

TARTU ÜLIKOOL  
Arvutiteaduse instituut  
Informaatika õppekava

Janar Saks

# Eestikeelsete tekstide sisukokkuvõtja EstSum edasiarendamine

Bakalaureusetöö (9 EAP)

Juhendaja: Kaili Müürisep, PhD

Tartu 2018

## **Eestikeelsete tekstide sisukokkuvõtja EstSum edasiarendamine**

### **Lühikokkuvõte:**

Tänapäevaste informatsioonihulkade juures on sageli vaja saada kiiresti ülevaade olulisest informatsioonist. Seepärast võiks rakenduse poolt automaatselt genereeritud sisukokkuvõte kui lühendatud kiiresti kättesaadav ülevaade algallikast olla oluline informatsiooni kogumise vahend. Kuid nagu iga keeletehnoloogiline rakendus, sõltub see sihtkeele omapäradest, mille jaoks see on disainitud. Inglise keele tarbeks loodud sisukokkuvõtja ei sobi eesti keele jaoks, sest eesti keelele omane sõnavormide rohkus vajab hoopis teistsugust lähenemist.

Kuigi eestikeelsete tekstide kokkuvõtja EstSum kasutab lause kaalu arvutamisel võtmesõnade põhise skoori, siis skoori arvutatakse sõnavormide, mitte sõnade algvormide pealt. EstSumi ühendamine lingvistilise mooduliga, mis suudab analüüsida sõnade algvorme, tõstis kokkuvõtja tulemuslikkust võrreldes EstSumiga, millel vastav moodul puudub. Töö tulemusena valminud automaatse kokkuvõtja uus versioon suudab eraldada rohkem olulist informatsiooni algallikast, kui seda tegi EstSumi vana versioon.

### **Võtmesõnad:**

Lemma, sisukokkuvõtja, EstSum, EstNLTK, lemmatiseerija, kokkuvõtja hindamine

**CERCS:** P175 Informaatika, süsteemiteooria

## **The Development of Estonian Texts' Summarizer EstSum**

### **Abstract:**

In today's vast information quantity, there is often a need for a quick overview of important information. Therefore, a summary as a shortened overview of the source material, could be an important source of information collection. But like any other language technology application, they also depend on the peculiarities of the language they are designed for. A summarizer created for the English language is not implementable for the Estonian language, because of the multitude of word forms that is typical for the Estonian language and therefore requires a completely different approach.

Although Estonian texts' summarizer EstSum uses a keyword-based score to calculate a sentence's weight, the score is calculated solely from word forms not the word's stem. The addition of a linguistic module, that can analyze word stems, did increase the evaluation score compared to the regular EstSum.

Furthermore, the new developed version of the Estonian texts' summarizer is capable of separating more important information from the source than the old version of EstSum.

### **Keywords:**

Word stem, summarizer, EstSum, EstNLTK, lemmatizer, summarization evaluation

**CERCS:** P175 Informatics, systems theory

# Sisukord

<b>1</b>	<b>Sissejuhatus</b>	<b>4</b>
<b>2</b>	<b>Eestikeelsete tekstide sisukokkuvõtja EstSum</b>	<b>5</b>
2.1	Sisukokkuvõtja . . . . .	5
2.2	Eesti keele sisukokkuvõtja . . . . .	5
2.3	EstSumi ülesehitus . . . . .	6
2.4	Lausete kaalu arvutamine . . . . .	6
2.5	Lausete valimine võrdluses SweSumiga . . . . .	7
<b>3</b>	<b>Sisukokkuvõtja tulemuslikkus</b>	<b>8</b>
3.1	Kokkuvõtja hindamine . . . . .	8
3.2	EstSumi tulemuslikkuse hinnangud . . . . .	9
<b>4</b>	<b>Lingvistilise mooduli ühendamine EstSumiga</b>	<b>10</b>
4.1	Lemmatiseerija rakendamine . . . . .	10
4.2	Sõnade algvormide sagedustabel ja stoppsõnad . . . . .	10
<b>5</b>	<b>Tulemuslikkuse hindamine</b>	<b>12</b>
5.1	Arendus- ja testkorpus . . . . .	12
5.2	Artiklite märgendamine . . . . .	13
5.3	Parameetrite uuendamine . . . . .	14
5.3.1	Parameeter $\alpha$ . . . . .	15
5.3.2	Parameeter $\beta$ . . . . .	16
5.3.3	Parameeter $\gamma$ . . . . .	16
5.3.4	Hinnang parameetritele . . . . .	16
5.4	Tulemuste kontrollimine . . . . .	17
5.5	Tulevikusuund . . . . .	18
<b>6</b>	<b>Kokkuvõte</b>	<b>19</b>
	<b>Viidatud kirjandus</b>	<b>21</b>
	<b>Lisad</b>	<b>22</b>
I.	Sõnade algvormide sagedustabel . . . . .	22
II.	Stoppsõnade loend . . . . .	29
III.	Näide algtekstis, märgendatud tekstist, käsitsi koostatud kokkuvõttest ja EstSumi genereeritud kokkuvõttest . . . . .	30
IV.	Parameetrite kolmikute keskmine kattuvus arneduskorpusel . . . . .	36
V.	Litsents . . . . .	38

# 1 Sissejuhatus

Tänapäevaste informatsioonihulkade juures on sageli vaja saada kiiresti ülevaade olulisest informatsioonist. Seepärast võiks rakenduse poolt automaatselt genereeritud sisukokkuvõtte kui lühendatud kiiresti kättesaadav ülevaade algallikast olla oluline informatsiooni kogumise vahend. Kokkuvõtte artiklitest, teadustöödest või dokumentidest võimaldab meil otsustada, kas soovime algtekstiga lähemalt tutvuda või saame selle kõrvale lükata. Samuti on väikestelt ekraanidelt ebamugav lugeda pikki tekste. Olukord muutuks oluliselt mugavamaks, kui oleks võimalik lühikese teksti põhjal otsustada, kas kogu teksti läbi töötamine on vajalik.

Automaatseid sisukokkuvõtjaid on loodud aastakümneid, neid on nii vabavaralisi kui ka kommertssüsteemide koosseisu kuuluvaid. Kuid nagu iga keeletehnoloogiline rakendus, sõltub see sihtkeele, millele see on disainitud, omapäradest. Inglise keele tarbeks loodud sisukokkuvõtja ei sobi eesti keele jaoks, sest eesti keelele omane sõnavormide rohkus vajab hoopis teistsugust lähenemist.

Eestikeelset sisukokkuvõtjat EstSum on arendatud mitu aastat. Pilleriin Mutso [P.05a] kirjeldas, kuidas EstSum määrab lausetele kaalu, mille põhjal koostatakse kokkuvõtte. Mutso leidis, et kaalu määramisel on tähtis lause positsioon, formaat ja lauses paiknevad võtmesõnad. Mutso kasutas võtmesõnade analüüsis sõnade erivorme. Kaili Müürisep ja Pilleriin Mutso [P.05b] kirjeldasid, kuidas määrati parameetrid, mis näitavad lause kaalus positsiooni, formaadi ja võtmesõnade tähtsuse osakaalu. Parameetrid määrati väikse korpuse põhjal, kuhu kuulus kakskümmend teksti, milles keskmiselt kaheksateist lauset. Käesoleva töö eesmärk on lisada eestikeelsele sisukokkuvõtjale EstSum lingvistiline moodul, mis võimaldaks analüüsida sõnade algvorme. Samuti on töö eesmärk uuele Estsumile leida sobivad parameetrite väärtused autori koostatud korpuse põhjal ja hinnata selle tulemuslikkust.

Käesolev töö koosneb neljast osast. Töö esimeses osas defineeritakse sisukokkuvõtja mõiste ja tuuakse välja, kuidas saab erinevaid kokkuvõtjaid liigitada. Samuti antakse ülevaade EstSumist ja võrreldakse EstSumi rootsikeelse sisukokkuvõtja SweSumiga, mis oli eeskujuks EstSumi loomisel. Töö teises osas kirjeldatakse, kuidas hinnatakse sisukokkuvõtja tulemuslikkust ja antakse ülevaade eelnevatest töödest, kus on hinnatud EstSumi funktsionaalsust. Töö kolmandas osas räägitakse, kuidas ühendati EstSumiga lingvistiline moodul ja kuidas koostati uus sõnade algvormide sõnasagedustabel ja stoppsõnade loend. Töö neljandas osas kirjeldatakse, kuidas koostati arendus- ja testkorpus. Samas töö osas antakse ülevaade ka tulemustest: kuidas leiti EstSumile uued parameetrite väärtused ja hinnatakse tulemuslikkust. Samuti tuuakse välja autori nägemus EstSumi arendussuunast.

## 2 Eestikeelsete tekstide sisukokkuvõtja EstSum

Selles peatükis antakse ülevaade, kuidas liigitatakse sisukokkuvõtjaid. Lisaks kirjeldatakse EstSumi arhitektuuri ning tuuakse täpsemalt välja, kuidas EstSum valib kokkuvõtetesse lauseid. Seejärel võrreldakse EstSumi lausete valimise protsessi teiste sisukokkuvõtjate lausete valimiste protsessidega.

### 2.1 Sisukokkuvõtja

Inderjeet Mani [I.01] annab sisukokkuvõtte laia definitsiooni: „Sisukokkuvõtja on süsteem, mille eesmärk on sisendist toota tihendatud kujutus, mis on mõeldud inimestele tarbimiseks.“ Definitsioonist saab järeldada, et sisukokkuvõtja eesmärk on algallikast eraldada oluline teave, nii et lugejal oleks sellest kiiresti hoomatav ülevaade.

Sisukokkuvõtjaid saab eristada selle järgi, kuidas need genereerivad väljundi [I.01]:

1. väljavõtte (ingl k *extract*) meetodit kasutavad sisukokkuvõtjad, mille väljund koosneb ainult sisendist kopeeritud osadest nt sisendteksti laused;
2. ülevaate (ingl k *abstract*) meetodit kasutavad sisukokkuvõtjad, mille väljundis mingi osa puudub sisendist nt sisendtekstist nime asendamine asesõnaga.

Samas, traditsiooniliselt on sisukokkuvõtjad jaotatud kaheks [I.01]:

1. osundav (ingl k *indicative*), kus kokkuvõtte annab ülevaate tervest sisendfailist;
2. informatiivne (ingl k *informative*), kus kokkuvõtte annab ülevaate kogu olulisest informatsioonist sisendfailis.

Käesolevas bakalaureusetöös käsitletav sisukokkuvõtja EstSum on osundav ja kasutab kokkuvõtete tegemiseks väljavõtte meetodit.

### 2.2 Eesti keele sisukokkuvõtja

EstSum on veebiuudistele ja elektroonilistele ajaleheartiklitele orienteeritud eestikeelne sisukokkuvõtja [K.06]. Kuigi EstSumi on arendatud mitu aastat, siis arendus on enamasti seotud ainult diplomi- ja bakalaureusetöödega [K.06].

Kuna EstSum kasutab kokkuvõtete genereerimiseks väljavõtte meetodit, siis genereeritud kokkuvõtte ei ole sidus [K.06]. EstSum on kirjutatud programmeerimiskeeles Perl ja koosneb kolmest moodulist [K.06].

Käesolevas töös valminud lähtekood uuest EstSumist on kättesaadav GitHubist<sup>1</sup>. Samuti on võimalik uut EstSumi katsetada veebis<sup>2</sup>.

<sup>1</sup>Koodi repositoorium: [https://github.com/janarsaks/EstSum\\_development](https://github.com/janarsaks/EstSum_development)

<sup>2</sup>EstSumi katsetamine: <http://kodu.ut.ee/~janar/EstSum2.0/>

## 2.3 EstSumi ülesehitus

Kaili Müürisep [K.06] järgi jaguneb EstSumi arhitektuur kolmeks osaks:

1. HTML-konverter, mille eesmärk on sisend viia SGML-formaati, kus sisendist on eemaldatud ebavajalikud elemendid ja lisatud vajalikud märgendid;
2. lausestaja, mis jagab regulaaravaldistega sisendi lauseteks;
3. lausete väljavalija, mis annab lausetele kaalu ja koostab kaalude järgi kokkuvõtte.

Käesolev töö keskendub EstSumi kolmandale osale, sest töö eesmärgiks on lisada lingvistiline moodul, mis võimaldaks lause kaalu hindamise protsessis arvestada sõnade algvormide esinemissagedust.

## 2.4 Lausete kaalu arvutamine

EstSum genereerib algtekstist kokkuvõtte, arvutades lausete kaalud ning nende tulemuste järgi valib laused, mille kaal on kõrgem. EstSumi lause kaalu arvutamise valem põhineb Edmundsoni paradigmat.

Inderjeet Mani [I.01] kirjeldas, kuidas Edmundson töötas teadusartiklite korpuse peal välja tingimused, mille abil saab määrata lause kaalu. See informatsioon aga võimaldas väljavalimismeetodi abil koostada algtekstist kokkuvõtte. Edmundson leidis, et lause kaalu mõjutavad märksõnad (ingl k *cue words*), võtmesõnad, lause positsioon ja pealkiri. Nende tingimuste alusel töötas Edmundson välja lause kaalu arvutamise valemi.

$$W(s) = \alpha C(s) + \beta K(s) + \gamma L(s) + \delta T(s) \quad (1)$$

Lause  $s$  kaalu leidmiseks tuleb liita kokku lause  $s$  märksõnade skoor  $C$ , võtmesõnade skoor  $K$ , positsiooni skoor  $L$  ja pealkirjaskoor  $T$ .  $\alpha$ ,  $\beta$ ,  $\gamma$  ja  $\delta$  on parameetrid, mille abil saab muuta tingimuse osakaalu lause kaalu arvutamises. Edmundson leidis, et lause positsioon oli parim tingimus ja võtmesõnad kõige nõrgem tingimus [I.01].

EstSum kasutab lause kaalu määramiseks kolme tingimust. Lause  $s$  kaal leitakse, kui liidetakse kokku lause  $s$  positsiooniskoor ( $P$ ), formaadiskoor ( $F$ ) ja sõnasageduste skoor ( $K$ ).

$$W(s) = \alpha P(s) + \beta F(s) + \gamma K(s) \quad (2)$$

$\alpha$ ,  $\beta$ ,  $\gamma$  on parameetrid, millega saab reguleerida, kui suur on kindla skoori osakaal. Lause positsioonipõhise skoori jaoks on tähtsad lõikude esimesed laused. Samas, kõige tähtsam on esimene lause, mis järgneb teksti pealkirjale.

Lause formaadipõhise skoori suurendab see, kui lause on kirjutatud paksus kirjas või kaldkirjas. Skoor langeb, kui lauses on jutumärgid või lause lõppeb hüüu- või küsimärgiga. Samuti väheneb skoor, kui on tegemist piltide all olevate tekstidega.

Lause sõnasageduste põhiskoori suurendavad sõnad, mis olid teksti pealkirjas. Skoor langeb, kui lauses on sõnad, mis paiknevad sõnavormide sagedustabelis või stoppsõnade loendis. Sagedustabelis oli 1057 sõnavormi ja stoppsõnade loendis 105 sõna. Pilleriin Mutso [P.05a] märkis, et sagedustabeli koostamiseks töödeldi läbi ajalehe „Postimees“ artiklitest koosnev korpus. Korpuses oli ligikaudu 390 000 sõna ning sagedustabelisse sisestati sõnavormide esinemissagedused 10 000 sõna kohta.

## 2.5 Lausete valimine võrdluses SweSumiga

EstSum on koostatud võttes eeskujuks sisukokkuvõtjat SweSum [K.06]. SweSum on rootsikeelne sisukokkuvõtja, mis on orienteeritud artiklite kokkuvõtmisele.

SweSum ja EstSum valivad kokkuvõtete jaoks lauseid erinevalt, kuigi mõlema kokkuvõtja jaoks on lause positsioon tähtis. SweSum annab lausetele skoori olenevalt sellest, millisel real nad paiknevad. Lause saab skoori vastavalt oma reanumbri pöördväärtusele [H.00]. Seevastu annab EstSum kõrgemaid skooresid lõikude esimestele lausetele. Mõlemad kokkuvõtjad peavad tähtsaks lauset, mis järgneb algteksti pealkirjale [M.03].

Mõlema kokkuvõtja jaoks on oluline lause formaat. Kui lause on paksus kirjas, suurendab see lause kaalu [H.00]. Erinevalt SweSumist suurendab EstSum lause kaalu, kui see on kaldkirjas ja vähendab lause kaalu, kui selles on jutumärgid või see lõppeb hüüu- või küsimärgiga.

Lause kaalu hindamisel uurivad mõlemad kokkuvõtjad ka sõnu. Erinevalt EstSumist suurendab SweSum lause kaalu, kui lauses leidub numbriline väärtus [H.00]. Ent mõlemad kokkuvõtjad peavad tähtsateks pealkirjas paiknevaid sõnu [M.03], mis tähendab, et lause kaal suureneb, kui selles paikneb mõni pealkirja sõna. Suureks erinevuseks on asjaolu, et SweSum analüüsib lause kaalu arvutamisel ka sõnade algvorme [H.00]. SweSum kasutab sõnade algvormide leidmiseks sõnastikku, kus võtmeteks on sõnade erivormid ning väärtusteks sõna algvorm [M.03]. Lisaks, erinevalt EstSumist peab SweSum tähtsateks lauseid, mis sisaldavad samu numbrilisi väärtusi ja arvestab kaalu määramisel lause pikkusega [M.03].

### 3 Sisukokkuvõtja tulemuslikkus

Selles peatükis kirjeldatakse, kuidas hinnata sisukokkuvõtja tulemuslikkust. Lisaks tuuakse ülevaade hinnangutest EstSumi tulemuslikkusele, milleni on jõudnud eelnevad tööd.

#### 3.1 Kokkuvõtja hindamine

Kokkuvõtete hindamine ei ole lihtne, sest pole võimalik objektiivselt väita, millised osad algtekstist on olulised ja millised mitte. Samuti pole olemas juhendit, mida järgides saaks anda hinnang kokkuvõtte ja seeläbi ka kokkuvõtja kohta. Üldiselt mõõdetakse kokkuvõtte hindamisel kahte tingimust [H.07]:

1. tihendus (ingl k *compression rate*) (CR), mis näitab, kui palju lühem on kokkuvõtte algtekstist;
2. andmepeetuse (ingl k *retention*) suhe (RR), mis näitab, kui palju informatsiooni on säilitatud.

Tihendus saadakse, kui kokkuvõtte pikkus jagatakse algteksti pikkusega. Käesolevas töös sellele suhtele ei keskenduta, sest EstSumi arendamiseks ja testimiseks automaatselt genereeritavad ja käsitsi tehtud kokkuvõtted moodustavad umbes 30% algteksti pikkusest.

$$CR = \frac{\text{kokkuvõttepikkus}}{\text{algtekstipikkus}} \quad (3)$$

Andmepeetuse suhe saadakse, kui kokkuvõttes paiknev informatsioon jagatakse algtekstis paikneva informatsiooniga. Seejuures tekib probleem, kuidas hinnata informatsiooni hulka kokkuvõttes või algtekstis.

$$RR = \frac{\text{informatsioonkokkuvõttes}}{\text{informatsioonalgtekstis}} \quad (4)$$

Kõige üldisemad meetodid, millega üritatakse anda hinnang andmepeetuse suhtele, saab jagada sisemisteks (ingl k *intrinsic*) ja välimisteks (ingl k *extrinsic*) [I.01]. Sisemised meetodid uurivad kokkuvõtte sidusust ja informatiivsust, mis tähendab, et hinnang antakse kokkuvõtte sisulisele kvaliteedile [H.07]. Selleks võrreldakse kokkuvõtet ideaaliga, mis üldjuhul koostatakse inimese poolt [I.01][H.07]. Välised meetodid hindavad kokkuvõtte asjakohasust ja loetavust, mis tähendab, et hinnang antakse kokkuvõtte vastuvõetavuse ja praktilisuse kohta [H.07].

EstSumi tulemuslikkust on ka varem hinnatud ning iga eelnev hinnang on antud kasutades sisemist hindamismeetodit, kus võrreldakse kokkuvõtja genereeritud kokkuvõtet käsitsi koostatud kokkuvõttega [P.05a][K.06][Sel08]. Seepärast kasutati käesolevas töös lingvistilise mooduliga EstSumi tulemuslikkuse hindamisel samasugust lähenemist.



## 3.2 EstSumi tulemuslikkuse hinnangud

EstSumi tulemuslikkusele on hinnangu andnud kolm autorit. Pilleriin Mutso [P.05a] leidis, et EstSum suutis ajaleheartiklist eraldada keskmiselt umbes 60% informatsiooni, mis oli Mutso jaoks oluline. Hinnangu andmiseks valis ta 20 ajalehe „Postimees“ artiklit ning koostas käsitsi nende põhjal kokkuvõtted. Artiklite keskmine pikkus oli 25 lauset. Järgmiseks häälestas ta EstSumi lause kaalu valemi parameetreid ning siis võrdles EstSumi genereeritud kokkuvõtteid käsitsi koostatud kokkuvõtetega.

Sama kõrge tulemuslikkuse hinnangu said ka Kaili Müürisep ja Pilleriin Mutso koostöös [P.05b]. Nad leidsid, et EstSum valis kokkuvõtteid genereerides keskmiselt 60% ulatuses samu lauseid, mis olid ka käsitsi koostatud kokkuvõtetes. Nende testkorpus koosnes 11-st artiklist, mille keskmine pikkus oli kakskümmend kolm lauset.

Kolmandana on EstSumi tulemuslikkust hinnanud Keili Sellik [Sel08]. Ta hindas EstSumi tulemuslikkust kahel viisil. Esimese viisi jaoks koostas ta viiekümnest artiklist, mille keskmine pikkus oli 10,16 lauset, käsitsi kokkuvõtte ja võrdles neid EstSumi genereeritud kokkuvõtetega. Sellik sai tulemuseks, et keskmiselt 65,29% ulatuses kattusid genereeritud ja käsitsi koostatud kokkuvõtted. Teise meetodina kasutas ta ROUGE programmi, et hinnata EstSumi tulemuslikkust. Kasutades programmi moodulit ROUGE-L, mis võrdleb lauseid, saadi tulemuseks, et EstuSumi kattuvus on 68,96%.

Eelnevad hinnangud EstSumi tulemuslikkusele saadi lühikeste artiklite põhjal. Käesolevas töös valiti korpustesse pikemaid artikleid. Peatükis viis kirjeldatakse, kuidas loodi arendus- ja testkorpus, mille uue versiooniga EstSumi tulemuslikkust hinnatakse. Arenduskorpus koosneb kahekümnest artiklist, mille keskmine pikkus oli 40,4 lauset. Testkorpus koosneb kümnest artiklist, mille keskmine pikkus oli 40,3 lauset.

## 4 Lingvistilise mooduli ühendamine EstSumiga

Selles peatükis kirjeldatakse, kuidas EstSumile lisati lingvistiline moodul, mis kokkuvõtete genereerimisel arvestab sõnade algvormidega. Samuti kirjeldatakse, kuidas koostati võtmesõnade analüüsiks vajalik sõnade sagedustabel, mis on lingvistilise mooduli edukaks kasutamiseks hädavajalik.

### 4.1 Lemmatiseerija rakendamine

EstSum oli algselt kirjutatud programmeerimiskeeles Perl. Et mugavamalt siduda programmi eesti keele keeletehnoloogiliste moodulitega, otsustati EstSumi lähtekood kirjutada ümber programmeerimiskeelde Python. Pythonis kirjutatud programmides on mugav kasutada teeki EstNLTK [HJ16], mis võimaldab EstSumil leida võtmesõnade algvorme ehk lemmasid.

EstSumi esialgne versioon kasutas võtmesõnade analüüsimiseks sõnavorme, mitte lemmasid. See tähendab, et sõna kaal sõltus sõna kirjapildist, mitte tähendusest. Näiteks kui artikkel rääkis metsandusest, siis sõnavormid mets, metsa, metsas, metsast jne loeti kõik erinevateks võtmesõnadeks. Selle puudujäägi eemaldamiseks lisati EstSumile juurde moodul, mis muudab lauses olevad sõnad nende algvormiks (eespool toodud näites lemmaks 'mets'). Selleks kasutati EstNLTK Text klassi lemmatiseerijat. Lemmatiseerija lisati meetoditesse „analyze\_title“, „analyze\_line“ ja „word\_based\_score“, kus analüüsitakse lausetes või pealkirjades paiknevaid sõnu. Töö eesmärgi täitmiseks kasutati Pythoni versiooni 3.5 ning EstNLTK versiooni 1.4.1.1.

EstNLTK lemmatiseerija kasutab sõnade algvormide leidmiseks Vabamorfi sõnastikupõhist morfoloogilist analüsaatorit [HJ16]. Vabamorfi morfoloogiline analüüs võrdleb sõnu sõnastikus paiknevate lekseemide kombinatsioonidega ja eemaldab sõnadelt liiteid ja lõppe, et kontrollida seda leksikoniga, kus paiknevad sõnatüved [T.16a]. Kui morfoloogiline analüüs annab tulemuseks mitu varianti, siis kasutab Vabamorf ühestajat, et valida õige ning kui ühestaja ei suuda õiget valida, väljastab analüüs mitu varianti [HJ16].

EstSum kasutab võtmesõnade analüüsimiseks sõnade sagedustabelit. Esialgses versioonis oli selleks sõnavormide sagedustabel. Uues versioonis loodi uus sõnade algvorme sisaldav sagedustabel. Uue versiooniga EstSumi lähtekoodiga saab tutvuda GitHubis<sup>3</sup>.

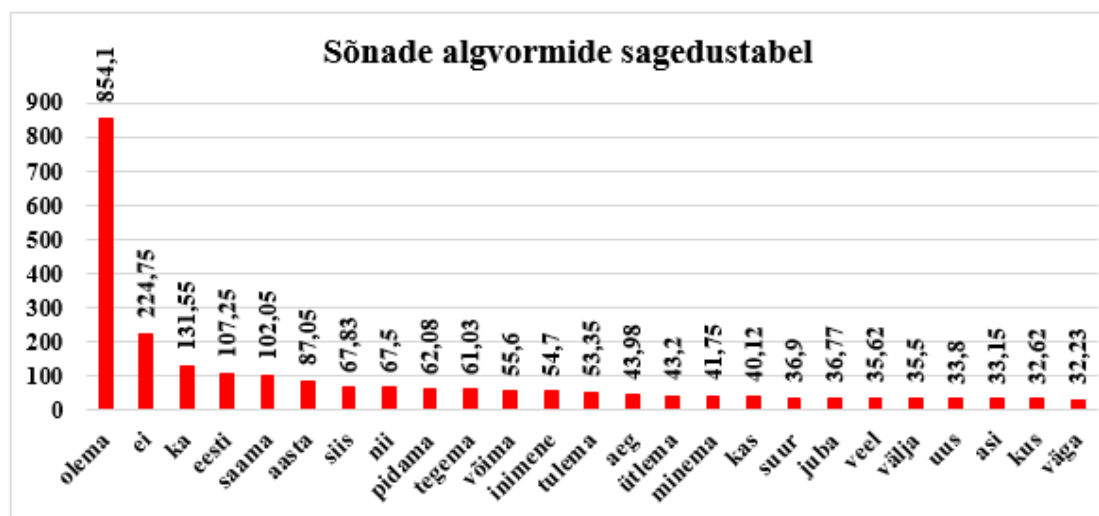
### 4.2 Sõnade algvormide sagedustabel ja stoppsõnad

EstSumi esialgne versioon kasutab võtmesõnade analüüsimiseks sõnade sagedustabelit ja stoppsõnade loendit. Kuid kuna nii sagedustabel kui ka loend koosnevad erinevates sõnavormidest, vajab uus, lingvistilise mooduliga ühendatud versioon võtmesõnade analüüsiks algvormidest koosnevat sagedustabelit ja stoppsõnade loendit. Seetõttu tuli

<sup>3</sup>Koodi repositoorium: [https://github.com/janarsaks/EstSum\\_development](https://github.com/janarsaks/EstSum_development)

koostada uus sõnade sagedustabel ning stoppsõnade loend, kus on ainult sõnade algvormid.

Sõnade algvormide sagedustabeli loomiseks kasutati Tartu ülikooli „Keeletehnoloogia“ kursuse eestikeelset korpust, kus on poole aasta jagu ajaleheartikleid. Korpus koosnes 894 737 sõnast. Korpus oli TEI-formaadis ning loeti sisse kasutades teicorpus klassi EstNLTK teegis. Esmalt korpuse tekstid lausestati ja seejärel lemmatiseeriti. Sõna algvormi sagedustabelis on toodud lemmade esinemissagedus 100 000 sõna kohta. Tabelisse lisati 1000 kõige kõrgema sagedusega lemmat. Samas välistati sõnade algvormid, mis paiknevad stoppsõnade loendis. Joonisel 1 on ära toodud 25 kõrgema sagedusega lemmat. Sagedustabeli sõnade loend ja sagedused on kogu suuruses toodud välja lisas I.



Joonis 1. 25 sagedamini esinevat sõna algvormi ajaleheartiklites.

Antud töös mõeldakse stoppsõnade all sõnu, millele antakse võtmesõnade analüüsi hetkel väärtuseks 0. See tähendab, et need sõnad ei mõjuta lause kaalu võtmesõnade skoori ehk need sõnad on tähenduse mõttes sisutühjad. Seepärast paiknevad uues stoppsõnade loendis ainult asesõnad ja sidesõnad. Nende loend saadi Eesti puudepangast ase- ja sidesõnade filtreerimise teel [T.16b]. Stoppsõnade loendis on 72 sõna ning need on välja toodud lisas II.

## 5 Tulemuslikkuse hindamine

Selles peatükis kirjeldatakse, kuidas koostati arendus- ja testkorpus, mis võimaldaks hinnata lingvistilise mooduliga ühendatud EstSumi tulemuslikkust. Samuti tuuakse välja, kuidas valiti uue mooduliga sobivad konstandid lause kaalu valemisse. Seejärel antakse hinnang uue mooduliga EstSumi tulemuslikkusele ja edastatakse autori nägemus, millises suunas peaks EstSumi arendamine tulevikus liikuma.

### 5.1 Arendus- ja testkorpus

Arenduskorpuse loomiseks valiti pisteliselt kaksikümne ajalehete „Postimees“ ja „Õpetajaleht“ veebiartiklit. Arenduskorpusesse valiti ainult artiklid, mille pikkus oli keskmisest kõrgem ehk umbes viissada kuni tuhat sõna. Esimesed kümme valitud artiklit olid arvamused ning ülejäänud kümme olid uudised. Pooled arvamused ja pooled uudised valiti ühest ajalehest ja ülejäänud teisest.

Testkorpuse loomisel lähtuti samadest tingimustest nagu arenduskorpuse loomisel. Kuid testkorpusesse valiti kõigest kümme artiklit, millest viis olid arvamused ja viis uudised. Arvamustest valiti kolm ajalehest „Õpetajaleht“ ja kaks ajalehest „Postimees“. Uudisetest valiti kaks ajalehest „Õpetajaleht“ ja kolm ajalehest „Postimees“.

Valitud artiklitest koostas autor väljavalmismetodit kasutades kokkuvõtted, mis moodustasid umbes 30% artikli kogupikkusest. Samas tuleb märkida, et kokkuvõtete tegemine on subjektiivne, sest ühe inimese tõlgendus artikli olulisest osast võib erineda teise inimese tõlgendusest. Seda näitas ka Martin Hassel [H.05], kes leidis, et kümne isiku koostatud artikli kokkuvõtetes oli igas kokkuvõttes ainult 33,9% sama informatsiooni. Tabelis 1 on ära toodud arendus- ja testkorpuse artiklite ja kokkuvõtete lausete ja sõnade arvud.

Tabel 1. Korpuses paiknevate artiklite ja kokkuvõtete pikkused

Arenduskorpus				
Artikkel	Laused	Sõnad	Kokkuvõtte laused	Kokkuvõtte sõnad
Arvamus1	52	666	15	198
Arvamus2	37	758	10	227
Arvamus3	51	897	16	270
Arvamus4	35	576	11	173
Arvamus5	43	816	12	239
Arvamus6	73	999	21	295
Arvamus7	46	711	15	214
Arvamus8	52	879	14	255
Arvamus9	34	578	11	176
Arvamus10	50	820	13	246
Uudis1	35	532	10	161

Uudis2	44	562	12	170
Uudis3	44	646	11	191
Uudis4	34	682	8	210
Uudis5	46	695	12	213
Uudis6	31	593	8	185
Uudis7	20	590	8	186
Uudis8	26	467	8	135
Uudis9	31	481	8	141
Uudis10	24	506	8	153
<b>Keskmine</b>	<b>40,4</b>	<b>672,7</b>	<b>11,55</b>	<b>201,9</b>
<b>Testkorpus</b>				
Arvamus1	61	958	20	289
Arvamus2	43	690	12	196
Arvamus3	44	535	12	164
Arvamus4	45	714	12	207
Arvamus5	57	960	19	278
Uudis1	36	530	9	156
Uudis2	31	535	9	153
Uudis3	22	498	8	141
Uudis4	32	514	10	148
Uudis5	32	741	9	217
<b>Keskmine</b>	<b>40,3</b>	<b>667,5</b>	<b>12,0</b>	<b>194,9</b>

Kuid korpuste koostamine ei piirdu ainult artiklitest kokkuvõtete tegemisega, sest kokkuvõtte genereerimiseks EstSumiga on vaja veel artiklid viia kindla märgendusega sisendfaili kujule. Lisas III on välja toodud näide artiklist, selle põhjal käsitsi koostatud kokkuvõttest ja EstSumi genereeritud kokkuvõttest. Korpuses paiknevad kokkuvõtted on kättesaadavad GitHubis<sup>4</sup>.

## 5.2 Artiklite märgendamine

Selleks, et viia artikkel EstSumi jaoks vajalikule sisendfaili kujule, on tarvis artikleid eeltöödelda. Eeltöötlusprotsess jagunes kaheks faasiks.

Eeltöötluse esimeses faasis tuli tekstis käsitsi ära märgendada kindlad teksti osad, et järgmises faasis oleksid need osad äratuntavad. Artikli teksti tuli lisada järgmised märgendid:

1. pealkiri pandi „<h1>“ ja „</h1>“ märgendite vahele;
2. alampealkirjad pandi „<h2>“ ja „</h2>“ märgendite vahele;

<sup>4</sup>Koodi repositoorium: [https://github.com/janarsaks/EstSum\\_development](https://github.com/janarsaks/EstSum_development)

3. paksus kirjas osa pandi „<b>“ ja „</b>“ märgendite vahele;
4. kaldkirjas osa pandi „<i>“ ja „</i>“ märgendite vahele.

Eeltöötuse teises faasis jagati artiklid lõikudeks ja lõigud lauseteks. Selleks kasutati EstNLTK Text klassi. See suudab teksti lõikudeks jagada, kuid eeldab, et lõigud on eraldatud kahekordse reavahetussümboliga<sup>5</sup>. Lõikude jagamisel lauseteks kasutati Text klassi lausestajat. Selline jagamine võimaldas lisada lõigud „<p>“ ja „</p>“ märgendite ning laused „<s>“ ja „</s>“ märgendite vahele. Samuti asendati eeltöötuse esimeses faasis lisatud märgendid sisendfaili jaoks sobilike märgenditega. Käsitsi lisatud märgendid võimaldasid algtekstis peal- ja alampealkirjad lisada „<head>“ ja „</head>“ märgendite vahele, kaldkiri lisada „<hi rend="italic">“ ja „</hi>“ märgendite vahele ja paks kiri lisada „<hi rend="bold">“ ja „</hi>“ märgendite vahele. Joonisel 2 on välja toodud kaks lõiku korrektselt märgendatud sisendfailist.

```
<div0><head><hi rend="bold">Riina Kabi: miks ei saa Eesti Punast Risti muuta projektipõhiseks?</hi></head>
<p>
<s><hi rend="bold">Eesti on alati pidanud tähtsaks elanike turvalisust.</hi></s>
<s><hi rend="bold">Enamik maailma riikidest on jõudnud äratundmisele, et turvalisuse tagamiseks on vaja veel midagi peale relvade ja kaitsejõudude.</hi></s>
</p>
<p>
<s>Tuleb hoolitseda haavatute eest, ühendada sõjas üksteist silmist kaotanud pereliikmeid, tegelda laste, orbude, vanurite ja invaliididega.</s>
<s>1864. aastal sõlmitud esimese Genfi konventsiooniga pandi nii alus humanitaarõigusele ja Punase Risti liikumisele.</s>
</p>
```

Joonis 2. Näide märgendatud tekstist.

Tuleb märkida, et Text klassi lausestaja teeb lausete eraldamisel vigu. Seetõttu tuli eeltöötuse tulemus käsitsi üle kontrollida. Lisas III on välja toodud näide märgendatud tekstist. Korpuses paiknevad märgendatud tekstid on kättesaadavad GitHubis<sup>6</sup>.

### 5.3 Parameetrite uuendamine

EstSumile, millega on ühendatud lingvistiline moodul, tuli leida uued parameetrid lause kaalu arvutamise valemis, sest eeldati, et uues versioonis töötab võtmesõnade arvestamise moodul efektiivsemalt. Valemis on kolm parameetrit ning sobivate kolmikute leidmiseks prooviti läbi kõik variandid, kus kolmiku summa andis tulemuseks ühe, sest nii oli võimalik analüüsida erineva parameetri osakaalu tähtsust.

Parameetrite väärtuste vahemikuks võeti 0,1 kuni 0,8. Kolmiku sobivuse hindamiseks genereeriti vastava kolmikuga kaksikümne kokkuvõtet artiklite põhjal, mis paiknesid arenduskorpuses. Iga genereeritud kokkuvõtet võrreldi käsitsi koostatud kokkuvõttega. Võrdluse aluseks oli kattuvus ning tulemuseks oli kattuvusprotsent. Kattuvuse hindamiseks vaadati läbi nii genereeritud kui ka käsitsi koostatud kokkuvõtete laused ning leiti

<sup>5</sup>EstNLTK Text klassi lähtekood. Kättesaadav: [https://estnltk.github.io/estnltk/1.4/\\_modules/estnltk/text.html](https://estnltk.github.io/estnltk/1.4/_modules/estnltk/text.html)

<sup>6</sup>Koodi repositoorium: [https://github.com/janarsaks/EstSum\\_development](https://github.com/janarsaks/EstSum_development)

laused, mis esinesid mõlemas kokkuvõttes. Kattuvusprotsendiks oli kattuvate lausete sõnade osakaal genereeritud kokkuvõtte kogupikkusest. Seega saadi iga kolmiku sobivuse hindamiseks kaksikümne kattuvusprotsenti ning nende abil arvutati keskmine kattuvus. Mida kõrgem oli keskmine kattuvus, seda sobilikum oli parameetrite kolmik. Tabelis 2 on toodud viis parameetrite kolmikut, mis saavutasid kõrgeimad kattuvusprotsendid. Tabel kõikide parameetrite kolmikutega on ära toodud lisas IV.

Tabel 2. Viis kõrgeimate keskmiste kattuvustega parameetrite kolmikut.

$\alpha$	0,4	0,5	0,4	0,4	0,7
$\beta$	0,4	0,4	0,3	0,5	0,2
$\gamma$	0,2	0,1	0,3	0,1	0,1
<b>Keskmine kattuvus</b>	62,24%	61,68%	61,65%	61,46%	61,39%

Tabelis 2 on välja toodud, et kõige kõrgema keskmise kattuvuse sai parameetrite kolmik 0,4, 0,4 ja 0,2. Kui analüüsiti tabelit 2, siis selgus, et lause kaalu arvutamise valemis on kõige olulisem parameeter  $\alpha$ , mis näitab lause positsiooni skoori osakaalu, ning kõige vähem olulisem on parameeter  $\gamma$ , mis näitab lause sõnasageduste põhise skoori osakaalu. Et saada parem ülevaade erinevate parameetrite olulisusest, vaadati iga parameetri puhul kõiki parameetrite kolmikuid, kus vaadeldava parameetri väärtus on suurem kui 0,5. See tähendab, et vaadeldi kolmikuid, kus üks parameeter moodustas üle poole lause valemi skoorist. Vastavaid kolmikuid oli igal parameetril kuus.

### 5.3.1 Parameeter $\alpha$

Parameeter  $\alpha$  näitab lause kaalu valemis, milline on lause positsiooni skoori osakaal. Kui Kaili Müürisep [K.06] otsis parameetritele väärtusi, siis leidis ta, et kõige olulisem neist kolmest oli parameeter  $\alpha$ . Tabelis 3 on välja toodud viis kõrgeimate keskmiste kattuvustega parameetrite kolmikut, kus  $\alpha$  parameetri väärtus on kõrgem kui 0,5.

Tabel 3. Parameetrite kolmikud, kus  $\alpha$  on suurem kui 0,5.

$\alpha$	0,7	0,8	0,6	0,7	0,6	0,6
$\beta$	0,2	0,1	0,3	0,1	0,2	0,1
$\gamma$	0,1	0,1	0,1	0,2	0,2	0,3
<b>Keskmine kattuvus</b>	61,39%	61,34%	61,32%	61,28%	61,27%	60,26%

Kui vaadeldi tabelis 2 keskmiseid kattuvusi ja võrreldi tabeliga 3, siis täheldati, et tulemused on väga sarnased. Samas on tabelist näha, et keskmise kattuvuse langusega on seotud parameetri  $\gamma$  suurenemine ja parameetri  $\alpha$  vähenemine.

### 5.3.2 Parameeter $\beta$

Parameeter  $\beta$  näitab lause kaalu valemis lause formaadi skoori osakaalu. Tabelis 4 on välja toodud 5 kõrgema keskmise kattuvusega parameetrite kolmikut, kus  $\beta$  parameetri väärtus on kõrgem kui 0,5.

Tabel 4. Parameetrite kolmikud, kus  $\beta$  on suurem kui 0,5.

$\alpha$	0,3	0,2	0,2	0,1	0,1	0,1
$\beta$	0,6	0,7	0,6	0,8	0,7	0,6
$\gamma$	0,1	0,1	0,2	0,1	0,2	0,3
<b>Keskmine kattuvus</b>	61,25%	60,74%	59,94%	58,55%	55,83%	51,63%

Kui vaadata tabeli 4 keskmiseid kattuvusi, siis on näha, et keskmise kattuvuse langus on seotud parameetri  $\alpha$  vähenemisega ja parameetri  $\gamma$  suurenemisega.

### 5.3.3 Parameeter $\gamma$

Parameeter  $\gamma$  näitab lause kaalu valemis, milline on lause sõnasageduste skoori osakaal. Tabelis 5 on välja toodud viis kõrgeimate keskmiste kattuvustega parameetrite kolmikut, kus  $\gamma$  parameetri väärtus on kõrgem kui 0,5.

Tabel 5. Parameetrite kolmikud, kus  $\gamma$  on suurem kui 0,5.

$\alpha$	0,3	0,2	0,2	0,1	0,1	0,1
$\beta$	0,1	0,2	0,1	0,2	0,3	0,1
$\gamma$	0,6	0,6	0,7	0,7	0,6	0,8
<b>Keskmine kattuvus</b>	55,02%	51,05%	50,4%	47,64%	47,46%	45,18%

Kui vaadeldi tabeli 5 keskmisi kattuvusi, siis on täheldati, et need on tunduvalt madalamad kui keskmised kattuvused tabelites 3 ja 4. Samuti tundus huvitav, et kolmiku puhul, kus  $\gamma$  väärtus on 0,8, on keskmine kattuvus kõige madalam. Nagu tabelite 3 ja 4 puhul, siis ka tabelil 5 kehtib eripära, et keskmine kattuvus langeb, kui parameeter  $\alpha$  väheneb ja parameeter  $\gamma$  suureneb.

### 5.3.4 Hinnang parameetritele

Lisas IV välja toodud tabel näitab, et kaheteistkümne parameetri kolmiku keskmine kattuvus oli suurem kui 60%. Nendes kolmikutes oli  $\alpha$  keskmine väärtus 0,5 ning  $\gamma$  keskmine väärtus 0,29. See tähendab, et positsioonipõhine skoor on kõige olulisem ning



sõnasageduste põhine skoor on kõige vähemolulisem. Väidet toetab ka asjaolu, et kui vaadati tabeleid 3, 4 ja 5, siis täheldati asjaolu, et kui positsioonipõhine skoor moodustab lause kaalust üle poole, siis saadud keskmised kattuvused olid kõige kõrgemad. Seevastu kui sõnasageduste põhine skoor moodustab lause kaalust üle poole, siis keskmised kattuvused olid kõige madalamad. Samuti kehtis tabelites 3, 4 ja 5 eripära, mis viitas, et keskmised kattuvused langesid, kui parameeter  $\alpha$  vähenes ja parameeter  $\gamma$  suurenes.

Seejärel leiti keskmine kattuvusprotsent arenduskorpuse artiklite kokkuvõtete ja EstSumi vanema versiooni genereeritud kokkuvõtete vahel. EstSumi vanem versioon kasutab lause kaalu valemis parameetreid 0,4, 0,4 ja 0,2. Keskmine kattuvusprotsent oli 58,08%. See viitas, et EstSumi vanema versiooni genereeritud kokkuvõtted sisaldavad vähem olulist informatsiooni kui tabelis 2 välja toodud kolmikutega genereeritud kokkuvõtted. Lisaks on EstSumi vanema versiooni keskmine kattuvus 4,16% madalam kui EstSumi uuemal versioonil, mis kasutab lause kaalu valemis samu parameetreid. Sellest järeldub, et lingvistilise mooduli ühendamine EstSumiga tõstab tulemuslikkust. Kindlasti tuleb aga mõõnda, et arenduskorpuses on kakskümmend artiklit ning käsitsi kokkuvõtete tegemine on subjektiivne. Samuti tuli saadud tulemusi kontrollida.

## 5.4 Tulemuste kontrollimine

Saadud tulemusi kontrolliti testkorpuse abil. Selleks kasutati testkorpuse artikleid, et leida käsitsi koostatud kokkuvõtete ja genereeritud kokkuvõtete kattuvusprotsendid ja arvutada nende põhjal keskmine kattuvusprotsent. Testkorpuse abil leiti keskmine kattuvus tabelis 2 väljatoodud kahe kõrgema keskmise kattuvusega parameetrite kolmikuga. Sama tehti ka tabelitega 3, 4 ja 5. Tabelis 6 on välja toodud saadud tulemused.

Tabel 6. Testkorpuse abil saadud keskmised kattuvused.

$\alpha$	$\beta$	$\gamma$	Keskmine kattuvus
0,4	0,4	0,2	53,91%
0,5	0,4	0,1	53,17%
0,7	0,2	0,1	50,35%
0,8	0,1	0,1	50,35%
0,3	0,6	0,1	53,81%
0,2	0,7	0,1	53,9%
0,3	0,1	0,6	51,26%
0,2	0,2	0,6	52,84%

Tabelist 6 on saab järeldada, et keskmine kattuvus oli kõrge parameetrite kolmikutel, kus parameeter  $\beta$  oli kõrge. Samuti seda, et parameetri  $\alpha$  tähtsus on langenud. Kui lähtuda ainult arenduskorpuse tulemustest, saanuks järeldada, et parameetri  $\alpha$  väärtus tagab

kõrge keskmise kattuvuse. Kuid testkorpuse tulemused seda väidet ei toeta, vaid näitavad, et ka parameeter  $\beta$  on tähtis. Kui lähtuda mõlema korpuse tulemustest, siis kõrgemad keskmised kattuvused saadi, kui parameetrid  $\alpha$  ja  $\beta$  olid osakaalult samatähtsad.

Seejärel leiti EstSumi eelmise versiooni kattuvusprotsendid ja keskmine kattuvusprotsent testkorpuse artiklitega, et saadud tulemusi võrrelda. Kui vaadata tulemusi parameetritega 0,4, 0,4 ja 0,2, siis EstSumi uuema versiooni keskmine kattuvus oli 53,91% ja EstSumi eelmise versiooni keskmine kattuvus oli 51,74%. Kuigi keskmiste kattuvuste erinevus vähenes 2,17%-ni, siis sellegipoolest olid EstSumi uuema versiooni genereeritud kokkuvõtted kõrgema keskmise kattuvusega. Arenduskorpuse abil saadud järeldus, et lingvistilise mooduli lisamine EstSumile tõstab tulemuslikkust, pidas paika ka testkorpuse puhul. Samuti tuleb märkida, et parameetrite kolmik 0,4, 0,4 ja 0,2 saavutas kõrgema keskmise kattuvuse arenduskorpuse artiklite seas. Lisaks saavutati kõrgeim keskmine kattuvus testkorpuse artiklite seas. Sellest järeldub, et keskmiselt kõige sobivam parameetrite kolmik EstSumi uuemale versioonile on 0,4, 0,4 ja 0,2. Ent siiski tuleb möönda, et käsitsi kokkuvõtete tegemine on subjektiivne ning üldiste järelduste tegemine oleks sobilikum suurema korpuse abil.

## 5.5 Tulevikusuund

Tehtud töö põhjal võib öelda, et lingvistilise mooduli lisamine EstSumile suurendab kokkuvõtja tulemuslikkust, kuid need tulemused on saadud väikeste korpuste pealt. Et saadud tulemusi kinnitada, tuleks tulevikus luua suurem korpus. Samuti tuleks korpuse loomisesse kaasata erinevaid inimesi, sest kui mitu inimest koostab ühele artiklile kokkuvõtte, on võimalik objektiivselt hinnata, mis on selles artiklis oluline enamiku inimeste jaoks. Sellise korpuse abil oleks võimalik täpsemalt leida lause kaalu valemile parameetrite väärtused, mis keskmiselt tagaksid kõrge kattuvusprotsendi.

Samuti tuleks korpusesse valida erinevat liiki artikleid, sest nii saaks kokkuvõtja parameetreid häälestada erinevat liiki artiklite jaoks. See võimaldaks aga EstSumile lisada artikliliigi seadistuse, millega saaks valida eelhäälestatud parameetrite kolmiku, mis keskmiselt tagaksid kõrge kattuvusprotsendi kindlale artikliliigile.

## 6 Kokkuvõte

Käesoleva töö eesmärgiks oli hinnata EstSumi tulemuslikkust, kui sellega on ühendatud lingvistiline moodul, mis võimaldaks kokkuvõtjal analüüsida sõnade algvorme.

Töös kirjeldatakse, kuidas kasutatakse EstNLTK lemmatiseerijat, et EstSumile lisada lingvistiline moodul. Lisaks on näidatud, kuidas koostati korpused, eesmärgiga leida EstSumi uuele versioonile uued parameetrid ja hinnata selle tulemuslikkust.

Töö tulemusena valmis uus versioon EstSumist, mis võtab arvesse sõnade algvorme. Samuti leiti, et kui võrrelda vana EstSumi genereeritud kokkuvõtet uue EstSumi genereeritud kokkuvõttega, siis uus EstSum eraldab rohkem autori jaoks olulist informatsiooni. EstSumi uus versioon eraldas arenduskorpuse peal 4,16% ja testkorpuse peal 2,17% rohkem olulist informatsiooni kui EstSumi vana versioon. Samuti leiti, et EstSumi uuele versioonile on kõige sobivam parameetrite kolmik 0,4, 0,4 ja 0,2. Samas tuleb mainida, et korpuste tegemine oli subjektiivne ja kolmekümne kokkuvõtte põhjal ei saa teha liiga üldiseid järeldusi.

EstSumi edasipidine arendussuund võiks olla seotud suurema korpuse loomisega, mille põhjal oleks võimalik anda üldisem järeldus EstSumi tulemuslikkuse kohta ja peenhäälestada parameetrite väärtuseid nii, et kasutajal oleks võimalik valida parameetrite väärtused vastavalt artikli liigile.

## Viidatud kirjandus

- [H.00] Dalianis H., 2000. SweSum - a text summarizer for Swedish. KTH-Stockholm. <ftp://ftp.nada.kth.se/IPLab/TechReports/IPLab-174.pdf> (08.05.2018).
- [H.05] Hassel M., Dalianis H., 2005. Generation of Reference Summaries. KTH-Stockholm. <https://pdfs.semanticscholar.org/a374/e9a8dbd0c395f28795625c8258e73bdbbb1a.pdf> (14.05.2018).
- [H.07] Hassel H. *Resource Lean and Portable Automatic Text Summerization*. PhD thesis, KTH School of Computer Science and Communication, 2007.
- [HJ16] Orasmaa S., Petmanson T., Tkachenko A., Laur S., Kaalep H.-J., 2016. EstNLTK-NLP Toolkit for Estonian. Tartu Ülikool, Arvutiteaduse Instituut. <https://pdfs.semanticscholar.org/d834/13f7f785aae067d49332239c3a36c346ba99.pdf> (13.05.2018).
- [I.01] Mani I. *Automatic Summarization*. Amsterdam: John Benjamins Publishing Company, 2001.
- [K.06] Müürisep K. Eestikeelsete tekstide sisukokkuvõtjast estsum. *Keel ja Arvuti*, Tartu Ülikooli üldkeeleteaduse õppetooli toimetised 6, 2006.
- [M.03] Dalianis H., Hassel M., 2003. From SweSum to ScandSum - Automatic text summarization for the Scandinavian languages. KTH-Stockholm. <https://people.dsv.su.se/~hercules/scandsum/ScandSumArsbog2002.pdf> (10.05.2018).
- [P.05a] Mutso P. Automaatne sisukokkuvõtete tegemine eestikeelsetest ajalehetekstidest: parameetrid ja hindamine, 2005. Diplomitöö. Tartu Ülikool, Arvutiteaduse Instituut.
- [P.05b] Müürisep K., Mutso P. Estsum - estonian newspaper texts summarizer. Proceedings of The Second Baltic Conference on Human Language Technologies, Tallinn 2005.
- [Sel08] Keili Sellik, 2008. Automaatse sisukokkuvõtja töö hindamine. Tartu Ülikool, Arvutiteaduse Instituut. [http://lepo.it.da.ut.ee/~kaili/juhendamised/Baka\\_Sellik.pdf](http://lepo.it.da.ut.ee/~kaili/juhendamised/Baka_Sellik.pdf) (12.05.2108).

- [T.16a] Kaalep H.-J., Vaino T. Complete morphological analysis in the linguist's toolbox. In *Congressus Nonus Internationalis Fenno-Ugristarum Pars V*, pages 9–16, 2016.
- [T.16b] Muischnek K., Müürisep K., Puolakainen T. Estonian dependency treebank: from constraint grammar tagset to universal dependencies. Proceedings of LREC, 2016.  
[http://www.lrec-conf.org/proceedings/lrec2016/pdf/411\\_Paper.pdf](http://www.lrec-conf.org/proceedings/lrec2016/pdf/411_Paper.pdf) (10.05.2018).

## Lisad

### I. Sõnade algvormide sagedustabel

olema 854.10	üle 28.80	näide 18.88	raamat 13.60
ei 224.75	kaks 28.30	ju 18.87	siiski 13.53
ka 131.55	ära 28.10	leidma 18.27	euroopa 13.40
eesti 107.25	andma 27.62	isegi 18.18	vana 13.40
saama 102.05	esimene 26.93	peale 18.12	kohta 13.40
aasta 87.05	tallinn 26.45	läbi 17.75	mõte 13.10
siis 67.83	eest 26.02	siin 17.72	vahel 12.97
nii 67.50	töö 25.18	firma 17.65	silm 12.93
pidama 62.08	elu 25.18	miljon 17.65	muidugi 12.40
tegema 61.03	kord 24.47	vaatama 16.95	maksma 12.40
võima 55.60	nägema 23.97	osa 16.90	mõtlemine 12.38
inimene 54.70	küll 23.63	ikka 16.75	jaoks 12.32
tulema 53.35	kroon 23.30	kokku 16.52	küsimus 12.25
aeg 43.98	kõige 22.70	keel 16.50	pilt 12.22
ütlemine 43.20	ainult 22.63	kolm 16.27	võimalus 12.20
minema 41.75	kuidas 22.55	juht 16.20	tartu 12.20
kas 40.12	nüüd 22.05	käsi 15.80	tundma 11.92
suur 36.90	tagasi 21.53	pool 15.75	lõpp 11.85
juba 36.77	maailm 21.15	kirjutama 15.75	lugema 11.78
veel 35.62	käima 21.02	i 15.70	linn 11.35
välja 35.50	raha 20.98	tegelikult 15.67	kunagi 11.30
uus 33.80	naine 20.90	elama 15.40	tooma 11.22
asi 33.15	viimane 20.77	auto 15.10	küsima 11.22
kus 32.62	sõna 20.73	nimi 15.05	suurem 11.20
väga 32.23	panema 20.60	maja 14.73	hästi 11.18
mitte 32.15	arvama 20.27	kasutama 14.63	alla 11.10
mees 32.15	enam 20.18	koos 14.55	eriti 11.08
tahtma 32.00	vastu 20.15	enne 14.50	pea 11.05
pärast 30.15	päev 20.13	lihtsalt 14.45	kohe 11.05
teadma 30.10	kogu 19.82	miks 14.43	võimalik 11.05
rääkima 30.07	rohkem 19.70	vene 14.30	eestlane 11.00
võtma 30.05	seal 19.15	noor 14.28	nädal 10.98
hakkama 29.55	laps 19.07	praegu 14.25	soome 10.97
jääma 29.48	riik 18.97	ette 13.68	all 10.97
hea 29.33	just 18.90	jõudma 13.62	ostma 10.88

hiljem 10.70	kuu 9.15	lubama 7.72	täiesti 6.88
rahvas 10.70	toimuma 9.00	vanem 7.70	kaasa 6.88
näitama 10.70	uskuma 8.90	ilma 7.67	jooksul 6.83
töötama 10.67	mängima 8.88	löpuks 7.67	täis 6.82
maa 10.58	järgmine 8.85	endine 7.67	kindel 6.80
sisse 10.55	protsent 8.82	tulemus 7.65	ühiskond 6.78
president 10.45	laskma 8.80	sajand 7.60	huvi 6.78
kunst 10.33	tähendama 8.75	isa 7.58	teema 6.77
maha 10.32	järgi 8.65	kindlasti 7.55	päris 6.75
parem 10.30	lisa 8.65	kinni 7.47	selge 6.72
paar 10.30	kümme 8.62	meenutama 7.45	ajalugu 6.72
viis 10.30	kiri 8.62	varem 7.40	pigem 6.72
pakkuma 10.22	vaja 8.57	kool 7.40	kohalik 6.72
film 10.22	samuti 8.55	meeldima 7.38	kuulma 6.72
hoopis 10.20	ühe 8.53	valitsus 7.37	kandma 6.70
alati 10.15	probleem 8.52	oskama 7.30	tegelema 6.60
suutma 10.10	ees 8.50	samas 7.28	kutsuma 6.60
edasi 10.07	üldse 8.47	plaat 7.28	ilmuma 6.60
poolt 9.98	saatma 8.40	õige 7.27	vähem 6.58
tunduma 9.95	jutt 8.35	idee 7.25	jälle 6.58
viima 9.93	muusika 8.30	neli 7.23	algama 6.57
liiga 9.93	hoidma 8.28	ekspress 7.20	poole 6.55
kodu 9.90	kuuluma 8.22	erinev 7.17	meel 6.55
müüma 9.88	ootama 8.12	the 7.12	teater 6.52
sõber 9.78	põhjus 7.97	kinnitama 7.12	vastama 6.50
juurde 9.75	asuma 7.93	püüdma 7.12	väike 6.50
näitus 9.75	alles 7.92	nõukogu 7.12	valge 6.47
sõitma 9.70	autor 7.92	enamik 7.10	koht 6.35
kunstnik 9.57	jätma 7.92	üsna 7.08	tüdruk 6.35
eelmine 9.55	ilmselt 7.90	hind 7.03	nimetama 6.33
algus 9.47	istuma 7.90	vähemalt 7.03	ainus 6.32
tegemine 9.40	ajama 7.87	sinna 7.00	saksamaa 6.30
umbes 9.30	hetk 7.85	juures 7.00	ajal 6.25
tee 9.30	parim 7.85	üumber 7.00	läti 6.23
üles 9.25	usa 7.83	pank 6.98	rootsi 6.22
aru 9.22	muutama 7.82	poiss 6.97	taga 6.20
ameerika 9.20	pikk 7.80	mart 6.93	andres 6.18
venemaa 9.20	ülikool 7.80	otsima 6.92	ema 6.18
lugu 9.18	raske 7.77	tekkima 6.90	kell 6.15
puhul 9.17	tund 7.75	raudtee 6.90	tänav 6.12

astuma 6.10	must 5.50	esitama 4.92	eriline 4.53
seisma 6.10	asemel 5.50	kirjanik 4.90	näitleja 4.53
häääl 6.08	otsus 5.47	müük 4.90	puuduma 4.53
eesmärk 6.05	õhtu 5.47	politsei 4.90	ajaleht 4.50
täpselt 6.05	otsustama 5.47	kasvama 4.90	paistma 4.50
oluline 6.05	lihtne 5.45	tulevik 4.88	kuidagi 4.50
toomas 6.05	selguma 5.45	jüri 4.88	loom 4.50
terve 6.00	ruum 5.45	kuhu 4.85	ajakiri 4.48
õppima 6.00	vesi 5.42	äri 4.85	jõud 4.47
vist 5.98	uks 5.40	sattuma 4.85	itaalia 4.45
meri 5.95	nägu 5.38	võrdlema 4.83	soovima 4.45
vein 5.95	kolmas 5.35	lõpetama 4.83	a 4.45
rahvusvaheline 5.92	teos 5.35	palk 4.82	lootma 4.45
lahti 5.90	praegune 5.32	üritama 4.80	tuhat 4.45
olukord 5.90	moskva 5.32	artikkel 4.80	avalik 4.42
nõudma 5.90	roll 5.25	kuulama 4.80	veidi 4.40
aitama 5.88	kohus 5.20	leht 4.78	mäng 4.40
tunnistama 5.82	tõttu 5.20	seega 4.78	märts 4.40
kultuur 5.80	hulk 5.20	võim 4.77	hommik 4.38
tänapäev 5.78	juhtuma 5.15	avaldama 4.75	new 4.38
tekst 5.77	teatama 5.15	dollar 4.75	internet 4.38
valima 5.75	mööda 5.12	näiteks 4.75	hoone 4.38
ikkagi 5.73	uurima 5.12	loom 4.75	tohtima 4.35
liige 5.73	pisut 5.12	tähtis 4.75	järele 4.35
muutma 5.70	peaaegu 5.10	surm 4.72	hulgas 4.33
turg 5.70	teenima 5.10	õigus 4.70	poliitik 4.33
täna 5.70	ometi 5.08	tuttav 4.70	kust 4.33
valmis 5.67	plaan 5.08	arvuti 4.67	saksa 4.33
tegevus 5.65	tunne 5.07	vastus 4.67	seni 4.33
ilus 5.65	leping 5.05	poliitiline 4.65	esinema 4.33
väitma 5.65	kaudu 5.05	vaba 4.65	tõeline 4.30
projekt 5.62	jalg 5.00	mood 4.60	paluma 4.30
tõesti 5.53	kõrval 5.00	inglise 4.60	natuke 4.28
mõistma 5.53	peeter 4.98	alustama 4.58	märk 4.28
ajakirjanik 5.53	tegelane 4.97	vähe 4.58	piir 4.27
ligi 5.53	huvitav 4.97	direktor 4.58	kartma 4.25
venelane 5.53	ettevõte 4.95	esi 4.57	sööma 4.25
töötaja 5.52	tavaline 4.95	suhe 4.55	s 4.23
omanik 5.52	liit 4.95	pärit 4.55	peal 4.22
suurim 5.50	jah 4.93	ärimees 4.53	klient 4.20



selg 4.18	materjal 3.88	juut 3.70	õnn 3.43
stiil 4.17	seletama 3.88	kiiresti 3.70	piisavalt 3.42
sada 4.17	korraldama 3.88	ringi 3.70	holland 3.42
partei 4.17	kuus 3.85	kolleeg 3.68	seisukoht 3.40
kujutama 4.17	seetõttu 3.85	eri 3.68	isik 3.40
tekitama 4.17	foto 3.85	isiklik 3.68	laul 3.40
loomulikult 4.15	muidu 3.83	album 3.67	suhtes 3.40
andmed 4.15	olev 3.83	maal 3.67	seadus 3.40
publik 4.15	tähelepanu 3.83	keha 3.67	väiksem 3.38
mäletama 4.15	üha 3.83	tegu 3.65	pärnu 3.38
perekond 4.10	looming 3.83	peaminister 3.65	riigikogu 3.38
korter 4.10	eks 3.82	tõmbama 3.65	sõda 3.38
pere 4.10	muuseum 3.82	koer 3.63	sügis 3.38
number 4.08	põlvkond 3.80	tulviste 3.62	rida 3.38
sobima 4.08	seas 3.80	selgitama 3.62	kohale 3.38
lava 4.08	sealt 3.80	minut 3.62	uudis 3.38
võibolla 4.05	jooma 3.80	punane 3.60	välismaa 3.37
miljard 4.05	valik 3.80	jne 3.60	surema 3.35
abi 4.05	suvi 3.78	sein 3.60	vorm 3.35
pääsema 4.03	kõlama 3.78	avama 3.57	arst 3.35
soov 4.03	vahele 3.77	lugeja 3.57	kavatsema 3.35
aprill 4.00	no 3.77	sees 3.55	meeter 3.35
paremini 4.00	taas 3.77	hiljuti 3.55	tühi 3.35
liikuma 4.00	ala 3.77	sageli 3.55	prantsuse 3.35
arvamus 4.00	esimees 3.75	saabuma 3.55	etendus 3.35
nimelt 3.98	keskmine 3.75	kuju 3.55	kasu 3.33
halb 3.98	tõstma 3.75	poeg 3.55	hotell 3.33
ehitama 3.98	lisama 3.75	mulje 3.55	tütar 3.32
jaanuar 3.98	puhas 3.73	york 3.55	mil 3.30
veebruari 3.95	enamasti 3.73	kogemus 3.55	poliitika 3.30
jaan 3.95	jagama 3.73	tase 3.53	tõsine 3.27
lööma 3.95	tõusma 3.73	arhitekt 3.53	sakslane 3.27
preemia 3.95	london 3.73	kõrge 3.52	pealt 3.25
näima 3.95	kaduma 3.72	estonia 3.50	areng 3.25
kuulus 3.95	juhtima 3.72	tugev 3.50	soovitama 3.25
tavaliselt 3.93	kirjandus 3.72	meelest 3.48	aktsia 3.25
sündima 3.93	vajama 3.72	pidevalt 3.45	korralik 3.25
lausa 3.93	mark 3.70	otse 3.45	klass 3.25
mõjuma 3.90	võitma 3.70	riie 3.43	esindaja 3.25
ajakirjandus 3.88	tegija 3.70	külg 3.43	leedu 3.25

lahkuma 3.25	kirjastus 3.08	meeskond 2.88	otsa 2.75
tükk 3.23	värske 3.08	protsess 2.88	nimel 2.75
tõenäoliselt 3.22	ooper 3.08	kohtuma 2.88	piirkond 2.75
vahe 3.20	minister 3.08	pudel 2.88	riia 2.75
juhtum 3.20	tihti 3.07	järjest 2.87	tehas 2.75
alates 3.20	viga 3.05	kõvasti 2.85	lahendus 2.75
elav 3.20	uuesti 3.05	ajalooline 2.85	elanik 2.73
sõnum 3.20	toit 3.05	tänu 2.85	klubi 2.73
reklaam 3.20	paik 3.05	kõne 2.85	pind 2.73
sadam 3.20	saade 3.02	saal 2.85	haigus 2.73
seltskond 3.18	kelam 3.02	leiduma 2.83	pealkiri 2.73
meedia 3.18	lavastus 3.02	hans 2.83	kaup 2.73
järel 3.18	kaotama 3.02	masin 2.83	päritolu 2.73
veelgi 3.18	suu 3.02	mets 2.83	siinne 2.72
vend 3.18	abil 3.00	tähendus 2.83	info 2.72
seitse 3.18	kiire 3.00	hindama 2.82	vedama 2.72
amet 3.17	hoolimata 3.00	juhataja 2.82	mullu 2.70
sarnane 3.15	süda 3.00	luule 2.82	korral 2.70
intervjuu 3.15	kontsert 3.00	avastama 2.82	laar 2.70
jumal 3.15	ülesanne 3.00	laulma 2.82	savisaar 2.70
süsteem 3.15	armastama 3.00	algul 2.82	pood 2.68
siia 3.15	pilk 2.98	kaheksa 2.80	õnnestuma 2.68
osalema 3.15	mägi 2.98	tuba 2.80	sotsiaalne 2.67
kallis 3.15	jälgima 2.98	valimine 2.80	uuring 2.67
jätkuma 3.15	seepärast 2.97	üritus 2.80	anne 2.67
vajalik 3.15	telefon 2.95	koostöö 2.80	paber 2.67
hansapank 3.15	erastamine 2.95	lääs 2.80	tiina 2.67
moodne 3.15	viskama 2.93	ammu 2.80	moodustama 2.67
kirik 3.15	sellepärast 2.92	vajadus 2.80	proovima 2.65
haigla 3.13	mõis 2.92	nali 2.78	kõrvale 2.65
väärtus 3.13	tegutsema 2.92	abikaasa 2.78	vahetama 2.65
kaua 3.13	taha 2.92	kerge 2.78	käigus 2.65
nr 3.13	paraku 2.92	tüüp 2.78	märkima 2.65
aken 3.12	laud 2.90	poolest 2.77	kalev 2.65
nn 3.10	sisu 2.90	võõras 2.77	tasuma 2.65
ameeriklane 3.10	sündmus 2.90	väide 2.77	ametnik 2.62
helistama 3.10	asuv 2.90	august 2.77	odav 2.62
eelkõige 3.10	tore 2.90	ühte 2.75	etv 2.62
täitma 3.10	arvestama 2.90	suhteliselt 2.75	jooksma 2.62
vabariik 3.10	kirjeldama 2.90	märkama 2.75	avalikkus 2.62

korrus 2.62	vabadus 2.50	majandus 2.38	grupp 2.30
lavastaja 2.62	hulka 2.50	linnapea 2.38	vahepeal 2.30
valmistama 2.62	nokia 2.50	loobuma 2.38	omavahel 2.30
huvitama 2.60	toetama 2.50	vale 2.38	suund 2.30
kõva 2.60	keeruline 2.50	advokaat 2.38	seadma 2.30
siit 2.60	of 2.50	kõigepealt 2.38	indrek 2.27
klaas 2.60	tingimus 2.50	mari 2.38	osutuma 2.27
rahvuslik 2.60	de 2.50	õnnelik 2.38	varsti 2.27
tarmo 2.60	võistlus 2.48	unt 2.38	is 2.27
värv 2.60	sposato 2.48	sisaldama 2.38	toode 2.25
mihkel 2.60	variant 2.48	konkreetne 2.37	barbie 2.25
õpetaja 2.60	inglismaa 2.48	võimaldama 2.35	investor 2.25
lai 2.58	sundima 2.47	riiklik 2.35	käik 2.25
tellima 2.58	suurepärane 2.47	tänane 2.35	kaunis 2.25
saar 2.58	narva 2.47	endiselt 2.35	lühike 2.25
kohaselt 2.58	edu 2.45	lähedal 2.35	õhk 2.25
maitse 2.58	öö 2.45	õpetama 2.35	keskus 2.25
alt 2.58	koguma 2.45	detsember 2.35	kodanik 2.25
reegel 2.58	külm 2.45	vaataja 2.35	normaalne 2.23
paul 2.58	von 2.43	vähene 2.35	professor 2.23
romaan 2.58	noormees 2.43	park 2.35	lk 2.23
õde 2.58	eelistama 2.43	julgema 2.33	mujal 2.23
restoran 2.58	teadlane 2.43	äkki 2.33	lause 2.23
viin 2.57	mootor 2.42	keskel 2.33	rikas 2.23
korruga 2.57	mõju 2.42	katse 2.33	armastus 2.22
punkt 2.57	kukkuma 2.42	põhjal 2.33	peamiselt 2.22
tõepoolest 2.55	komme 2.42	tähtsam 2.33	teenus 2.20
hiina 2.55	valgus 2.42	nõunik 2.33	pealinn 2.20
kasutamine 2.55	haridus 2.40	keerama 2.32	arhitektuur 2.20
loodus 2.53	tõsiselt 2.40	edukas 2.32	tarvis 2.20
häda 2.52	koguni 2.40	vaene 2.32	tunduvalt 2.20
seejärel 2.52	kestma 2.40	kriitik 2.32	määrama 2.20
langema 2.52	kilomeeter 2.40	ots 2.32	valitsema 2.20
kuulutama 2.52	juhatus 2.40	pealegi 2.32	staar 2.20
lennuk 2.52	sõltuma 2.40	tänavu 2.30	hinnang 2.20
veri 2.52	tootmine 2.40	kummaline 2.30	n 2.20
taust 2.52	samm 2.40	noorem 2.30	karu 2.20
arv 2.52	soe 2.40	unustama 2.30	fotograaf 2.18
lennart 2.50	väljas 2.40	võitlus 2.30	uurimine 2.17
soomlane 2.50	tasu 2.38	tiit 2.30	oht 2.17

luuletus 2.17	prantsusmaa 2.15	jätkama 2.10	niivõrd 2.08
t 2.17	tootma 2.15	saladus 2.10	ilves 2.08
kätte 2.17	jürgenson 2.15	seekord 2.10	paks 2.08
erakond 2.17	kuluma 2.13	laev 2.10	klassikaline 2.08
linnavalitsus 2.17	lökkama 2.13	küla 2.10	mure 2.08
rein 2.17	raadio 2.12	kõrgem 2.10	tõnu 2.08
st 2.17	juuni 2.12	vaatamata 2.10	kangelane 2.08
paavo 2.17	suhtlema 2.12	pikem 2.10	siiani 2.08
kohtumine 2.15	õpilane 2.12	helsingi 2.10	saavutama 2.07
nina 2.15	itaallane 2.12	magama 2.10	kujunema 2.07
vaim 2.15	ühtlasi 2.12	esindama 2.10	objekt 2.07
lootus 2.15	kahjuks 2.12	kuum 2.10	süüdistama 2.07
ohver 2.15	urmas 2.12	naerma 2.10	osakond 2.07
kommenteerima	hull 2.10	teadmine 2.08	
2.15	kvaliteet 2.10	lõppema 2.08	
ühine 2.15	tegelik 2.10	kohal 2.08	

## II. Stoppsõnade loend

aga	igasugune	omaenese
ega	igäüks	omasugune
ehk	ise	palju
elik	iseenese	säärane
ent	keegi	sama
ja	kes	samasugune
kui	kõik	see
kuid	kumbki	seesama
kuni	milline	seesamune
nagu	mina	selline
ning	mingi	sihuke
vaid	mingisugune	sina
või	mis	sinusugune
ehkki	miski	teiesugune
et	missugune	teine
justkui	mitmesugune	teineteise
kuigi	mitu	teistsugune
kuna	mõlema	tema
nagu	mõni	temasugune
olgugi	muu	too
otsekui	nemad	toosama
selmet	niipaljuke	üks
sest	niisugune	ükski
iga	oma	üksteise

### **III. Näide algtekstis, märgendatud tekstist, käsitsi koostatud kokkuvõttest ja EstSumi genereeritud kokkuvõttest**

#### **Algtekst**

##### **Brüssel tuleb täna välja uue eelarvekavaga**

Täna Euroopa Komisjoni avaldatav järgmine pikaajaline eelarvekava on seotud vähemalt kahe suure probleemiga: Suurbritannia ehk ühe suurema netomaksja lahkumine Euroopa Liidust ning Brüsseli soov kasutada suuremat tükki eelarvest tsentraalselt.

Brüsseli soov peegeldab tahet järgida Prantsusmaa presidendi Emmanuel Macroni üleskutset muuta Euroopa Liit tulevaste kriiside puhuks rahanduslikult võimekamaks, kirjutas Financial Times, kellel on õnnestunud eelarveperspektiivi mustand kätte saada.

Tegemist on 2021.–2027. aastani hõlmava eelarvekava esimese ettepaneku ehk perspektiiviga, kus arvatavasti konkreetseid arve veel ei ole. Kindlad summad võivad tulla juuni keskpaigas, aga riikidevahelised läbirääkimised nii omavahel kui ka Euroopa Komisjoniga kestavad ilmselt kaks aastat.

Vahepeal toimuvad aga Euroopa Parlamendi valimised, mille tulemusena vahetub Euroopa Komisjoni president ja ilmselt ka paljud volinikud. Vahepeal toimuvad mitmes riigis, näiteks Eestis, ka parlamendivalimised, mis samuti võib jõuvahekordi muuta. Eelarvekava vastuvõtmiseks peavad sellega nõustuma kõik Euroopa Liidu liikmesriigid. Riikidevahelised läbirääkimised eelarvekava üle on proovikivi riigipeadele. Valitsusjuhid teavad väga hästi, et nende läbirääkimisvõimekust mõõdetakse eurodes, õigemini miljonites ja miljardites eurodes. Seitsmeaastase eelarve kogumaht on triljon eurot.

«See on alati komplitseeritud. See võtab alati eeldatust rohkem aega. Ja seal on alati suur hulk dramaatilisust. Alati,» ütles eelarveläbirääkimistega tihedalt seotud Euroopa Komisjoni ametnik. «On hämmastav, et me üldse sellega toime tuleme. Aga me tuleme,» lisas ta.

Juba praegu väljendub kahe suurriigi Prantsusmaa ja Saksamaa suurim lahkeli suhtumises vajadusse Euroopa kriisifondide puhvreid suurendada.

«Vajame suuremat manööverdamisruumi,» ütles Euroopa Komisjoni ametnik Financial Timesile, nimetades põhjusena ohtu, et Euroopa fondid, nagu Euroopa Finantsstabiilsusmehhanism (EFSM), võivad kaotada kõrgeima (AAA) krediidireitingu.

Brüsseli ametnikud püüavad Euroopa laenusüsteemi muuta nii, et see ei põrkuks mitme riigi, eelkõige Saksamaa ja Hollandi vastuseisule. Saksamaa kantsleri Angela Merkeli sõnum teistele valitsusjuhtidele märtsikuisel tippkohtumisel oli, et fondide fiskaalvõimekus tuleb hoida nii madalal kui võimalik ning seda tuleks kasutada pigem investeringuteks kui lihtsalt majandustoetusteks.

Nagu öeldud, lööb Suurbritannia lahkumine Euroopa Liidust eelarvesse märgatava eelarveaugu – 10–15 miljardit eurot aastas. Selle katmiseks on plaanis suurendada sissemakset praeguselt ühelt protsendilt 1,2–1,3 protsendile sisemajanduse kogutoodangust. Kaalu-

takse maksubaasi laiendamist, näiteks ettevõtte tulumaksu ühtlustamist.

Samas on esimest korda oodata ühtekuuluvusfondide kahanemist. Seni on eurotoetus- te jagamise aluseks peamiselt üks kriteerium. Kuuldavasti kaalutakse uue perspektiivi väljatöötamisel veel lisakriteeriume: näiteks pagulaste vastuvõtmisest tingitud koor- mus ja noorte tööpuudus. Peale selle plaanitakse euroraha jagamisel arvesse võtta ka majandusreformide ja seaduste vastamist õigusriigi kriteeriumitele.

Sellega saavad ilmselt suurima löögi Poola ja Ungari, kes on eelmistel aastatel olnud suu- rimad Euroopa Liidu toetuste saajad. Bloombergi andmetel sai Poola aastatel 2014–2016 Euroopa Liidult netotoetusi keskmiselt 10,1 miljardit eurot aastas.

Samas on diplomaadid hoiatanud, et nende karistamine võib avada Pandora laeka ja pädida veel mõne riigi lahkumisega Euroopa Liidust. Igatahes Poola on väljendanud sellele kavatsusele juba vastuseisu.

«Näeme selles soovi avaldada mõnele riikidele enne läbirääkimisi poliitilist survet,» lausus uudisteagentuurile AP Poola välisminister Jacek Czaputowicz. «Seepärast suhtume taolistesse ideedesse väga negatiivselt.»

Bloombergi andmetel plaanitakse vähendada põllumajandustoetusi, seni nullilähedasi kaitsekulutusi aga märkimisväärselt suurendada. Märkimisväärselt on kavas tugevdada ka piirivalvet, suurendades töötajate arvu praegusega võrreldes enam kui viis korda.

«See on hädavajalik Euroopa Liidu piiri mitme lõigu, näiteks Kreeka saarte tõttu,» rääkis Euroopa Komisjoni allikas Reutersile. Uudisteagentuuri teatel kulutatakse käimasoleval seitseaastakul Euroopa Liidu piiride kaitseks neli miljardit eurot, piirivalvurite arvu suurendamiseks vähemalt 3000 võrra tuleb eelarvet kasvatada 25 miljardi euroni.

## Märgendatud tekst

<div0><head><hi rend="bold>Brüssel tuleb täna välja uue eelarvekavaga</hi></head>

<p>

<s><hi rend="bold>Täna Euroopa Komisjoni avaldatav järgmine pikaajaline eelarvekava on seotud vähemalt kahe suure probleemiga: Suurbritannia ehk ühe suurema netomaks- ja lahkumine Euroopa Liidust ning Brüsseli soov kasutada suuremat tükki eelarvest tsentraalselt.</hi></s>

</p>

<p>

<s>Brüsseli soov peegeldab tahet järgida Prantsusmaa presidendi Emmanuel Macroni üleskutset muuta Euroopa Liit tulevaste kriiside puhuks rahanduslikult võimekamaks, kirjutas Financial Times, kellel on õnnestunud eelarveperspektiivi mustand kätte saa- da.</s>

</p>

<p>

<s>Tegemist on 2021.–2027. aastani hõlmava eelarvekava esimese ettepaneku ehk pers- pektiiviga, kus arvatavasti konkreetseid arve veel ei ole.</s>

<s>Kindlad summad võivad tulla juuni keskpaigas, aga riikidevahelised läbirääkimised nii omavahel kui ka Euroopa Komisjoniga kestavad ilmselt kaks aastat.</s>

</p>

<p>

<s>Vahepeal toimuvad aga Euroopa Parlamendi valimised, mille tulemusena vahetub Euroopa Komisjoni president ja ilmselt ka paljud volinikud.</s>

<s>Vahepeal toimuvad mitmes riigis, näiteks Eestis, ka parlamendivalimised, mis samuti võib jõuvahekordi muuta.</s>

<s>Eelarvekava vastuvõtmiseks peavad sellega nõustuma kõik Euroopa Liidu liikmesriigid.</s>

</p>

<p>

<s>Riikidevahelised läbirääkimised eelarvekava üle on proovikivi riigipeadele.</s>

<s>Valitsusjuhid teavad väga hästi, et nende läbirääkimisvõimekust mõõdetakse eurodes, õigemini miljonites ja miljardites eurodes.</s>

<s>Seitsmeaastase eelarve kogumaht on triljon eurot.</s>

</p>

<p>

<s>«See on alati komplitseeritud. See võtab alati eeldatust rohkem aega. Ja seal on alati suur hulk dramaatilisust. Alati,» ütles eelarveläbirääkimistega tihedalt seotud Euroopa Komisjoni ametnik.</s>

<s>«On hämmastav, et me üldse sellega toime tuleme. Aga me tuleme,» lisis ta.</s>

</p>

<p>

<s>Juba praegu väljendub kahe suurriigi Prantsusmaa ja Saksamaa suurim lahkeli suhtumises vajadusse Euroopa kriisifondide puhvreid suurendada.</s>

</p>

<p>

<s>«Vajame suuremat manööverdamisruumi,» ütles Euroopa Komisjoni ametnik Financial Timesile, nimetades põhjusena ohtu, et Euroopa fondid, nagu Euroopa Finantsstabiilsusmehhanism (EFSM), võivad kaotada kõrgeima (AAA) krediitireitingu.</s>

</p>

<p>

<s>Brüsseli ametnikud püüavad Euroopa laenusüsteemi muuta nii, et see ei põrkuks mitme riigi, eelkõige Saksamaa ja Hollandi vastuseisule.</s>

<s>Saksamaa kantsleri Angela Merkeli sõnum teistele valitsusjuhtidele märtsikuisel tippkohtumisel oli, et fondide fiskaalvõimekus tuleb hoida nii madalal kui võimalik ning seda tuleks kasutada pigem investeringuteks kui lihtsalt majandustoetusteks.</s>

</p>

<p>



<s>Nagu öeldud, lööb Suurbritannia lahkumine Euroopa Liidust eelarvesse märgatava eelarveaugu – 10–15 miljardit eurot aastas.</s>

<s>Selle katmiseks on plaanis suurendada sissemakset praeguselt ühelt protsendilt 1,2–1,3 protsendile sisemajanduse kogutoodangust.</s>

<s>Kaalutakse maksubaasi laiendamist, näiteks ettevõtte tulumaksu ühtlustamist.</s>

</p>

<p>

<s>Samas on esimest korda oodata ühtekuuluvusfondide kahanemist.</s>

<s>Seni on eurotoetuste jagamise aluseks peamiselt üks kriteerium.</s>

<s>Kuuldavasti kaalutakse uue perspektiivi väljatöötamisel veel lisakriteeriume: näiteks pagulaste vastuvõtmisest tingitud koormus ja noorte tööpuudus.</s>

<s>Peale selle plaanitakse euroraha jagamisel arvesse võtta ka majandusreformide ja seaduste vastamist õigusriigi kriteeriumitele.</s>

</p>

<p>

<s>Sellega saavad ilmselt suurima löögi Poola ja Ungari, kes on eelmistel aastatel olnud suurimad Euroopa Liidu toetuste saajad.</s>

<s>Bloombergi andmetel sai Poola aastatel 2014–2016 Euroopa Liidult netotoetusi keskmiselt 10,1 miljardit eurot aastas.</s>

</p>

<p>

<s>Samas on diplomaadid hoiatanud, et nende karistamine võib avada Pandora laeka ja päädida veel mõne riigi lahkumisega Euroopa Liidust.</s>

<s>Igatahes Poola on väljendanud sellele kavatsusele juba vastuseisu.</s>

</p>

<p>

<s>«Näeme selles soovi avaldada mõnedele riikidele enne läbirääkimisi poliitilist survet,» lausus uudisteagentuurile AP Poola välisminister Jacek Czaputowicz.</s>

<s>«Seepärast suhtume taolistesse ideedesse väga negatiivselt.»</s>

</p>

<p>

<s>Bloombergi andmetel plaanitakse vähendada põllumajandustoetusi, seni nullilähedasi kaitsekulutusi aga märkimisväärselt suurendada.</s>

<s>Märkimisväärselt on kavas tugevdada ka piirivalvet, suurendades töötajate arvu praegusega võrreldes enam kui viis korda.</s>

</p>

<p>

<s>«See on hädavajalik Euroopa Liidu piiri mitme löögu, näiteks Kreeka saarte tõttu,» rääkis Euroopa Komisjoni allikas Reutersile.</s>

<s>Uudisteagentuuri teatel kulutatakse käimasoleval seitseaastakul Euroopa Liidu piiride kaitseks neli miljardit eurot, piirivalvurite arvu suurendamiseks vähemalt 3000 võrra tuleb eelarvet kasvatada 25 miljardi euronni.</s>

</p>

</div0>

## **Käsitsi koostatud kokkuvõtete**

### **Brüssel tuleb täna välja uue eelarvekavaga**

Täna Euroopa Komisjoni avaldatav järgmine pikaajaline eelarvekava on seotud vähemalt kahe suure probleemiga: Suurbritannia ehk ühe suurema netomaksja lahkumine Euroopa Liidust ning Brüsseli soov kasutada suuremat tükki eelarvest tsentraalselt.

Tegemist on 2021.–2027. aastani hõlmava eelarvekava esimese ettepaneku ehk perspektiiviga, kus arvatavasti konkreetseid arve veel ei ole.

Vahepeal toimuvad aga Euroopa Parlamendi valimised, mille tulemusena vahetub Euroopa Komisjoni president ja ilmselt ka paljud volinikud.

Riikidevahelised läbirääkimised eelarvekava üle on proovikivi riigipeadele.

Juba praegu väljendub kahe suurriigi Prantsusmaa ja Saksamaa suurim lahkeli suhtumises vajadusse Euroopa kriisifondide puhvreid suurendada.

Brüsseli ametnikud püüavad Euroopa laenusüsteemi muuta nii, et see ei põrkuks mitme riigi, eelkõige Saksamaa ja Hollandi vastuseisule.

Samas on esimest korda oodata ühtekuuluvusfondide kahanemist.

Sellega saavad ilmselt suurima löögi Poola ja Ungari, kes on eelmistel aastatel olnud suurimad Euroopa Liidu toetuste saajad.

Poola on väljendanud sellele kavatsusele juba vastuseisu.

Bloombergi andmetel plaanitakse vähendada põllumajandustoetusi, seni nullilähedasi kaitsekulutusi aga märkimisväärselt suurendada.

## **EstSumi genereeritud kokkuvõte**

<s><hi rend="bold">Täna Euroopa Komisjoni avaldatav järgmine pikaajaline eelarvekava on seotud vähemalt kahe suure probleemiga: Suurbritannia ehk ühe suurema netomaksja ja lahkumine Euroopa Liidust ning Brüsseli soov kasutada suuremat tükki eelarvest tsentraalselt.</hi></s>

<s>Brüsseli soov peegeldab tahet järgida Prantsusmaa presidendi Emmanuel Macroni üleskutset muuta Euroopa Liit tulevaste kriiside puhuks rahanduslikult võimekamaks, kirjutas Financial Times, kellel on õnnestunud eelarveperspektiivi mustand kätte saada.</s>

<s>Vahepeal toimuvad aga Euroopa Parlamendi valimised, mille tulemusena vahetub Euroopa Komisjoni president ja ilmselt ka paljud volinikud.</s>

<s>Riikidevahelised läbirääkimised eelarvekava üle on proovikivi riigipeadele.</s>

<s>Juba praegu väljendub kahe suurriigi Prantsusmaa ja Saksamaa suurim lahkeli suhtumises vajadusse Euroopa kriisifondide puhvreid suurendada.</s>

<s>Brüsseli ametnikud püüavad Euroopa laenusüsteemi muuta nii, et see ei põrkuks mitme riigi, eelkõige Saksamaa ja Hollandi vastuseisule.</s>

<s>Nagu öeldud, lööb Suurbritannia lahkumine Euroopa Liidust eelarvesse märgatava eelarveaugu – 10–15 miljardit eurot aastas.</s>

<s>Sellega saavad ilmselt suurima löögi Poola ja Ungari, kes on eelmistel aastatel olnud suurimad Euroopa Liidu toetuste saajad.</s>

#### IV. Parameetrite kolmikute keskmine kattuvus arnedus-korpusel

$\alpha$	$\beta$	$\gamma$	Keskmine kattuvus
0,4	0,4	0,2	62,24%
0,5	0,4	0,1	61,68%
0,4	0,3	0,3	61,65%
0,4	0,5	0,1	61,46%
0,7	0,2	0,1	61,39%
0,8	0,1	0,1	61,34%
0,6	0,3	0,1	61,32%
0,7	0,1	0,2	61,28%
0,6	0,2	0,2	61,27%
0,3	0,6	0,1	61,25%
0,5	0,2	0,3	61,08%
0,5	0,3	0,2	60,85%
0,2	0,7	0,1	60,74%
0,3	0,5	0,2	60,7%
0,3	0,4	0,3	60,27%
0,6	0,1	0,3	60,26%
0,2	0,6	0,2	59,94%
0,4	0,2	0,4	59,41%
0,1	0,8	0,1	58,55%
0,3	0,3	0,4	58,2%
0,2	0,4	0,4	57,94%
0,2	0,5	0,3	57,47%
0,5	0,1	0,4	57,21%
0,3	0,2	0,5	56,49%
0,1	0,7	0,2	55,83%
0,4	0,1	0,5	55,65%
0,2	0,3	0,5	55,3%
0,3	0,1	0,6	55,02%
0,1	0,6	0,3	51,63%
0,2	0,2	0,6	51,05%
0,1	0,5	0,4	50,44%
0,2	0,1	0,7	50,4%
0,1	0,2	0,7	47,64%
0,1	0,4	0,5	47,56%

0,1	0,3	0,6	47,46%
0,1	0,1	0,8	45,18%

## **V. Litsents**

### **Lihtlitsents lõputöö reprodutseerimiseks ja lõputöö üldsusele kättesaadavaks tegemiseks**

Mina, **Janar Saks**,

1. annan Tartu Ülikoolile tasuta loa (lihtlitsentsi) enda loodud teose

#### **PEALKIRI**

mille juhendaja on Kaili Müürisep

- 1.1 reprodutseerimiseks säilitamise ja üldsusele kättesaadavaks tegemise eesmärgil, sealhulgas digitaalarhiivi DSpace-is lisamise eesmärgil kuni autoriõiguse kehtivuse tähtaja lõppemiseni;
  - 1.2 üldsusele kättesaadavaks tegemiseks Tartu Ülikooli veebikeskkonna kaudu, sealhulgas digitaalarhiivi DSpace'i kaudu kuni autoriõiguse kehtivuse tähtaja lõppemiseni.
2. olen teadlik, et punktis 1 nimetatud õigused jäävad alles ka autorile.
3. kinnitan, et lihtlitsentsi andmisega ei rikuta teiste isikute intellektuaalomandi ega isikuandmete kaitse seadusest tulenevaid õigusi.

Tartus, 14.05.2018