandus*. 2022, vol. 65, p. 526–545. https://doi.org/10.54013/kk774a3

[43] Roosve G. Ajakirjaniku sekkuja rolli esinemine eesti uudistekstides. Bachelor's thesis, University of Tartu, Institute of Social Studies. 2021. https://dspace.ut.ee/server/api/core/bitstreams/14a20dcb-9c0f-4134-8b13-04b6374109a7/content (15.05.2024)

[44] VanderPlas J. Python Data Science Handbook: Essential Tools for Working with Data. Sebastopol: O'Reilly Media, Inc. 2016.

[45] Bisong E. Building Machine Learning and Deep Learning Models on Google Cloud Platform: A Comprehensive Guide for Beginners. Berkeley: Apress. 2019.

[46] McKinney W. Data Structures for Statistical Computing in Python. *Proceedings of the 9th Python in Science Conference (SciPy 2010)*. 2010, p. 56–61. https://doi.org/10.25080/Majora-92bf1922-00a

[47] Laur S, Orasmaa S, Särg D, Tammo P. EstNLTK 1.6: Remastered Estonian NLP Pipeline. *Proceedings of The 12th Language Resources and Evaluation Conference.* Marseille: European Language Resources Association. 2020, vol. 12, p. 7154–7162. https://www.aclweb.org/anthology/2020.lrec-1.884 (15.05.2024)

[48] Pedregosa F, Varoquaux G, Gramfort A, Michel V, Thirion B, Grisel O, Blondel M, Prettenhofer P, Weiss R, Dubourg V, Vanderplas J, Passos A, Cournapeau D, Brucher M, Perrot M, Duchesnay E. Scikit-learn: Machine Learning in Python. *Journal of Machine Learning Research*. 2011, vol. 12, p. 2825–2830. https://jmlr.csail.mit.edu/papers/volume12/pedregosa11a/pedregosa11a.pdf (15.05.2024)

[49] Hunter J. Matplotlib: A 2D Graphics Environment. *Computing in Science & Engineering.* 2007, vol. 9, p. 90-95. https://doi.org/10.1109/MCSE.2007.55

[50] Wolf T, Debut L, Sanh V, Chaumond J, Delangue C, Moi A, Cistac P, Rault T, Louf R, Funtowicz M, Davison J, Shleifer S, von Platen P, Ma C, Jernite Y, Plu J, Xu C, Le Scao T, Gugger S, Drame M, Lhoest Q, Rush A. Transformers: State-of-the-Art Natural Language Processing. *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing: System Demonstrations*. Online: Association for Computational Linguistics. 2020, p. 38–45. https://doi.org/10.18653/v1/2020.emnlp-demos.6

[51] University of Hull. Public Communication: Newspaper Article. https://libguides.hull.ac.uk/public-comm/newspaper-article (15.05.2024)

[52] Merriam-Webster Article. An Encyclopaedia Britannica Company. https://www.merriam-webster.com/dictionary/article

[53] Merriam-Webster Lexical. An Encyclopaedia Britannica Company. https://www.merriam-webster.com/dictionary/lexical

[54] Nivre J, Hall J, Nilsson J. Maltparser: A data-driven parser-generator for dependency parsing. *LREC 2006*. 2006, vol. 6, p. 2216–2219. http://lrec-conf.org/proceedings/lrec2006/pdf/162_pdf.pdf (15.05.2024)

[55] Nivre J. Dependency Parsing. *Language and Linguist Compass*. 2010, vol. 4, p. 138–152. https://doi.org/10.1111/j.1749-818X.2010.00187.x

[56] Zamanian M, Heydari P. Readability of Texts: State of the Art. *Theory and Practice in Language Studies.* 2012, vol. 2, p. 43–53. https://doi.org/10.4304/tpls.2.1.43-53

[57] Jurafsky D, Martin J. Sequence Labeling for Parts of Speech and Named Entities. *Speech and Language Processing*. 2024, vol. 3, p. 162–185. https://web.stanford.edu/~jurafsky/slp3/8.pdf (15.05.2024)

[58] Tkachenko A, Petmanson T, Laur S. Named Entity Recognition in Estonian. *Proceedings of the 4th Biennial International Workshop on Balto-Slavic Natural Language Processing*. 2010, p. 78–83. https://aclanthology.org/W13-2412.pdf (15.05.2024)

[59] Maide R. Eesti keele nimeolemite märgendaja analüüs ja pirandamine. Bachelor's thesis, University of Tartu, Institute of Computer Science. 2020. https://dspace.ut.ee/server/api/core/bitstreams/babadb6f-e74f-4f13-a365-a4baa8bed379/content (15.05.2024)

[60] Erelt M. Nimisõnafraasi sõnajärjest. *Oma Keel*, 2013, nr 26, lk 56–60.

[61] Merriam-Webster Adjective. An Encyclopaedia Britannica Company. https://www.merriam-webster.com/dictionary/adjective

[62] Schlechtweg M, Härtl H. Quotation marks and the processing of irony in English: evidence from a reading time study. *Linguistics*. 2023, vol. 61, p. 355–390. https://doi.org/10.1515/ling-2021-0079

[63] Erelt M, Erelt T, Ross K. Eesti keele käsiraamat. Tallinn: Eesti Keele Sihtasutus. 2020.

[64] Van Den Berg E, Markert K. Context in Informational Bias Detection. *Proceedings of the 28th International Conference on Computational Linguistics*. Barcelona: International Committee on Computational Linguistics. 2020, p. 6315–6326. https://doi.org/10.18653/v1/2020.coling-main.556

[65] Raadik M. Lugu „ärakaotatud" jutumärkidest. *Keelenõuanne soovitab 4*. Eesti Keele Instituut. Tallinn: Eesti Keele Sihtasutus. 2008, lk 164–183.

[66] Ehala M. Linguistic strategies and markedness in Estonian morphology. *Language Typology and Universals*. 2009, vol. 62, p. 29–48. https://doi.org/10.1524/stuf.2009.0003

[67] Argus R. Acquisition of Estonian: some typologically relevant features. *Language Typology and Universals*. 2009, vol. 62, p. 91–108. https://doi.org/10.1524/stuf.2009.0006

[68] Liivak M. KE(NE)-LIITELISED DEMINUTIIVID EESTI SUULISES ARGISUHTLUSES. Master's thesis, University of Tartu, Institute of Estonian and General Linguistics. 2023. https://dspace.ut.ee/server/api/core/bitstreams/a68f5d64-3921-4739-a2b0-687802644ba5/content (15.05.2024)

[69] Kasik R. Sõnamoodustus. Tartu: Tartu Ülikooli Kirjastus. 2015.

[70] Pai K. TRANSLATIIVNE JA ESSIIVNE PREDIKATIIVADVERBIAAL EESTI KIRJAKEELES. Master's thesis, University of Tartu, Institute of Philosophy. 2001. https://dspace.ut.ee/server/api/core/bitstreams/d4b3bb9d-8aca-4ad3-91bf-20d368c7f2a8/content (15.05.2024)

[71] Lee J, Pinker S. Rationales for Indirect Speech: The Theory of the Strategic Speaker. *Psychological Review*. 2010, vol. 117, p. 785–807. https://doi.org/10.1037/a0019688

[72] Van Vleet J. Informal Logical Fallacies: A Brief Guide. Lanham: University Press of America. 2011.

[73] Jurafsky D, Martin J. N-gram Language Models. *Speech and Language Processing*. 2024, vol. 3, p. 32–58. https://web.stanford.edu/~jurafsky/slp3/3.pdf (15.05.2024)

[74] Morstatter F, Wu L, Yavanoglu U, Corman S, Liu H. Identifying Framing Bias in Online News. *ACM Transactions on Social Computing*. 2018, vol. 1, p. 1–18. https://doi.org/10.1145/3204948

[75] Taboada M. Sentiment Analysis: An Overview from Linguistics. *Annual Review of Linguistics*. 2016, vol. 2, p. 325–347. https://doi.org/10.1146/annurev-linguistics-011415-040518

[76] Küçük D, Can F. A Tutorial on Stance Detection. *Proceedings of the Fifteenth ACM International Conference on Web Search and Data Mining*. 2022, p. 1626–1628. https://doi.org/10.1145/3488560.3501391

[77] Luukas R. Tartu Ülikooli õppeainete tagasiside meelsusanalüüs. Bachelor's thesis, University of Tartu, Institute of Computer Science. 2023. https://dspace.ut.ee/server/api/core/bitstreams/60fc63c2-3caa-4181-98cc-3293a737ab11/content (15.05.2024)

[78] Mohammad S, Turney P. Crowdsourcing a Word–Emotion Association Lexicon. *Computational Intelligence*. 2013, vol. 29, pp. 436–465. https://doi.org/10.1111/j.1467-8640.2012.00460.x

[79] Institute of the Estonian Language. https://eki.ee/EN/ (15.05.2024)

[80] Pajupuu H, Altrov R, Pajupuu J. Identifying Polarity in Different Text Types. *Folklore: Electronic Journal of Folklore*. 2016, vol. 64, p. 125–142. https://doi.org/10.7592/FEJF2016.64.polarity

[81] Devlin J, Chang M-W, Lee K, et al. BERT: Pre-training of Deep Bidirectional Transformers for Language Understanding. *arXiv preprint arXiv*. 2018. https://doi.org/10.48550/arXiv.1810.04805

[82] Tanvir H, Kittask C, Eiche S, Sirts K. EstBERT: A Pretrained Language-Specific BERT for Estonian. *Proceedings of the 23rd Nordic Conference on Computational Linguistics*

*(NoDaLiDa).* Reykjavik: Linköping University Electronic Press, Sweden. 2021, p. 11–19. https://aclanthology.org/2021.nodalida-main.2 (15.05.2024)

[83]   Barbieri F, Espinosa Anke L, Camacho-Collados J. XLM-T: Multilingual Language Models in Twitter for Sentiment Analysis and Beyond. *Proceedings of the Thirteenth Language Resources and Evaluation Conference*. Marseille: European Language Resources Association. 2022, p. 258–266. https://aclanthology.org/2022.lrec-1.27.pdf (15.05.2024)

[84]   Sadeniemi M, Kettunen K, Lindh-Knuutila T, Honkela T. Complexity of European Union Languages: A comparative approach∗. *Journal of Quantitative Linguistics*. 2008, vol. 15, p. 185–211. https://doi.org/10.1080/09296170801961843

[85]   Than K. Hungary's Orban says his anti-immigration stance not rooted in racism after backlash. *Reuters*, 2022. https://www.reuters.com/world/europe/hungarys-orban-says-his-anti-immigration-stance-not-rooted-racism-after-backlash-2022-07-28/ (15.05.2024)

[86]   Sõnaveeb. Lumehelbeke. https://sonaveeb.ee/search/unif/dlall/dsall/lumehelbeke/1

[87]   Dalianis H. Clinical Text Mining. Cham: Springer International Publishing. 2018.

# Appendix

## I.       Code Availability

The code written for the purposes of this thesis is publicly available in in the following GitHub repository: https://github.com/laurilyysi/EstonianStanceDetection.

## II. Licence

**Non-exclusive licence to reproduce the thesis and make the thesis public**

I, Lauri Lüüsi

1. grant the University of Tartu a free permit (non-exclusive licence) to

reproduce, for the purpose of preservation, including for adding to the DSpace digital archives until the expiry of the term of copyright, my thesis

**Political Stance Detection in Estonian News Media**

supervised by Uku Kangur and Roshni Chakraborty.

2. I grant the University of Tartu a permit to make the thesis specified in point 1 available to the public via the web environment of the University of Tartu, including via the DSpace digital archives, under the Creative Commons licence CC BY NC ND 4.0, which allows, by giving appropriate credit to the author, to reproduce, distribute the work and communicate it to the public, and prohibits the creation of derivative works and any commercial use of the work until the expiry of the term of copyright.

3. I am aware of the fact that the author retains the rights specified in points 1 and 2.

4. I confirm that granting the non-exclusive licence does not infringe other persons' intellectual property rights or rights arising from the personal data protection legislation.

*Lauri Lüüsi*
**15/05/2024**