

LEKSIKOGRAAFILISE TARKVARA SKETCH ENGINE EESTI KEELE MOODUL

Jelena Kallas, Maria Tuulik ja Madis Jürviste

Eesti Keele Instituut

Kokkuvõte. 2010. aasta sügisel alustas Eesti Keele Instituut koos ettevõttega Lexical Computing Ltd. leksikograafilise tarkvara Sketch Engine (Kilgarriff jt 2004) eestikeelse mooduli väljatöötamist. Artiklis kirjeldatakse programmi põhifunktsioone. Põhjalikumalt käsitletakse funktsiooni Word Sketch (ee *sõnavisand*) võimalusi. Tutvustatakse sõnavisandite grammatika koostamise põhimõtteid, vaadeldakse eraldi substantiivide, adjektiivide ja verbide sõnavisandites esitatud süntagmaatilisi seoseid (st grammatilisi ja leksikaalseid kollokatsioone) ning arutletakse mooduli edasiarendusvõimaluste üle. Lisaks analüüsitakse, mil määral saab sõnavisandeid kasutada verbide lausemallide tuvastamisel.

Märksõnad: korpuseleksikograafia, konkordantsid, sõnavisand, tesaurus, eesti keel

1. Sketch Engine'i põhifunktsioonid

2002. aastal loodud tarkvara Sketch Engine (SkE)¹ on tekstikorpuste töötlemiseks mõeldud programm, mille põhifunktsioonid on konkordantside koostamine ja nende mitmekülgne töötlemine, korpusest sagedusloendite koostamine, korpusest automaatne sõnavisandite genereerimine ja distributiivne tesaurus (Kilgarriff jt 2004). Kokku on programmis esindatud 42 keele korpused. Tänapäeval kasutavad programmi sõnaraamatute loomisel sellised tuntud kirjastused nagu Oxford University Press, Cambridge University Press, Collins, Le Robert ja Cornelsen Verlag ja seda rakendati ka Dante leksikaalse andmebaasi (Atkins jt 2010) väljatöötamisel. Eestis kasutatakse prog-

¹ Sketch Engine vt <http://the.sketchengine.co.uk/auth/corpora/>. Eesti Keele Instituut sõlmis lepingu 2010. a. sügisel.

rammi praegu ühekõitelise seletava sõnaraamatu (Langemets jt 2010) ja eesti keele põhisõnavara sõnastiku (Kallas ja Tuulik 2011) koostamisel.

Programmi eestikeelse mooduli sisendiks on eesti keele koondkorpus² (u 250 mln sõnet), mille OÜ Filosoft on morfoloogiliselt märgendanud, osaliselt ühestanud ja osalausestanud. Peale selle on kättesaadav ka eesti lastekeele korpus (u 400 000 sõnet).


1.1. Konkordantsid ja nende töötlus

Konkordants on sõnavormide loend koos nende tekstilise ümbrusega. Konkordantse kasutatakse mitmes rakenduslingvistika valdkonnas, eelkõige leksikograafias (Atkins jt 2003: 252–253) ja keeleõppes (Kitsnik 2006: 97–98). Adam Kilgarriff (2009) nimetab korpuste tugevuseks just autentset keelt ja ülisuurt näidete arvu. Konkordantsitööriist pakub õppijale suuremat hulka valikuid kui sõnaraamatud, näiteks saab otsida eri parameetrite põhjal ja valida ka viisi, kuidas tulemusi kuvatakse. Sõnaraamatu puhul on õppijal endal valikuvõimalusi paratamatult vähem, kuna leksikograaf on juba ära otsustanud, mida iga sõna puhul teada on vaja.

Sketch Engine on oma olemuselt midagi korpuse ja sõnaraamatu vahepealset. Programmis saab konkordantse koostada sõnavormist, lemmast või fraasist lähtudes või programmi pärin-gukeelt³ kasutades. Konkordantsi põhisõnale (ingl *keyword*) klõpsates näeb seda sõna lauses koos laiema kontekstiga, mis kuvatakse sõnast vasakul ja paremal. Joonisel 1 on esitatud sõna *armastus* konkordants koos morfoloogiliste märgenditega (sõnaliik, lemma, muutevormi tunnus).

2 Eesti keele koondkorpus vt <http://test.cl.ut.ee/korpused/segakorpus/>. Korpuse märgendamise kirjeldust vt http://www.filosoft.ee/html_morf_et/morfoutinfo.html

3 vt <http://trac.sketchengine.co.uk/wiki/SkE/CorpusQuerying>



[About](#) [Home](#) [Settings](#) [Change password](#) [Log out](#)

user: Jelena Kallias corpus: EstonianRC

Search in EstonianRC

[? Help on main menu](#)

[? Help on Conc. menu](#) [Switch menu position](#)

[Concordance](#) [Word List](#) [Word Sketch](#) [Thesaurus](#) [Find X](#) [Sketch-Diff](#) [Sketch-Eval](#)

[Save](#) [View options](#) [Sort](#) [Sample](#) [Filter](#) [Frequency](#) [Collocations](#) [ConcDesc](#)

Corpus: **EstonianRC**

Hits: **18576** (74.4 per million)

Page 1 of 929 [Next](#) | [Last](#)

3	lähedale jäävad . Ja eraldi muidugi on mul armastus /S/armastus	3	Riigikaitse Akadeemia kuldude vastu . Ma ning sõidavad poliitilise valge hobuse
3	puhtalt poliitiline . Ühed lubavad õnne ja armastust /S/armastus	4	. Kui me tahame tema kõiki külgi kirjeldada
4	, et see on umbes samasugune mõiste kui armastus /S/armastus	7	eesmärgil . See on tagumise poole vastus
7	korrastada ja teeb seda n-õ vastastikuse armastuse /S/armastus	14	on meil üks . " Usun , et see kehtib ka
14	Teed võivad meil olla küll erinevad , aga armastuse /S/armastus	15	kasvatuse tulemust saab ilmselt mõõta pigem armastusega /S/armastus
27	kohta väidetakse , et nad tunnevad sügavat armastust /S/armastus	32	kohaliku etu vastu , mis ei lase neil mitte
32	on , et ettepaneku on põhjustanud rohkem armastus /S/armastus	32	valija vastu kui armastus kalapüügi vastu
32	põhjustanud rohkem armastus valija vastu kui armastus /S/armastus	32	kalapüügi vastu . See on kõik . Tänan !
33	See eeldab täielikku pühendumist , suurt armastust /S/armastus	33	Ja kõiki muid häid asju , millest meil
33	toetab teda lugupidamisega , truudusega ja armastusega /S/armastus	33	. Palju õnne , härra Vabariigi President
33	Sa oled väärt , et sind armastataks . " Armastus /S/armastus	33	tähendab ennekõike pühendumist . Selleks
33	nagu ju kutsus Aino Jarvesoo üles - see on armastuse /S/armastus	33	defitsiidi probleemi üks seige väljendus

Page 1 of 929 [Next](#) | [Last](#)

Joonis 1. Sõna *armastus* konkordants.

Menüüst näeb konkordantside töötlemise võimalusi. Konkordantse saab salvestada (ingl *Save*) eri vormingutes, nt *.txt*- ja *.csv*-vormingusse. Vaate profiili (ingl *View options*) sätete muutmiseks saab lisaks põhisonale vaadata ka kõigi konkordantsides olevate sõnade morfoloogilisi märgendeid, saab määrata kuvatava konteksti pikkust ning konkordantsi ridade arvu. Kuna esialgu oli SkE loodud leksikograafide töövahendiks, siis prooviti parimate näitelauseste valimist hõlbustada. Sel eesmärgil lisati Sketch Engine'isse aastal 2008 uus funktsioon nimega GDEX (ingl *Good Dictionary Example*) (Kilgarriff jt 2008). Eestikeelses moodulis on praegu GDEX-i nn *vanilla*-versioon, mis on keelest sõltumatu (olulisteks kriteeriumiteks on lause pikkus ning see, et lause alguses oleks suurtäht ja lõpus punkt, hüüumärk või küsimärk). Mitme keele jaoks on loodud ka keelespetsiifilised GDEX-i rakendused, mis arvestavad ka seda, mil määral sisaldub lauses konkreetse keele põhisonavara (ingl *core vocabulary*). Vaate profiilis saab määrata, mitut parimat näitelausest soovetakse näha.

Sorteerimise abil (ingl *Sort*) saab konkordantse lahterdada põhisona või parema ja vasaku naabri järgi. Sagedus (ingl *Frequency*) võimaldab koostada sagedusloendeid lähtuvalt põhisona lemmast ja sõnavormist, nt saab teada, mis käändes esineb mingi substantiiv kõige sagedamini.

Erinevate statistikute abil (ingl *T-score*, *MI*, *MI3*, *log likelihood*, *min. sensitivity*, *logDice*) saab otsida põhisona kollokaate (ingl *Collocates*).

1.2. Tesaaurus

Distributiivne tesaaurus on funktsioon, mille abil saab automaatselt genereerida konkreetse lemma tesaauruse. Tegemist ei ole semantilise andmebaasiga, mis keskenduks mõistele ja semantiliste suhete kaudu tema semantilisele väljale (vrd nt Vider 1999), vaid süsteemiga, mis toob statistika põhjal esile sõnu, millel on sarnane grammatiline ja kollokatiivne käitumine. Kuid siiski võib sarnast käitumist põhjustada ka samasse tähendusvälja kuulumine (nt *armastus*, *kirg*, *sõprus* joonisel 2).

armastus		
EstonianRC freq = 18576		
Lemma	Score	Freq
<u>ilu</u>	0.208	7174
<u>usk</u>	0.201	11245
<u>vabadus</u>	0.196	19110
<u>sõprus</u>	0.194	5372
<u>õnn</u>	0.192	29198
<u>king</u>	0.189	5126
<u>rõõm</u>	0.186	17878
<u>tunne</u>	0.184	29587
<u>lootus</u>	0.183	31330

Joonis 2. Lemma *armastus* tesaurus.

1.3. Sõnavisand

Sõnavisand on automaatne üheleheline kokkuvõte sõna süntaktilisest ja kollokatiiivsest käitumisest, mis tugineb korpusele (Kilgarriff jt 2004). SkE-s saab näha ühendit või konstruktsiooni moodustavate sõnade koosinemiste arvu (ingl *raw frequency*). Kuid sõnadevahelise seose tugevuse mõõtmiseks on rakendatud ka esilduvust (ingl *salience*), mille annab *logDice*-statistik. Kristel Uiboaed (2010: 307) peab statistikute kasutamise eeliseks seda, et peale sõnade koosinemise võetakse arvesse ka ühendit moodustavate sõnade eraldi esinemise sagedusi.

Sõnavisandi funktsiooni rakendamine eeldab seda, et oleks kirjutatud sõnavisandite grammatika (ingl *Sketch Grammar*), mis määrab, milliseid grammatilisi suhteid (ingl *grammatical relations*) programm otsima hakkab. Süsteem eristab grammatilisi suhteid, mille määrab sõnavisandite grammatika, ja alles seejärel mõõdetakse nende esinemissagedust ja esilduvust.

Sõnavisandite päringu kasutajaliideses toimub otsing lemma kaudu. Laiendatud päringuga (ingl *Advanced options*) on võimalik (vt joonist 3):

- luua isiklik allkorpus (ingl *Subcorpus*), st valida olemasoleva korpuse failidest need, mis hakkavad kuuluma allkorpusesse;
- määrata suhte minimaalne esinemissagedus (ingl *Minimum frequency*);
- määrata grammatilise suhte minimaalne esilduvus (ingl *Minimum salience*);
- määrata ühe kategooria kuvatavate üksuste arv (ingl *Maximum number of items in grammatical relations*);
- sortida kollokaate esilduvuse määra või koosinemiste arvu järgi (ingl *Sort collocations according to salience/ raw frequency*);
- kasutada funktsiooni “Tickbox Lexicography template”, mis pakub eri kollokatsioonide kohta (GDEX-i sätete kohaselt) välja teatud arvu nn parimaid näitelauseid;
- klasterdada kollokatsioone (ingl *Cluster collocations*);
- määrata minimaalne sarnasus klasterdavate üksuste vahel (ingl *Minimal similarity between cluster collocations*).

<p>Concordance</p> <p>Word List</p> <p>Word Sketch</p> <p>Thesaurus</p> <p>Find X</p> <p>Sketch-Diff</p> <p>Sketch-Eval</p> <p>? Help on main menu</p> <hr/> <p>? Help on Word Sketches</p> <p>Advanced options</p> <p>Switch menu position</p>	<h3>Word Sketch Entry Form</h3>
<p>Lemma: <input type="text" value="armastus"/></p>	
<h4>Advanced options</h4>	
Subcorpus:	<input type="text" value="None (whole corpus)"/> · info create new
Minimum frequency:	<input type="text" value="1"/>
Minimum salience:	<input type="text" value="0.0"/>
Maximum number of items in a grammatical relation:	<input type="text" value="10"/>
Sort collocations according to:	<input type="radio"/> Salience <input checked="" type="radio"/> Raw frequency
Tickbox Lexicography template:	<input type="text" value="None"/> · Examples per collocate: <input type="text" value="10"/>
Cluster collocations	<input type="checkbox"/>
Structure word sketch by gramrels	<input checked="" type="checkbox"/>
Minimum similarity between cluster items:	<input type="text" value="0.15"/>
<input type="button" value="Show Word Sketch"/> <input type="button" value="Save Options"/>	

Joonis 3. Sõnavisandi kasutajaliides.

2. Eesti keele sõnavisandid

Artikli teises osas tutvustame eesti keele sõnavisandite grammatika koostamise põhimõtteid ja näitame sõnaliikide kaupa süntagmaatiliste seoste esitamise printsiipe. Illustreerime ka sõnavisandite võrdlemise funktsiooni kasutust.

2.1. Sõnavisandite grammatika koostamise põhimõtted

Grammatika kirjutamise formalismiks on regulaaravaldised, mis põhinevad sõnaliigil ja muutetunnustel⁴. Lähtuvalt sõnaliigist ja muutetunnustest esitab süsteem sellised grammatilised klassid nagu grammatiline predikaat (öeldis), grammatiline subjekt (alus), grammatiline objekt (sihitis), predikatiiv (öeldistäide), adverbiaal (määrus). Grammatilistest seostest on keskendunud rinnastusseoste ja alistusseoste esiletoomisele. Rinnastusseostest toob süsteem esile konstruktsioonid sidesõnadega *ja/või*, *kui/nagu*. Alistusseostest esitatakse võimalike substantiivi, adjektiivi, verbi ja adverbi reksioonistruktuuride enamik (vt lähemalt Kallas ja Tuulik 2011: 70).

Grammatilisi suhteid on sõnavisandite grammatikas nelja liiki (vt ka Svensén 2009: 416–418):

- a) üheliikmeline suhe (UNARY) on defineeritud ühe konkreetse morfoloogilise kategooria kaudu ning see suhe annab infot ühe kindla grammatilise vormi kohta, nt mis käändes esineb konkreetse noomeni lemma kõige sagedamini;
- b) sümmeetriline suhe (SYMMETRIC) puudutab eelkõige rinnastusseoses olevaid ühendeid (eesti moodulis on nendeks *ja/või* ja *kui/nagu* suhted);
- c) kaheliikmelised suhted (DUAL) võimaldavad otsida kahe sõna ühendeid ja on määratletud kahe morfoloogilise kategooria kaudu (nt *subjekt*, *objekt*);

4 St noomeni arvu- ja käändetunnustel, võrdlusastmete tunnustel, verbi arvu-, isiku-, aja-, kõneviisi-, tegumoe- ja kõneliigitunnustel, samuti infiniitsete verbivormide tunnustel.

- d) kolmeliikmelised suhted (TRINARY) võimaldavad leida, millised nimisõnad esinevad substantiivi, adjektiiviga ja verbi kaassõnafrasides, ja ühendverbide puhul nende laiendeid.

Sõnavisandite grammatika koostamisel on kasutatud ka „Eesti keele formaalses grammatikas” (Roosmaa jt 2001) kirjeldatud kitsendusi. Kokku on programmiga võimalik näidata 38 tüüpi grammatilisi seoseid. Peale selle võimaldab süsteem koostada ühendverbe moodustavate afiksaaladverbide ja verbi laienditega koos esinevate kaassõnade sagedusjärjestusi ning väljendverbide nominaalsete komponentide sagedusjärjestusi tingimusel, et nominaalne komponent on märgendatud kui „X”. (X-iga on märgendatud verbi juurde kuuluv sõna, millel eraldi sõnaliigi tähistus puudub, nt *plehku*.)

2.2. Substantiivide sõnavisandid

Substantiivide visandites tulevad esile sellised kategooriad nagu subjekt (mida X teeb?), objekt (mida X-iga tehakse?), adjektiivatribuut (milline X on?) ja eri tüüpi aktantsed atribuudid, mille moodustavad ühelt poolt semantilised käanded ja kaassõnad, teiselt poolt genitiiv. Rektsioonistruktuuridest esinevad visandites käände- (usk *kellesse-millesse*), kaassõna- (viha *kelle-mille vastu*) ja tegevusnimerektsioon (tahe *mida teha*). SkE abil saab mingil määral vaadelda ka sõna paradigmaatilisi tähendussuhteid, näiteks atributsiooni (*hall, valge hiir*), mero-nüümiat (*auto rool, mootor, ratas*) ja komplementaarseid vastandusi (*poiss ja tüdruk*), kuna seotud lemmad on tihtipeale kollokatsioonidest väljanõpitud. Allpool on esitatud substantiivide *päike* ja *usk* sõnavisandid.

Päikese sõnavisandis esinevad sellised kategooriad nagu subjekt (*päike paistab, loojub, särab...*), objekt (*päikest nau-tima, võtma*), adjektiivatribuut (*loojuv, tõusev, lõõskav päike*), genitiivatribuut (*päikese loojumine, kiir, aktiivsus*), aktantne atribuut (*päikese käes*) ja otsitava sõnaga rinnastusseoses olevad substantiivid (*päike ja vihm, päike ja kuu...*). Programm

päike ⁽¹⁾ EstonianRC freq = 16463 (65.9 per million)								
object of	338	3.1	subject of	3106	7.2	a modifier	2234	2.3
nautima	<u>29</u>	7.19	paistma	<u>948</u>	10.52	loojuv	<u>143</u>	10.92
võtma	<u>76</u>	3.42	loojuma	<u>147</u>	10.48	tõusev	<u>271</u>	10.6
			särama	<u>116</u>	9.39	lõõskav	<u>110</u>	10.56
			kõrvetama	<u>66</u>	9.17	ere	<u>164</u>	10.39
			lõõskama	<u>49</u>	8.97	kõrvetav	<u>102</u>	10.38
gen modifies	1615	1.3	gen modifier	498	0.4			
loojumine	<u>41</u>	9.64	lõunamaa	<u>25</u>	9.07			
ultraviolettkiirgus	<u>39</u>	9.27						
aktiivsus	<u>91</u>	8.52						
kiir	<u>42</u>	8.5	ja/või	792	3.2			
kiirgus	<u>34</u>	8.19	vihm	<u>52</u>	7.83			
			tuul	<u>79</u>	7.02			
			vesi	<u>29</u>	4.02			
noun seesütlev	230	1.8	kuu	<u>40</u>	3.66			
käsi	<u>82</u>	4.83						
noun alalütlev	175	1.5						
maa	<u>56</u>	3.46						

Joonis 4. Substantiivi *päike* sõnavisand.

näitab iga üksuse juures ka koosinemiste arvu ja esilduvuse määra.

Usu sõnavisandis tulevad esile subjekt (*usk kaob, puudub, aitab...*), objekt (*usku kaotama, sisendama, andma...*), adjektiiv-atribuut (*hea, kindel, uus usk ...; katoliku, luteri, vene usk...*), genitiivatribuut (*usu põhimõte, esindaja, puudumine...*), öeldistide (*usk on oopium...*) ja otsitava sõnaga rinnastusseoses olevad substantiivid (*usk ja lootus, usk ja rahvus...*).

Süsteem käsitleb atributiivset suhet kaheliikmelise paarina, seega esitatakse mõned ühendid puudulikult (nt *usu kirik*), kuna neil on veel omakorda laiendid, mis jäävad visandist kõrvale (*luteri usu kirik*). Kollokatsiooni kasutusjuhud koos kontekstiga avanevad koosinemisarvu alt.

usk ()

EstonianRC freq = 11245 (45.0 per million)

a_modifier	1623	2.2	subject_of	582	1.8	object_of	473	5.8
hea	<u>319</u>	5.59	puuduma	<u>42</u>	5.57	andma	<u>47</u>	2.76
kindel	<u>124</u>	6.38	kaduma	<u>39</u>	5.52	sisendama	<u>46</u>	9.69
uus	<u>113</u>	3.28	aitama	<u>34</u>	4.87	lisama	<u>33</u>	3.99
suur	<u>74</u>	2.71	lubama	<u>23</u>	3.72	kaotama	<u>28</u>	5.04
pime	<u>45</u>	7.01	tulema	<u>20</u>	1.08	väljendama	<u>21</u>	6.23
eriline	<u>45</u>	5.6	jääma	<u>19</u>	1.62	kinnitama	<u>19</u>	2.66
kristlik	<u>39</u>	7.33	andma	<u>19</u>	1.45	avaldama	<u>16</u>	3.64

a_modifier_comp	49	1.1
suurem	<u>16</u>	2.42

a_modifier_ordinal	72	1.1
teine	<u>71</u>	2.83

adj_modifier_käandumatu	237	8.9
katoliku	<u>81</u>	10.33
luteri	<u>74</u>	11.33
vene	<u>18</u>	3.36
eesti	<u>18</u>	2.93
muhamedi	<u>17</u>	10.62
buda	<u>10</u>	9.71

gen_modifier	1050	1.2
islam	<u>153</u>	10.14
inimene	<u>135</u>	3.55
rahvas	<u>50</u>	4.39
eestlane	<u>24</u>	3.41
juut	<u>19</u>	5.7
moslem	<u>16</u>	7.01
kodanik	<u>15</u>	3.7

gen_modifies	762	0.8
põhimõte	<u>35</u>	4.9
jumal	<u>17</u>	5.05
esindaja	<u>16</u>	2.59
kirik	<u>12</u>	3.2
puudumine	<u>11</u>	3.55
küsimus	<u>11</u>	0.98

olema_noun	80	8.6
uskmatust	<u>11</u>	11.14
oopium	<u>10</u>	9.72

noun_sisseütlev	417	32.6
jumal	<u>66</u>	7.04
õiglus	<u>15</u>	7.14
ime	<u>12</u>	6.38

ja/või	953	5.1
lootus	<u>34</u>	5.82
keel	<u>33</u>	3.71
rahvus	<u>31</u>	6.35
rass	<u>30</u>	8.48
enesekindlus	<u>24</u>	7.79
armastus	<u>23</u>	5.51
teadus	<u>18</u>	5.57

Joonis 5. Substantiivi *usk* sõnavisand.

2.3. Adjektiivide sõnavisandid

Adjektiivifraasi laiendiliikmeks võib olla a) substantiiv(iframe); b) kaassõnafraas; c) infinitiiv(iframe); d) adjektiiv(iframe), e) kvantoriframe; f) adverb(iframe), g) kõrvallause (EKG 1993: 130). Sõnavisandid näitavad substantiivi, kaassõnafraasi, infinitiivifraasi ja adverbi. Rektsioonistruktuuridest on esitatud käände- (kindel *kelles-milles*), kaassõna- (kade *kelle-mille peale*) ja tegevusnimerektsioon (julge *mida tegema*).

Toome näiteks adjektiivide *hea* ja *rikas* sõnavisandid.

hea (EstonianRC freq = 224446 (898.7 per million))		
modifies	146325	7.9
meel	<u>10497</u>	10.85
tulemus	<u>3711</u>	8.87
kolleeg	<u>3244</u>	9.22
võimalus	<u>3094</u>	8.71
sõber	<u>2660</u>	8.88
subj_olema	5466	8.2
vorm	<u>147</u>	6.54
enesetunne	<u>146</u>	8.98
tulemus	<u>116</u>	5.0
da_infinitiiv	1853	0.7
olema	<u>249</u>	3.82
elama	<u>143</u>	6.17
teadma	<u>119</u>	4.53
adverb	62517	5.4
väga	<u>20295</u>	11.36
nii	<u>4993</u>	9.4
päris	<u>2460</u>	9.8
kui	<u>2203</u>	7.74
eriti	<u>1605</u>	8.92

Joonis 6. Adjektiivi *hea* sõnavisand.

Lemma *hea* visandis tulid esile järgmised kategooriad: sagedasemad substantiivid, mille fraasis esineb *hea* adjektiivatribuudina (*hea meel, tulemus, sõber*...). Sellest võib järeldada, et *hea meel* on lemma *hea* sagedasem *substantiiv+adjektiiv*-tüüpi kollokatsioon. Sageduselt teisel kohal on adverbid (*väga, päris, eriti hea*...). Peale selle tulid esile predikatiiv (*enesetunne on hea, vorm on hea, tulemus on hea*) ja tegevusnimerektsioon (*hea elada, hea teada, hea olla*).

rikas () EstonianRC freq = 12796 (51.2 per million)

<u>modifies</u>	<u>4546</u>	<u>6.1</u>	<u>subj_olema</u>	<u>242</u>	<u>9.0</u>	<u>adverb</u>	<u>3872</u>	<u>8.3</u>
riik	<u>665</u>	5.95	riik	<u>13</u>	0.29	nii	<u>684</u>	6.93
inimene	<u>469</u>	5.34	ajalugu	<u>7</u>	2.28	väga	<u>426</u>	6.17
mees	<u>379</u>	5.75	keel	<u>6</u>	1.26	kui	<u>262</u>	4.94
maa	<u>139</u>	4.74	linn	<u>6</u>	0.28	ka	<u>189</u>	4.77
ärimees	<u>72</u>	7.0	maa	<u>6</u>	0.24	ratsa	<u>152</u>	10.27
pere	<u>67</u>	5.21	vanem	<u>4</u>	1.11	ainult	<u>121</u>	6.81
naine	<u>65</u>	3.88	loomaaed	<u>3</u>	4.08	vaid	<u>107</u>	6.43

<u>adjective_pp_poolest</u>	<u>57</u>	<u>4778.0</u>
kaalium	<u>3</u>	7.45
antioksidant	<u>3</u>	7.22
maavara	<u>3</u>	6.93
loodusvara	<u>3</u>	6.76
flavonoidi	<u>2</u>	9.46
B-vitamiin	<u>2</u>	8.87
kiudaine	<u>2</u>	7.51

Joonis 7. Adjektiiv *rikas* sõnavisand.

Lemma *rikas* visandis esinesid järgmised kategooriad: sagedasemad substantiivid, mille fraasis esineb *rikas* adjektiivatribuudina (*rikas riik, inimene, mees...*). Sageduselt teisel kohal on adverbid (*nii, väga, piisavalt rikas*). Lisaks tulid esile predikatiiv (*riik on rikas, maa on rikas*) ja kaassõnareksioon (*rikas mille poolest*).

2.4. Verbide sõnavisandid

Verbide sõnavisandites toob süsteem peale subjekti ja objekti esile ka adverbiaale ehk määrusi. Rektsioonistruktuuridest on esindatud järgmised rektsioonid: objekti- (toetama *kedamida*), käände- (tutvuma *kellega-millega*), kaassõna- (võitlema

kelle-mille vastu) ja tegevusnimereksioon (tahtma *mida teha*, jätma *mida tegemata*).

Suur leksikograafiline probleem on aga ühendverbide esiletoomine ning nende süntagmaatiliste seoste selgitamine. Kuna afiksaaladverbid on märgendatud korpuses kui adverbid, siis koostasime EKSS-i ühendverbide põhjal sagedasemate afiksaaladverbide loendi ja lisasime selle sõnavisandite grammatikasse. Tulemuseks on see, et süsteem otsib korpusest ja esitab sagedusloendina ka iga verbi ühendverbe, samuti ühendverbide objekte ja adverbiaale. Teine eesmärk oli tõsta visandites esile, millised kaassõnafraasid laiendavad üht või teist verbi (nt hoolitsema *kelle-mille eest*). Neid suhteid oli võimalik määratleda kolmeliikmeliste kategooriate kaudu.

Toome näiteks verbide *hoolitsema* ja *tooma* sõnavisandid.

hoolitsema ()		EstonianRC freq = 9376 (37.5 per million)	
subject	1133 4.6	adverbial_kaasaütlev	55 4.0
riik	<u>78</u> 2.87	välimus	<u>26</u> 6.45
firma	<u>31</u> 3.24		
mees	<u>27</u> 1.96	verb_pp_eest	483 446.9
valitsus	<u>26</u> 2.12	laps	<u>41</u> 2.81
inimene	<u>24</u> 1.06	maja	<u>9</u> 1.3
ema	<u>22</u> 3.69	skoor	<u>8</u> 7.42
naine	<u>22</u> 2.36	loom	<u>8</u> 2.97
		adverb	969 2.8
		hästi	<u>45</u> 5.56
		rohkem	<u>39</u> 5.11
		ka	<u>37</u> 2.44
		paremini	<u>32</u> 6.16
		nüüd	<u>24</u> 3.51
		siis	<u>24</u> 2.87
pp_verb	2685 26.1	modifies_past_participle	269 4.1
eest	<u>2574</u> 10.04	välimus	<u>43</u> 7.15
kõrval	<u>8</u> 3.69	aed	<u>18</u> 5.54
		käsi	<u>16</u> 2.47
		muru	<u>10</u> 5.98

Joonis 8. Verbi *hoolitsema* sõnavisand.

Verbi *hoolitsema* visandis on näha sellised grammatilised klassid nagu subjekt (*kes hoolitseb*), sõltuvusmäärus (*hoolitseb kelle eest*) ja viisimäärus (*kuidas hoolitseb*). SkE tõi välja ka

tud-partitsiibi sagedasemad kollokaadid (*hoolitsetud välimus, aed* jne). Nende grammatiliste klasside (uuritava verbi puhul subjekt ja adverbiaalid) alusel võib öelda, et sagedasemad laiendid on: *kes (riik, firma, valitsus)* hoolitseb *kelle-mille eest (laste, inimeste eest)* ja *kuidas (hästi, korralikult)*.

tooma ()			EstonianRC freq = 180229 (721.7 per million)					
subject	23212	5.5	object	6608	6.3	adv ühendverb	32416	21.6
aasta	<u>243</u>	3.04	näide	<u>840</u>	7.07	kaasa	<u>13557</u>	12.09
tulevik	<u>223</u>	6.42	selgus	<u>177</u>	8.76	välja	<u>6604</u>	9.63
päev	<u>221</u>	4.81	edu	<u>158</u>	7.43	esile	<u>3090</u>	11.03
mees	<u>213</u>	4.8	kasu	<u>140</u>	6.42	sisse	<u>2959</u>	10.24
inimene	<u>198</u>	4.02				ära	<u>1899</u>	8.4
naine	<u>117</u>	4.56				tagasi	<u>1615</u>	8.62
elu	<u>114</u>	4.51				juurde	<u>659</u>	8.29
adverbial sisseütlev 2875 13.4			verb pp ette 536 126.1					
meel	<u>260</u>	7.63	avalikkus	<u>149</u>	7.84			
kodu	<u>260</u>	6.43	adverbial rajav 413 24.5					
haigla	<u>102</u>	5.6	vaataja	<u>151</u>	8.32			
adverbial saav 3414 20.8			adverbial alaleütlev 6276 15.3					
näide	<u>1571</u>	8.0	turg	<u>747</u>	8.01			
põhjus	<u>411</u>	6.91	lava	<u>374</u>	8.01			
põhjendus	<u>254</u>	9.32	päevavalge	<u>202</u>	9.45			
võrdlus	<u>164</u>	8.88	võitja	<u>165</u>	7.11			
lisa	<u>141</u>	4.91	koha	<u>161</u>	5.08			
ettekääne	<u>126</u>	9.32	ilm	<u>160</u>	7.21			
eeskuju	<u>125</u>	7.74	ekraan	<u>143</u>	7.65			
modifies past participle 10399 9.1								
näide	<u>334</u>	5.72						
tabel	<u>295</u>	8.5						
andmed	<u>196</u>	5.25						
väide	<u>180</u>	6.49						
nõue	<u>115</u>	5.83						

Joonis 9. Verbi *tooma* sõnavisand.

Verbi *tooma* visandis on eristatud subjekt (*kes* toob), objekt (*mida* toob), adverbiaalid: 1) latiivne kohamäärus (*tooma kuhu* (*adverbial_alaleütlev*, *adverbial_sisseütlev*)); 2) latiivne valdajamäärus (*tooma kelleni*, *kelle ette* (*adverbial_rajav*; *verb_pp_ette*), sagedasemad on *tooma vaatajani*, *tooma avalik-kuse ette*) ja 3) otstarbemäärus (*tooma milleks* (*adverbial_saav*), sagedasem on *tooma näiteks*). Seega verbi *tooma* sagedasemateks laienditeks on subjekt, objekt ja adverbiaalid, mis markeerivad sihtkohta, valdajat ning otstarvet. Peale selle näeb verbi *tooma* visandis ka sagedasemaid ühendverbe: *kaasa tooma*, *välja tooma*, *esile tooma* ja *sisse tooma*.

2.4.1. Verbide sõnavisandid kui lausemallide tuvastamise vahend

Kas SkE annab infot otseselt lausemallide kohta või pakub ta vaid sagedasemate laiendite tüüpe? Kas lisaks sõnasemantika võimaldab visand saada aimu ka konkreetse verbi lausesemantikast? Kas võib oletada, et verbide visandites esile toodud grammatilised klassid (subjekt, objekt, määrus jt) võimaldavad rekonstrueerida verbiga tähistatud tegevussituatsiooni ja selle komponente?

Sõnatähenduse ja tegevussituatsiooni vahekorra üle on arutlenud Huno Rätsep (1978: 237–243), kes on tõdenud, et keeliliit võivad tegevussituatsiooni komponendid lauses avalduda situatsiooni kesksest sõnast sõltuvalt eraldi sõnades. Sellisel juhul on lause või sõnaühendi keskse sõna tähendus üldisem ja tegevussituatsiooni mõistele kõige lähedasem. Seejuures rõhutab Rätsep (1978: 241), et eesti keeles esitatakse tegevussituatsiooni komponendid pindstruktuuris eraldi sõnadena tegevussituatsiooni mõistet märkiva sõna kõrval.

Teiselt poolt arvatakse, et semantilisi rolle tuleb grammatilistest rollidest ehk lauseliikmetest lahus hoida, sest seost grammatilise käände ja semantilise rolli vahel ei ole ning nimi-sõna käändevormist ei saa kunagi üheselt järeldada tema semantilist vormi (vt nt Pajusalu 2009: 81–82). Siiski märgib Renate

Pajusalu (*ibid.*), et kuigi semantilistel rollidel ei ole grammatilise vormistusega üksühest seost, pole vorm kunagi tähendusetu.

Allpool proovime lemma *jalutama* näitel rekonstrueerida konkreetse verbi tüüpilise lausemalli, toetudes sõnavisandi statistilistele andmetele. Seejuures eristame lausemalli tuletamisel järgmisi etappe: 1) semantiliste rollide tuvastamine grammatiliste suhete põhjal, mida sõnavisand ette annab; 2) rollide taga olevate entiteetide semantiline analüüs ja kategoriseerimine; 3) lausemalli tuletamine.

jalutama ()			EstonianRC freq = 9114 (36.5 per million)					
subject	612	3.1	adv ühendverb	1129	16.3	pp verb	933	11.3
mees	<u>50</u>	2.85	ringi	<u>493</u>	9.8	mööda	<u>193</u>	8.25
inimene	<u>43</u>	1.9	mööda	<u>135</u>	7.72	vahel	<u>83</u>	6.19
isa	<u>14</u>	3.44	läbi	<u>122</u>	5.44	poole	<u>56</u>	6.7
naine	<u>13</u>	1.6	välja	<u>94</u>	3.67	ääres	<u>53</u>	8.46
koer	<u>12</u>	4.13	edasi	<u>58</u>	4.88	peale	<u>49</u>	4.81
noormees	<u>10</u>	3.89	vastu	<u>51</u>	3.69	läbi	<u>43</u>	3.94
poiss	<u>7</u>	2.25	ütle	<u>32</u>	2.98	koos	<u>42</u>	4.95
PP_X			adverbiaal seesütlev			adverbiaal alalütlev		
	183			204	11.9		162	10.3
verb_pp_mööda	<u>51</u>	513.3	park	<u>40</u>	6.58	tänav	<u>36</u>	3.27
verb_pp_poole	<u>26</u>	75.0	mets	<u>21</u>	4.45			
verb_pp_vahel	<u>21</u>	140.2	linn	<u>19</u>	1.95			
verb_pp_läbi	<u>20</u>	66.9	vanalinn	<u>11</u>	4.6			
verb_pp_juurde	<u>14</u>	35.7						
verb_pp_ääres	<u>14</u>	207.4						
verb_pp_koos	<u>12</u>	21.6						
adverb			adverbiaal kaasütlev			ÜHENDVERB		
	1034	3.7		104	9.4		66	
lihtsalt	<u>54</u>	5.86	koer	<u>20</u>	4.9	pp_ringi_ühendverb	<u>45</u>	474.6
kui	<u>51</u>	2.59	laps	<u>10</u>	0.77	pp_välja_ühendverb	<u>21</u>	15.3
väljas	<u>35</u>	8.17	poeg	<u>7</u>	3.02			
seal	<u>31</u>	4.52						
niisama	<u>28</u>	7.49						
rahulikult	<u>26</u>	7.71						
siis	<u>21</u>	2.68						

Joonis 10. Verbi *jalutama* sõnavisand.

Sõnavisand aitab tuvastada eelkõige verbiga seotud sagedasemad grammatilised klassid, milleks on subjekt (612 esinemisjuhtu) ja eri tüüpi määrused: kaasnemismäärus (104 esinemisjuhtu); lokatiivne kohamäärus (204 esinemisjuhtu); latiivne kohamäärus, kuhu kuuluvad ka *jalutama millegi ääres-tüüpi määrused*; prolatiivne kohamäärus. Seega lähtuvalt esile tulnud grammatilistest klassidest võib selle verbi tegevussituatsioonis eristada järgmisi komponente ehk semantilisi rolle: AGENT, KAASLANE, KOHT, TEE.

Entiteetide analüüs kategooriate järgi näitab seda, et AGENDIKS on *jalutama*-verbi puhul enamasti *inimene*, harvem *loom*; KAASLASEKS on samuti elusolend, kas *loom* või *inimene*; KOHA sagedasemateks täitjateks on substantiivid *park*, *mets*, *linn ja meri*; TEED märkis sageli substantiiv *tänav*.

Niisiis on statistiliste andmete põhjal *jalutama*-verbi sagedasemad lausemalli komponendid AGENT, KOHT ja KAASLANE, seega võiks sõnastikus välja tuua vastavad verbilaiendid.

2.5. Sõnavisandite võrdlus

SkE lisafunktsioon *Sketch-Diff* võimaldab võrrelda kahe lemma sõnavisandeid. Eelkõige on see funktsioon mõeldud osasünonüümide võrdlemiseks ja eristamiseks. Süsteem toob esile sarnasusi ja erinevusi kahe semantiliselt seotud sõna süntagmaatilises käitumises, st esitab nii sõnade ühised kui ka vaid ühe sõnaga kokkukuuluvad kollokaadid. Näitena esitame katke komplementaarsete vastandite *naine* ja *mees* sõnavisandite võrdlusest (joonis 11).

SkE eristab kollokaate, mis kuuluvad sagedamini kokku lemmaga *naine* (näiteks *rase naine*, *naine sünnitab*, *naist võtma*), ja kollokaate, mis esinevad tihedamini koos lemmaga *mees* (näiteks *tundmatu mees*, *meest süüdistama*, *mees ja naine*). Heledal taustal on atribuudid, mis esinesid enam-vähem võrdselt mõlema lemmaga. Erinevuste esiletoomisest arvestab programm ka sõnade üldist esinemissagedust korpuses.

naine/mees					EstonianRC freqs = 179449/307624				
naine	6.0	4.0	2.0	0	-2.0	-4.0	-6.0	mees	
a_modifier	23848	43752	1.9	2.2	subject_of	26160	52914	4.6	5.9
tundmatu	122	905	6.8	9.1	tunnistama	240	619	6.5	7.5
vana	781	1441	8.0	8.6	istuma	163	485	6.7	7.7
keskealine	381	708	8.9	9.0	surema	265	590	7.7	8.1
noor	3647	4435	10.0	10.0	rääkima	645	929	7.4	7.7
rased	348	0	8.8	0.0	sünnitama	203	0	7.8	0.0
object_of	2034	2789	1.4	1.3	ja/või	7299	3437	2.3	0.7
süüdistama	0	116	0.0	8.0	naine	0	1753	0.0	8.6
võtma	130	0	4.2	0.0	laps	212	0	5.1	0.0
					mees	5225	0	9.5	0.0

Joonis 11. Lemmade *naine* ja *mees* sõnavisandite võrdlus.

3. Sketch Engine'i eestikeelse mooduli edasiarendused

Sõnavisandite kvaliteet sõltub vahetult morfoloogilise ühestamise täpsusest. Praegu esineb morfoloogilise märgendamise vigu näiteks sõnaliikide tasandil ning käändsõnade puhul ka käände tasandil. Samuti on üldnimisõnadeks märgendatud isikunimesid (nt Pilv, Jänes, Rebane). Vigu tuleb ka sellest, et sisendkorpuse mitme tõlgendusega sõnade puhul jäi korpuse SkE jaoks ettevalmistamise järgus sisse ainult see tõlgendus, mis oli esikohal. Näiteks esialgu oli verbil *kukutama* kaks vormi: *kukku+tama*; *kukuta+ma*, kuid SkE-sse on jäänud neist vaid esimene.

Väljundit mõjutab oluliselt ka korpuse tekstide valik (koondkorpuses on praegu u 75% tekstidest ajakirjandustekstid).

Seega tuleb Sketch Engine'it edasi arendades täiendada korpuse sisu ja parandada ühestamise kvaliteeti. Tulevikus võiks korpus olla rohkem tasakaalus ja seeläbi representatiivsem.

Teine arengusuund on sõnavisandite grammatika täiendamine ja selle metakeele (näiteks grammatiliste kategooriate nimetuste) lihtsustamine. Visandite kvaliteeti võiks parandada ka süntaktiliste märgendite kasutamine sõnavisandite gramma-

tika kirjutamisel. See aitaks lahendada nt vaba sõnajärje probleemi. Tuleks püüelda ka semantilise ühestaja kasutamise poole, mida eesti keele jaoks praegu veel ei ole, kuid mida ingliskeelses Sketch Engine'is juba katsetatakse (McCarthy ja Reddy 2010).

Kuna sõnavisandid annavad ülevaate sõna süntagmaatilistest suhetest, saaks Sketch Engine'it edukalt kasutada eesti keele kui teise keele õppimisel ja õpetamisel. Sarnaselt õppesõnastikega toob SkE eksplitsiitselt esile sagedasemad rektsioonistruktuurid ja leksikaalsed kollokatsioonid ning aitab seega õppijat eelkõige just eestikeelse teksti koostamisel. Sõnavisandites näidatud rektsioonistruktuurid peaksid aitama vältida teise keele mõjudest tingitud vigu. Ingliskeelseid sõnavisandeid on juba proovitud automaatselt veebisõnastikuks pöörata (Hvelplund 2011). Tulemuseks on tavakasutajale lihtsasti mõistetav kasutajaliides, mis võimaldab sõnavisandeid vaadata elektroonilise sõnastiku vormingus. Sõnavisandite ja parimate näitelauseste valimise funktsiooni kombineerides toimib SkE kollokatsioonide täisautomaatse sõnastikuna. Sellest lähtudes on meie eesmärk ka programmi eestikeelse mooduli veebisõnastikuks pööramine. Sellisel moel saadud n-õ automaatsõnastiku sihtgrupiks oleksid nii eesti keele kui teise ja/või võõrkeele õppijad kui ka tõlkijad, eesti keele uurijad ja kõik teised keelehuvilised.

4. Kokkuvõte

Artiklis tutvustasime leksikograafilise tarkvara Sketch Engine eestikeelse mooduli põhifunktsioone ja rakendamise võimalusi. Konkordantse ja sõnavisandeid saab kasutada eesti keele uurimisel ja kirjeldamisel, (elektrooniliste) sõnastike ja sagedusgrammatikate koostamisel ning eesti keele kui teise ja/või võõrkeele õpetamisel. Sõnavisandites tulevad esile eesti substantiivide, adjektiivide ja verbide sagedasemad süntagmaatilised seosed (rektsioonistruktuurid ja leksikaalsed kollokatsioonid), verbivisandeid saab kasutada ka lausemallide tuletamisel. Sketch Engine on üks esimesi eesti keelele kohandatud korpusleksikograafilisi programme, mille kasutajaskond

ulatub leksikograafide hulgast väljapoole. Kuna korpuspõhisus on tänapäeva keeleteaduse levinumaid arengusuundi, siis on selliste programmide edendamine eesti keele uurimise seisukohast väga oluline.

Address:

Jelena Kallas, Maria Tuulik, Madis Jürviste
Eesti Keele Instituut
Roosikrantsi 6
10119 Tallinn, Eesti

E-mail: jelena.kallas@eki.ee, maria.tuulik@eki.ee,
madis.jyrviste@eki.ee

Kirjandus

- Atkins, Sue, Charles J. Fillmore, and Christopher R. Johnson (2003) “Lexicographing evidence: selecting information from corpus evidence”. *International Journal of Lexicography* 15, 251–280.
- Atkins, Sue, Adam Kilgarriff, and Michael Rundell (2010) “The DANTE database (Database of analysed texts of English)”. *Proceedings of the XII EURALEX International Congress* (Leeuwarden, 6–10 July 2010), 167. Leeuwarden: Fryske Akademy.
- EKG 1993 = Erelt, Mati, Reet Kasik, Helle Metslang, Henno Rajandi, Kristiina Ross, Henn Saari, Kaja Tael ja Silvi Vare (1993) *Eesti keele grammatika II. Süntaks*. Tallinn: Eesti Teaduste Akadeemia Keele ja Kirjanduse Instituut.
- Hvelplund, Holger (2011) Using Sketch Engine with IDM’s DPS for online dictionaries. Seminari materjalid <<https://trac.sketchengine.co.uk/wiki/SKEW-2/Program>>. Vaadatud 01.05.2011.
- Kallas, Jelena ja Maria Tuulik (2011) „Eesti keele põhisõnavara sõnastik: ajalooline kontekst ja koostamispõhimõtted”. *Eesti Rakenduslingvistika Ühingu aastaraamat* 7, 59–75.
- Kilgarriff, Adam, Pavel Rychly, Pavel Smrž, and David Tugwell (2004) “The Sketch Engine”. *Proceedings of the 11th EURALEX International Congress* (Lorient, 6–10 July 2004), 105–117. Lorient: Université de Bretagne-Sud.
- Kilgarriff, Adam, Milos Husák, Katy McAdam, Michael Rundell, and Pavel Rychly (2008) “GDEX: Automatically finding good dictionary examples in a corpus”. *Proceedings of the XIII EURALEX International Congress* (Barcelona, 15–19 July 2008), 105–117. Barcelona: IULA, Documenta Universitaria.

- Kilgarriff, Adam (2009) "Corpora in the classroom without scaring the students". *Proceedings 18th International Symposium on English Teaching*. Taipei.
- Kitsnik, Mare (2006) „Keelekorpused ja võõrkeeleõpe”. *Eesti Rakenduslingvistika Ühingu aastaraamat 2*, 93–107.
- Langemets, Margit, Mai Tiits, Tiia Valdre ja Piret Voll (2010) „In spe: üheköiteline eesti keele sõnaraamat”. *Keel ja Kirjandus* 11, 793–810.
- McCarthy, Diana and Siva Reddy (2010) Semantic tagging. Seminari materjalid <<https://trac.sketchengine.co.uk/wiki/SKEW-2/Program>>. Vaadatud 01.05.2011.
- Pajusalu, Renate (2009) *Sõna ja tähendus*. Tallinn: Eesti Keele Sihtasutus
- Roosmaa, Tiit, Mare Koit, Kadri Muischnek, Kaili Müürisep, Tiina Puolakainen ja Heli Uibo (2001) *Eesti keele formaalne grammatika*. Tartu: Tartu Ülikooli Kirjastus.
- Rätsep, Huno (1978) *Eesti keele lihtlausete tüübid*. Tallinn: Valgus.
- Svensén, Bo (2009) *A handbook of lexicography: the theory and practice of dictionary-making*. Cambridge: Cambridge University Press.
- Uiboed, Kristel (2010) „Statistilised meetodid murdecorpuse ühendverbide tuvastamisel”. *Eesti Rakenduslingvistika Ühingu aastaraamat 6*, 307–326.
- Vider, Kadri (1999) *Sagedamad eesti verbid semantilises andmebaasis*. Käsitirjaline magistritöö. Tartu: Tartu Ülikool. Eesti keele osakond.

Abstract. Jelena Kallas, Maria Tuulik, and Madis Jürviste: Estonian language module of lexicographic software Sketch Engine.

The Sketch Engine is a software aiming to analyse corpus data. Its main functions being the compilation and varied analysis of concordances, creating frequency lists, automatic compilation of word sketches and distributive thesauruses.

In the autumn of 2010 the Institute of the Estonian Language started to develop, in collaboration with Lexical Computing Ltd., the Estonian language module for the Sketch Engine software. The Estonian Sketch Engine consists of a 250-million word Estonian Reference Corpus.

The article gives an overview of Estonian Word Sketch grammar compilation principles. The authors demonstrate word sketches for nouns, verbs and adjectives. Also, an insight is given into the possibilities of further development of the Estonian language module.

The corpus-based approach being one of the most widespread tendencies in modern linguistics, it is of utmost importance for Estonian language research to promote the usage of programs such as the Sketch Engine.

Keywords: corpus lexicography, concordance, word sketch, thesaurus, Estonian

