

Korpuslingvistiline lähenemine eesti internetikeele automaatsele morfoloogilisele analüüsile

Heiki-Jaan Kaalep, Kadri Muischnek, Raul Sirel

Ülevaade. Käesolevas artiklis analüüsitakse eesti uue meedia keelekasutuse e internetikeele automaatse morfoloogilise analüüsi kerkivaid probleeme ja esitatakse meetodid nende lahendamiseks.

Võtmesõnad: arvutilingvistika, korpuslingvistika, automaatne morfoloogia, morfoloogia, ortograafia, sõnaliigid, eesti keel, jututoad

Sissejuhatus

Morfoloogiline analüüs on eestikeelsete tekstide korpuslingvistilisel ja/või automaatse analüüsil oluline etapp. Selle käigus lisatakse igale tekstisõnale tema algvorm e lemma ja info tema grammatiliste kategooriate kohta. Lemma kaudu saab tekstisõna ühendada sõnastikuga, grammatiliste kategooriate kaudu võrrelda eri tekste ja tekstiliike omavahel. Morfoloogiliselt analüüsitud tekst on sisendiks ka automaatse lingvistilise analüüsi järgmistele etappidele: süntaktilisele ja semantilisele analüüsile. Samuti on teadmine lemmade kohta abiks automaatsele infootsingule.

Käesoleva artikli eesmärgiks on analüüsida uue meedia keele erijooni, mis muudavad ta ebamugavaks kirjakeele tarvis loodud morfoloogilise analüsaatori jaoks. Uue meedia keele all on siin mõeldud interaktiivset internetikeelt ehk võrgusuhtluse keelt. Enam-vähem samas tähenduses on eesti keeles kasutatud ka terminit internetikeel (nt Soodla 2010) ja me kasutame selles artiklis neid väljendeid sünonüümsetena.

Analüüsime morfoloogiaanalüsaatorile tundmatuks jäävaid sõnavorme uue meedia korpusel ja esitame oma meetodi analüsaatori kohandamiseks uuele tekstitüübile.

Internetikeelt on Eesti keeleteaduses käsitlenud Anni Oja (2006, 2010), kelle põhitähelepanu on pööratud selle keelevariandi sotsiolingvistilistele aspektidele. Uue meedia allkeeltest on jututubade keelt ja selle taustaks olevaid suhtlusreegleid ning –tavasid analüüsinud Sigrid Salla (2002) ning veebikommentaare tekstilingvistika vaatepunktist Krista Kerge (2004). Foorumitekstide morfoloogilisi, morfosüntaktilisi ja sõnamoodustuslikke erijooni on uurinud Karin Soodla (2010).

Artikli ülesehitus on järgmine: esimeses alaosas kirjeldame uue meedia korpusi, nende koostist ja märgendust ning selgitame lühidalt morfoloogilise analüüsi protsessi ja sellega seotud mõisteid. Teises osas analüüsime morfoloogiaanalüsaatorile uue meedia tekstides tundmatuks jäänud sõnavorme. Artikli kolmandas osas käsitleme tundmatute sõnade analüüsi võimalikke meetodeid ja neljandas osas uurime täiustatud morfoloogilise analüsaatori töö tulemust.

1. Materjal: uue meedia korpus ja morfoloogiline analüsaator

1.1. Uue meedia korpus

Tartu Ülikooli Koondkorpus on umbes 200 miljoni sõna suurune tänapäeva kirjaliku eesti keele elektrooniline tekstikogu. Selle üks allosa on 22 miljoni sõna suurune uue meedia keelekasutuse allkorpus. Uue meedia korpus sisaldab omakorda nelja allkorpusel: jututubade tekste u 7 miljonit sõna, uudisgruppide tekste u 8 miljonit sõna, foorumitekste u 5 miljonit sõna ja kommentaaride tekste u 2 miljonit sõna.

Nimetatud allkorpuste märgendus on erinev. Muudes Koondkorpuse allkorpustes (peamiselt ajalehtede, ilukirjanduse ja teaduse tekstid) on tekstiüksustena märgendatud lõigud, nende allosadena laused ja omaette üksustena veel pealkirjad ja autorikirjed.

Kommentaari allkorpus on märgendatud enam-vähem samade põhimõtete järgi, jututubade, foorumite ja uudisgruppide allkorpuste märgendamine lähtus aga tõdemusest, et jututoasalvestus või uudisgrupi ja foorumi arhiiv on nagu näidendi üleskirjutus: tegelased tulevad lavale, esitavad oma repliigid ja lahkuvad sealt. Kõik kõnelejad st kasutajanimed ehk pseudonüümid on tähistatud märgendiga <speaker> ja kasutaja kirjutatud tekst on ümbritsetud lõigumärgenditega <p>. Kõneleja kasutajanimi ja tema toodetud tekst on ümbritsetud veel märgendiga <sp> (1). Jututoad sisaldavad ka automaatteateid jututoas osalejate liitumise või lahkumise kohta (2) ja jututoas osalejate kommentaare oma tegevuse kohta (3). Nii jututoas osaleja lausung kui ka automaatsete ja kasutajate loodud kommentaaride juures on märgendiga <time> märgendatud lausungi või kommentaari lisamise aeg.

(1) <p> <time> 16:26 </time> </p><sp> <speaker> rohi </speaker> <p> ok aga nyid minek </p> </sp>

(2) <p> <time> 16:32 </time> </p> <stage> * tessa has joined #Kreisiradio </stage>

(3) <p> <time> 13:20 </time> </p> <stage> * Camille muigab </stage>

Kommentaari, foorumite ja uudisgruppide allkorpused on automaatselt lausestatud (st lisatud on lausepiiride märgendid) kirjakeele normidest lähtuvalt, mis tähendab seda, et kui lause ei alga suure algustähega, siis ei ole seda eraldi lausena märgendatud.

Omaette probleemiks on aga laused jututubade tekstides. Automaatselt lausestada pole neid tekste õnnestunud, sest jututubade tekstides ei kasutata lause alguses suurtähte ja ka lauselõpumärgid võivad puududa. Morfoloogiline ühestaja vajab edaspidi aga teadmist lausepiiride kohta. Otsustasime sel otstarbel võrdsustada lõigu lausega. See lihtsus vastab sageli tõele, jututubade lõik, st ühe kasutaja postitus, on tüüpiliselt lühike lause või (nimisõna)fraas. Kuna jututoa lausega võrdsustatud üksus ei vasta tüüpilise lause tunnustele (vt EKK: 429-430), kasutame selle kohta siin artiklis nimetust lausung.

Kõigis uue meedia allkorpustes on nendes esinenud meili- ja internetiaadressid asendatud märgendiga <gap desc='hüperlink' />. Korpust tehes üritati ka võõrkeelne tekst korpusest välja jätta ja asendada märgendiga <gap desc='võõrkeelne_tekst' />, kuid nagu edaspidises analüüsis selgub, on tekstidesse ikkagi jäänud hulgaliselt võõrkeelset materjali.

Automaatse morfoloogilise analüüsi sisendiks on nendes korpustes ainult kasutajate tekst, mitte kasutajanimed või automaatteated (2). Samas kuuluvad kasutajate kommentaarid oma tegevuse kohta (3) samuti kasutaja loodud teksti hulka, kuid need on käesoleval etapil morfoloogilisest analüüsist välja jäänud.

Jututubade korpuses oli morfoloogilise analüüsi tarvis vaja teha veel eeltööd, et sisend oleks mõistlik. Sõnade eraldamiseks on jututubades kasutusel tühikud, aga ka kirjavahemärgid. Viimased on sageli järgneva ja emotikonid eelneva sõnaga kokku kleepunud; emotikonid ja kirjavahemärgid on samuti sageli koos. Seetõttu tuli enne morfoloogilist analüüsi sõnad kirjavahemärkidest ja emotikonidest eraldada.

Üks jututubades esinev nähtus on läbustajad. Need on osalejad, kelle eesmärgiks näib olevat teiste osalejate segamine ja vihastamine, milleks nad sisestavad võimalikult väikese vahega mõttetuid märgijadasid ja/või fraaside mitmekordseid kordusi. Tavaliselt õnnestub läbustajal oma postitust paar korda korrata, enne kui ta jututoast välja visatakse. Automaatselt on võimalik selliseid teateid ära tunda selle järgi, et nad on tavalisest postitusest pikemad, sisaldavad pikki korduvaid fraase ja/või korduvad ebatavaliselt palju (lõika-ja-kleebi meetodit kasutavad oma postitustes küll ka jututubades osalejad ise, aga läbustajad kordavad end

oluliselt rohkem). Jätsime sellised postitused morfoloogilisest analüüsist kõrvale ning ei arvestanud neid ka korpuse mahu arvutamisel.

1.2. Morfoloogiline analüsaator

Tekstide morfoloogiliseks analüüsiks kasutasime OÜ Filosofti morfoloogilist analüsaatorit *etmrf*. Tegemist on sama firma programmi ESTMORF (Kaalep, Vaino 2000) edasiarendusega. Võrreldes ESTMORFiga on *etmrf*il suurem leksikon – 71 000 sõna – ja ka mõneti kvaliteetsem liitsõnade ja produktiivsete tuletiste analüüsialgoritm. Programmi demoversioon on tasuta kasutatav veebiaadressil www.filosoft.ee.

Morfoloogiline analüsaator „näeb“ tekstis korraga ainult ühte sõnavormi ja lisab sellele kõik võimalikud analüüsid konteksti arvestamata. Antud kontekstis ainuõige tõlgenduse väljavalimist nimetatakse morfoloogiliseks ühestamiseks ja sellest käesolevas artiklis juttu ei tule.

Meie jaoks oli oluline, et *etmrf*i käitumist morfoloogilise analüsaatorina saab muuta, kui anda talle sobiv kasutajasõnastik. Kasutajasõnastik on tekstifail, milles igal real on nii analüüsiv sõnavorm kui ka väljund, mille *etmrf* meie arvates peaks antud sõnale andma. Iga analüüsimist vajava sõna korral kontrollib *etmrf* kõigepealt kasutajasõnastikust, kas sobiva analüüsi saaks võtta otse sealt. Alles juhul, kui sealt vastust ei saa, minnakse tegema morfoloogilist analüüsi. Seega saab kasutajasõnastikku panna sõnu, mida *etmrf* muidu analüüsida ei suudaks, aga ka sõnu, mis meie arvates peaks konkreetsetes tekstis saama teistsuguse analüüsi kui tavaliselt.

2. Esimene katse: lihtne morfoloogiline analüüs

Esimesel katsel analüüsisime uue meedia keelt ilma analüsaatorit kohandamata. Teostasime morfoloogilise analüüsi ilma oletamise ja ilma ühestamiseta. Tundmatuks jäänud sõnade osakaal allkorpuste kaupa on esitatud tabelis 1.

Tabel 1. Tundmatu sõna analüüsi saanud märgijadade hulk ja osakaal allkorpuste kaupa.

	Maht	T %
Jututoad	7016800	27,2
Foorumid	4981121	10,3
Kommentaariumid	1987304	5,6
Uudisgrupid	6850873	11,7

Järgnevalt analüüsisime tundmatu sõna analüüsi saanud märgijadasid ja arutleme võimalike analüüsimeetodite ning morfoloogiliste märgendite üle. Kuna kõige rohkem jäi sõnavorme tundmatuks jututubade tekstides, siis moodustab nende analüüs ka suurima osa järgnevast käsitlesest, kuid tähelepanu on pööratud ka teiste allkorpuste – foorumitele, uudisgruppidele ja kommentaaridele – tundmatuks jäänud sõnavormidele.

Morfoloogilisel analüüsil tundmatuks jäänud märgijadade põhiosa jaguneb seitsme grupi vahel: partiklid, emotikonid, täheasenduste jms sihiliku kirja pildi muutustega

sõnavormid, pärisnimed, võõrkeelsed sõnavormid ja toorlaenuid, normeeritud kirjakeele seisukohalt valed (kuid kõnekeeles/murretes esinevad) sõnavormid ja lihtsalt trükivigadega sõnavormid. Järgnevalt analüüsime neid grupe lähemalt.

2.1. Partiklid

Eelkõige jututubade tekstides moodustavad tundmatuks jäänud sõnade sagedusloendi tipu lühikesed, tihti lühenenud sõnavormid, mis sageli moodustavad üksi lausungi ja mis lausungi osana toimivad üldlaiendina, nt *tre, irw, ok, kle, we* jt. Tegemist on suulise kõne sõnaliikide süsteemist tuttava sõnaliigi (või sõnaliikide komplekti) – partikliga.

Tiit Hennoste (2002: 63) on suulises keekekasutuses esinevat partiklit defineerinud kui süntaktiliselt sõltumatut sõna, mis on semantiliselt sisutühi, mis paikneb lause propositsioonilisest sisust väljaspool ja millel on eelkõige suhtluslik ja emotiivne funktsioon. Hennoste (2002: 71-72) liigitab suulises eesti keeles esinevad partiklid üksiesinevateks partikliteks, tekstipartikliteks ja toimetamispartikliteks. Üksiesinevad partiklid jagunevad Hennoste järgi omakorda dialoogipartikliteks, afektiivseteks partikliteks ja aktiivseteks suhtluspartikliteks (nn tähelepanupüüdjad). Tekstipartiklid ehk üksi mitteesinevad partiklid jagab Hennoste edasi veel piiripartikliteks ja vabalt liikuvateks partikliteks.

Jututubade tekstide morfoloogilisel märgendamisel kasutusele võetud partikli sõnaliik hõlmab kõiki Hennoste poolt välja toodud partiklite alaliike, välja arvatud vabalt liikuvad partiklid, mille alla Hennostel on liigitatud traditsioonilised lausemodaalid (nt *ikka, alles, vist, kindlasti, veel, muidugi, väga, täitsa*), mis sõnaliigiliselt kuuluvad adverbide või modaaladverbide hulka. Sellesse rühma kuuluvad sõnad saavad morfoloogiliselt analüsaatorilt adverbilise analüüsi ja me ei pidanud otstarbekaks seda muuta, vähemalt mitte ilma põhjaliku analüüsita.

Partikli sõnaliigi märgendi saavate sõnavormide hulgas on seega:

dialoogipartiklid, mis osutavad kuuldel olemist või sõnumi vastuvõtmist, nt *aa, asoo, jaja, jep, nunuh, ok* jt,

afektiivsed partiklid – keekekasutaja reaktsioonid, väljendavad kasutaja tundeid ja meeleolusid, nt *auts, irw, icc, krt, oih, wau* jt

aktiivsed suhtluspartiklid (tähelepanupüüdjad), nt *tre, tsau* (ka *sau, saux* jms), *kle* (ka *kule* ja *kuule*), *vata* (*vaata*), *ota* (*oota*) jne. Tähelepanupüüdjad on Hennoste järgi üksiesinevad partiklid, kuid jututubades kasutatakse neid peamiselt koos üttega (*kle* + kasutajanimi) ja tihti pikema lausungi sees (kuigi partikkel pole neis süntaktiliselt lauseliige), nt (4). Meie korpuse andmetel ei saa nõustuda ka Sigrid Salla (2002: 134) väitega, et jututoavestluses puudub sageli tervitus: jututoakorpuses on lisaks partiklile *tre* väga sagedased ka tervituspartiklid *tsau, sau, tsauks*, jms. Meie ja Salla uurimused põhinevad erinevate jututubade tekstidel ja on võimalik, et erinevates jututubades käibivad erinevad suhtlustavad.

(4) AL_Capone *kle* ma tulen Pärnu mul jube nälg ;D

piiripartiklid – sõnad lausungi piiri lähedal, mis osutavad, kuidas nendega algav üksus on seotud käsiloleva tekstiga, nt küsipartikkel *ve* (*we*),

toimetamispartiklid, nt *ee, hmm* jt; põhiliselt kasutatakse jututubades takerduse või mõttepöörde edasiandmiseks küll ainult mõttepunkte (5 vs 6).

(5) viimane mäng mus mu arvutis oli... hmm see oli 3 kuud tagasi

(6) ära seleta nii palju .. ehk siis kriba aeglasemini

Tuleb rõhutada, et uue meedia tekstide morfoloogilisel analüüsil saavad kasutajasõnastikust partikli analüüsi peamiselt need sõnavormid, mis ilma oletamiseta said tundmatu sõna analüüsi ja tundmatute sõnade sagedusloendis piisavalt kõrgel kohal paiknesid.

Analoogia alusel suulise keelega kontrolliti veel mõne sõnavormi kasutust ja vajadusel lisati partikli märgend (nt *kuule*, *vaata*, *oota*, *ütleme*, *ütleks*), kuid süstemaatiline jututubade (ja teiste internetiregistrite) partiklite uurimine on tegemata.

Partiklit omaette sõnaliigina on varem eristatud ka suulise kõne morfoloogilisel märgendamisel (Hennoste jt 2002), kuid suulise kõne korpuse märgendamisel lisati partikli analüüs käsitsi, mitte automaatselt. Samuti on partikli sõnaklassi kasutatud eesti murrete korpuse morfoloogilisel märgendamisel (Lindström jt 2006).

Tüüpiliselt on partikli funktsioonis olev sõnavorm morfoloogiliselt ühene, st tal on võimalik ainult üks morfoloogiline analüüs. Kuid korpuses esines ka selliseid sõnavorme, nii mitte-kirjakeelseid kui ka kirjakeelseid, mida kasutati paralleelselt nii partikli kui ka mingi muu sõnaliigi liikmena, näiteks *ütleme* (*ytleme*) või *ommik*.

(7) *ytleme* mina teen kodus lan game

(8) *ytleme* siis kenasti suure pika lause kokku

(9) *ommik* kõigile

(10) ma tavaliselt *ommik* vara läen ja öösel tulen

Sellistel nii partikli kui muu sõnaliigina kasutatavatel sõnavormidel võib olla mitu kirjavarianti, tavaliselt kirjakeelne vorm ja selle baasil tekkinud lühivorm või –vormid. Vaatame näiteks sõnavorme *vata* ja *vaata*, mis võivad esineda nii „tähelepanupüüdjana“ kui ka täistähendusliku verbivormina. Võiks oletada, et lühenenud vormi *vata* kasutatakse pigem partiklina ja kirjakeelset vormi *vaata* verbina, kuid esialgne korpuseuring seda oletust ei kinnitanud: nii *vata* kui *vaata* võivad esineda nii partikli kui ka verbivormina (11-14). Selliseid kirjakeelse vormi ja tema lühenenud variantide komplekte esines tekstides veelgi: *ota* ja *oota*; *kle*, *kule* ja *kuule* jms.

(11) *vata* siin põlvkondade konfliktiks kisup

(12) ei suuda veel nii et ei *vata* seda klaviatuuri

(13) nuh *vaata* mul on siin mitu inimest

(14) ära minu otsa *vaata*

Partikleid esineb vähemal määral ka teiste uue meedia allkorpuste tekstides. Analüüsidest uue meedia tekste partikleid sisaldava kasutajasõnastikuga, saavad partikli analüüsi 5,8% jututubade korpuse tekstisõnadest, 0,6% foorumite tekstisõnadest, 0,24% kommentaariumide tekstisõnadest ja 0,25% uudisgruppide tekstisõnadest.

2.2. Emotikonid

Emotikonid on kirjavahemärkidest kombineeritud ikoonilised märgid, mida kasutatakse internetikeeles emotsioonide väljendamiseks ja ka üneemide ja hüüundite asendajatena, nt : P , :) , :D (vt ka Salla 2002: 138). Kuna emotikon panustab internetikeele teksti oma tähendusnüansi ja võib moodustada üksi lausungi, siis oleks mõistlik neid analüüsida ja märgendada nagu sõnu, luues ka nende jaoks uue sõnaliigi. Emotikonid moodustavad jututubade tekstides 3,2% tekstisõnadest, foorumites 0,16% tekstisõnadest, kommentaaride tekstides 0,28% ja uudisgruppides 0,37% tekstisõnadest.

2.3. Täheasendused jm sihilikud kirjpildi muutmised

Jututubade keeles esinevad sümboliasendused, mille puhul ühe tähe või tähekombinatsiooni asemel kasutatakse teisi tähti või numbreid. Levinumad asendused on ff hv asemel (nt *raffas*), x ks asemel (nt *näitex*), y ü asemel (nt *küll*), c ts asemel (nt *täica*), 2 ä asemel (nt *h2sti*), 6 õ asemel (nt *h6be*), 8 ö asemel (nt *t88*). Vähemal määral esineb tekstides ka z kasutamist s asemel (nt *meez*).

Probleemi lahendamiseks on hea teada, kas sellised asendused piirduvad ainult kindla väikese hulga sagedaste sõnavormidega või on produktiivsed, st selliste asendustega on kirjutatud suur hulk väikese sagedusega sõnavorme. Esimesel juhul võime need sõnavormid lihtsalt lisada kasutajasõnastikku, teisel juhul on vaja mingit eeltöötlust, mis tundmatu sõna analüüsi saanud märgijadas asendaks ühe sümboli teisega ja prooviks uuesti automaatselt morfoloogiliselt analüüsida.

hv asendamine ff –ga näib olevat kinnistunud ainult teatud sagedaste sõnavormide eripäraks; ff-ga hv asemel kirjutatakse peamiselt sõnavorme *raffas* (mis on jututubade tekstides sage populaarse tervituse *tere raffas* tõttu), *vaffa* (*waffa*) ja *aff*, harvem ka *koff* (omastav *koffi*).

Seevastu ts asendamine c-ga ei toimu mitte ainult kindlates, sagedastes sõnavormides (nt ei *viici*, *täica*), vaid ka kogu korpusel vaid üks kord esinevates sõnavormides (nt *ülbicema*, *veremaice*). ks asendatakse x-ga samuti nii sagedastes sõnavormides *mix*, *olex*, *ex*, *plix* 'plika' kui ka korpusel 1-2 korda esinevates sõnavormides, eriti käändsõna translatiivi (*kirurgix*, *ajatäitex*) ja verbi konditsionaali muutelõpuna (*ärkax*, *tõestax*).

Täheasenduste hulka võib lugeda ka h ärajätmise sõna algusest (*ommik*, *uvitav*, *ullem* jpt).

Veel oleme siia gruppi kuuluvateks lugenud sõnumi emotsionaalse sisu rõhutamiseks kasutatud sihilikud tähe või silbi kordused, nt *hahahahaha*, *ehhhhh*.

Kirjeldatud täheasendused ja silbi- või tähekordused esinevad ka foorumite-, kommentaaride ja uudisgruppide tekstides, kuid neid on kasutatud vähemal määral.

2.4. Pärinimed

Kuigi, nagu juba öeldud, on morfoloogilise analüüsi sisendiks ainult kasutajate loodud tekst, mitte kõnelejateks või autoriteks märgendatud sõnavormid, sisaldavad tekstid ikkagi palju pärisnimesid. Jututubade tekstides on tüüpiline tundmatu sõna analüüsi saanud pärisnimi kasutajanimi, mida on kasutatud ütte funktsioonis (15). Jututubade tekstisõnadest moodustavad pärisnimed 6,5%. Väikese algustähega kirjutatud pärisnimesid kohtab sageli ka kommentaaride tekstides (16), vähem foorumites ja uudisgruppides.

(15) krizzy tule privva

(16) tead sokrates sinuga on ükskõik kellel täiesti mõtetu vaielda

Suur- ja väiketähe eristusel ei järgita jututubade tekstides normeeritud kirjakeele reegleid, selle funktsiooniks on hoopis rõhutamine ja emotsioonide väljendamine; nii kirjutatakse pärisnimed reeglina väikese algustähega. Väiketähelisi kasutajanimisid ei suudaks pärisnimedeks määrata ka mitte oletaja, st programm, mis annab morfoloogiaanalüsaatorile tundmatuks jäänud sõnavormile analüüsi lähtudes sõnavormi kirjapildist ja sõnaliikide ning muutevormide sagedusest tekstides. Isegi kui väikesetäheline pärisnimi sisaldub morfoloogiaanalüsaatori sõnastikus (suurtähelisena) jääb ta ikkagi ära tundmata (17). Ka uudisgruppides jääb palju pärisnimesid ära tundmata väikese algustähe tõttu (18). Viimastes, nagu ka foorumites esineb eriti palju võõrpärisnimesid (19).

(17) ma kuulsin et abja paluoja kultuurimaja on kuum koht

(18) Küsimus, et kas keegi eestis ka müüb (vahendab)?

(19) Ma saan aru, et minu Sierra kapoti alt võib leida mitu ...

2.5. Võõrkeelsed sõnavormid ja toorlaenud

Foorumid ja uudisgrupid on peamiselt mingi kindla huvi- või erialaga tegelevate keelekasutajate info-ja arvamusevahetuse paikadeks. Nii kasutatakse seal palju erialakeelt ja – slängi, mis on tugevalt mõjutatud vastava ala ingliskeelsetest tekstidest ja terminoloogiast.

Foorumite ja uudisgruppide tekstides moodustavadki automaatsel morfoloogilisel analüüsil tundmatuks jäänud sõnavormidest põhiosa inglise keele sõnavormid ning toorlaenulised erialaslängi sõnad ja erialaterminid. Ingliskeelne võib olla terve lausung, tavalisemad on siiski ingliskeelsed sõnavormid või –fraasid toorlaenudena ning tsitaatsõnade või –fraasidena eestikeelses lauses (20-22).

(20) kui sa diskilt bootida tahad

(21) saab ka nii kui panna SpyBotil advanced mode- > ignore products- > linnuke ette DSO Exploit ja exit

(22) osta originaalid ehk MITTE autojupp, vaid The Real Thing

Jututubade ja kommentaaride tekstides esineb samuti võõrkeelseid, peamiselt ingliskeelseid sõnavorme, kuid nende osakaal tundmatuks jäänud sõnavormide hulgas on väiksem kui foorumites ja uudisgruppides.

Võõrkeelsed sõnavormid võivad nii jututubades kui ka teiste allkorpuste tekstides esineda nii terve lausungi ulatuses (23), üksiksõnadena eestikeelse teksti sees (24); esineb ka lausungeid, kus poole peal keel vahetub (25). Kasutatakse, eriti jututubades, ka inglise keelest pärinevaid partikleid, nt *hello, fuck*, samuti on eriti jututubades kirjutatud ingliskeelset teksti häälduspäraselt (26).

(23) Ezkimo no offence but that is not true

(24) help ma nii vesinud

(25) Metaxa ma juba lootsin et sa teistsugune ja puha ,, but no !! :(

(26) luk huus taking :D

Eesti keele morfoloogiaanalüsaator ei ole mõeldudki võõrkeelsete sõnade analüüsiks, nii et lahendus oleks võõrkeelse teksti eelnev äratundmine ja sellisena märgendamine, nii et morfoloogiaanalüsaator saaks selle vahele jätta.

2.6. Normeeritud kirjakeele seisukohalt valed sõnavormid: kõnekeelsused, murdevormid, slängisõnad, lühenenud sõnad ja sõnakatked, kirjakeele vormimoodustusnõuetele mittevastavad sõnavormid

Siia rühma kuuluvateks võib lugeda järgmised grupid:

uue meedia allkeelele eripärased (uudis)sõnad, nt *loogish* (ka *logish*), *friik*, *tydo*, verbide *privama* ja *ruulima* vormid jms.

iseseisvate sõnadena käibivad sõnakatked (vt ka Salla 2002: 133 ja Soodla 116 jj), nt *suht*, *tegelt*, *norm*, *aint* jms.

gi-ki –liite kirjakeele normist erinev asetus sõna muutelõppude seas, nt *kellegil*, *kellegile*, *kellegiga*, *millegist*, *millegiga* jt

nud-partitsiibi kõnekeelne nd-lõpuline variant, nt *läind*, *surnd* jpt (vt ka Soodla 2010: 59-61)

õigekeelenormidele mittevastavad kokkukirjutised: *eiole*, *midaiganes*, *niiet*, *eksole*, *minuteada* jpt (vt ka Soodla 2010: 87-114 ja Salla 2002: 145).

muud õigekeelenormidele mittevastavad keelendid: *midagist*, *kudagi*, *mudu*, *ikkagist*, *kussa*, *mingine*, *lissalt* 'lihtsalt', *ikki*, *prääga*, *õhtast*, *jummala* jpt.

2.7. Trükivigadega sõnad

Uue meedia keelekasutus on huvitav just oma spontaansuse tõttu. Spontaansusega seotud kiirustamisega kaasnevad rohked trükivead, mida on eriti palju jututubade kui sünkroonse keskkonna tekstides, kus kiirus on oluline. Vead on aga juhuslikud ja ebasüsteemilised ja seetõttu on neid raske automaatselt tuvastada või parandada.

2.8. Vahekokkuvõte

Eelpoolkäsitletud grupe üldistades võib öelda, et uue meedia keelekasutus erineb normeeritud kirjakeelest nii oma leksika kui ka ortograafia poolest. Leksikaalsete omapärade hulka kuuluvad partiklid, emotikonid, allkeelespetsiifilised uudissõnad, lühendid ja toorlaenuid ning nn kõnekeelsused. Leksikoni koostise seisukohalt on eripärane pärisnimede suur hulk tekstides.

Ortograafiliste eripärade hulka kuuluvad nn ortograafiamängud – ühe tähe või tähejärjendi asendamine teise või teistega ning suur- ja väiketähtede kasutamine mitte ortograafiareeglitest lähtuvalt, nt pärisnime või lause alguse tähistamiseks, vaid emotsioonide väljendamiseks.

Vincent Ooi (2002: 96-98) kirjeldab ingliskeelsete jututubade tekstide automaatse morfoloogilise analüüsi katset ja rühmitab probleemseid sõnavorme, mis annab meile hea võimaluse võrrelda eesti ja inglise vastavat keelekasutust. Ingliseelsete jututubade automaatse morfoloogilise analüüsi jaoks probleemsed sõnade rühmad sarnanevad meie osades 2.1. – 2.7. kirjeldatud rühmadega. Nii nimetab Ooi sagedasemate probleemidena emotikone, diskursuspartikleid, väikese algustähga kirjutatud pärisnimesid, žanrispetsiifilisi lühendeid, mittestandardset ortograafiat ning tähekorpusi.

Analüüsides jututubade korpusid kasutatud sõnavara ja selle erinevust kirjakeelsega, jääb mulje, et valdav osa erinevustest on teadliku keelemängu tulemus. Kasutajad justkui ütleksid endale, et „siin me kirjutame nii, nagu räägime, mitte nii, nagu koolis õpetatakse“ ja kirja pildist häälduse väljalugemine on osa mängu võlust. Sellest siis vormid *nääd näed* asemel, *kellelgi kellelgi* asemel, *vata vaata* asemel, aga ka *enivei anyway* ja *jummala jumala* asemel. Ka täheasendused x ks asemel, ff hv asemel, c ts asemel, y ü või j asemel ei tundu olevat põhjustatud kiire trükkimise vajadusest, vaid samuti pigem keelemängu lustist.

3. Kasutajasõnastik ja selle automaatne täiendamine

Nagu osas 2.2. öeldud, on morfoloogilist analüsaatorit *etmrf* võimalik allkeelespetsiifiliseks kohandada kasutajasõnastiku abil, mis annab analüüsi muidu tundmatuks jäävatele sõnavormidele ja mille abil saab anda üldkeelest erineva tõlgenduse allkeelespetsiifilise kasutuse ja funktsiooniga sõnavormidele.

Lisaks kasutajasõnastikule kasutasime uue meedia korpusete analüüsi hõlbustamiseks ka teksti eeltöötlust. Näiteks kasutatakse jututubade tekstides millegi rõhutamiseks sageli tähe või silbi kordamist, nt. *eieieieiei, teereeeee*. Ehkki kordus ise on tahtlik ja tal on kommunikatiivne funktsioon, ei ole korduvate silpide või tähtede täpne arv arvatavasti oluline. Seetõttu jätsime enne morfoloogilist analüüsi teksti ühtlustamisel korduvuse küll alles, kuid teisendasime kõik pikemad kui kolm kordust kolmeks, saades seega *eieiei* ja *teereeee*.

Sarnane probleem oli emotikonide kumuleerimine, st nende mõne koostisosa mitmekordistamine väljendamaks emotsiooni tugevust, nt :))))), kusjuures korduste arv võis olla isegi suurem kui 100. Need kordused normaliseeriti samuti eeltöötluste käigus.

Oleme seisukohal, et kasutajasõnastiku loomisel tuleb arvestada tekstikorpuse ja tema sõnavara statistiliste karakteristikutega. On teada, et sõnade sagedused korpusis järgivad Zipfi seadust, mille kohaselt (lihtsustatult ja ligikaudselt) on sõna sagedus pöördvõrdelises seoses

selle sagedusega sõnade arvuga. Ehk teiste sõnadega, väga väike hulk sõnu on väga sagedased ja väga paljud sõnad esinevad väga harva. Näiteks 7 miljoni sõnalise jututubade korpuse erinevate sõnavormide arv on ümmarguselt 350 000 ja neist esineb korpuses üks kord 220 000. Ka jututubade korpuses morfoloogilise analüsaatori jaoks tundmatute sõnavormide sagedusjaotus on sarnane: erinevaid tundmatuid sõnavorme on kokku 240 000, neist üks kord esineb 160 000. Võrdluseks: eesti kirjakeele sagedussõnastiku (Kaalep, Muischnek 2002) aluseks olevas 1 miljoni sõnalises korpuses on ümmarguselt 155 000 erinevat sõnavormi, millest üks kord esineb 95 000.

Võib arvata, et internetikeele kui suhtlusvahendi sõnavara peaks järgima põhimõtet, et kui sõnavorm erineb kirjakeelsest ja on samas haruldane, siis on ta kirjakeelsest vormist mingi regulaarse teisenduse abil tuletatud. Teiselt poolt, kui sõnavorm ei ole kirjakeelsest vormist tuletatav regulaarse teisenduse abil, siis peaks ta olema sageli kasutatav, et tema tähendus ja funktsioon kasutajatel meeles püsiks. Morfoloogilise analüsaatori kohandamine seisneb siis selles, et sagedased, ebaregulaarsed sõnavormid lisada kasutajasõnastikku käsitsi, haruldasemad ja regulaarselt kirjakeelest tuletatavad aga automaatselt.

Tundmatute sõnade sagedusloendi tipus, nagu öeldud osas 2.1, olid partiklid. Nende jaoks lisati morfoloogiaanalüsaatori sõnaliikide süsteemi uus sõnaliik – partikkel märgendiga *_B_*. Partikkel on muutumatu sõna, tema algvorm on tema tekstis esinemise kuju, nt *jap*, *jep* ja *jup* on kolm erineva algvormiga tagasisidepartiklit. Sellest reeglist on ka erand: partiklites esineb sageli tähe- ja silbikordusi, need on algvormile viimisel eemaldatud, nt tekstis esines partikkel kujul *irwwwwwww*, eeltöötuse käigus sai sellest *irwww* ja selle partikli algvorm on *irw*.

Kasutajasõnastiku abil anti morfoloogilise analüüsi märgend ka emotikonidele, mille jaoks samuti lisati morfoloogilise analüsaatori sõnaliikide süsteemi uus märgend – emotikon *_E_*. Kuna emotikone on lõplik hulk ja nad on alati ühesed, on lihtne neile kasutajasõnastiku abil morfoloogiline tõlgendus anda. Kasutajasõnastikus on 100 erinevat emotikoni. Selline suhteliselt suur arv on tingitud sellest, et eeltöötuse käigus jäeti alles kuni kolme sümbolikordusega emotikonid ja need lisati kasutajasõnastikku eraldi, nii on kasutajasõnastikus eraldi kirjed :) , :) ja :)) jaoks.

Kasutajasõnastiku abil said analüüsi ka sagedasemad toorlaenud (nt adjektiiv *cool*, partiklid *ok* ja *bye*), kuid võõrkeelsete sõnade probleemi ei pidanud me õigeks kasutajasõnastiku abil lahendada, vaid tulevikus luua või leida parem keele tuvastaja, mis mitte-eestikeelsed sõnad tekstis ära tunneks ja vastavalt märgendaks.

Normeeritud kirjakeele seisukohalt valedele sõnavormidele püüti kasutajasõnastiku abil anda kirjakeelne algvorm, nt sõnavormi *mudu* algvorm on *muidu* ja vormi *kellelile* analüüsitakse nagu vormi *kellelegi*. Probleemiks on siin piiri tõmbamine ühelt poolt allkeespetsiifiliste sõnade ja teiselt poolt nn valede sõnavormide vahele, nt kas sõnavorm *plix* on kirjakeele normile mittevastav variant sõnast *plika* (millele kasutajasõnastik peaks andma algvormiks *plika*) või selle allkeele sõna (mille algvorm on *pliks*).

Ülejäänud tundmatute sõnade grupid – pärisnimed ja muudetud kirja-pildiga sõnavormid – on gruppides olulised, kuid iga grupp koosneb hulgast keskmise või madala sagedusega sõnavormidest, mille käsitsi kasutajasõnastikku lisamine on töömahukas. Nii katsetasime kasutajasõnastiku automaatset täiendamist.

Kasutajasõnastiku automaatne täiendamine toimus järgmiselt. Algul leidsime sõnad, mida *etmrf* ei suutnud analüüsida, nt *kick* ja *viici*. Seejärel teisendasime need sõnad mingil (osas 2.3. kirjeldatud) regulaarsel moel, nt asendades kõigis sõnas *c* *ts*-iga (saime *kitsk* ja *viitsi*) ja lasime *etmrf*il neid sõnu uuesti analüüsida. Kui *etmrf* sai analüüsiga hakkama, lisasime kasutajasõnastikku esialgse sõnavormi ja uue analüüsi, seega *viici viitsi+0 //_V_ o,*

//

Teine hea näide kasutajasõnastiku automaatselt täiendamisest on nud-partitsiibi kõnekeelse nd-lõpulise variandi automaatne äratundmine ja kasutajasõnastikku lisamine. Esimesel katsel tundmatuks jäänud sõnavormide hulgast eraldati nd-lõpulised, nt *istund* ja *around*. nd-lõpp teisendati nud-lõpuks ja tulemust (*istunud* ja *arounud*) lasti jälle uuesti analüüsida. Kui analüüs õnnestus, lisati sõnavorm ja tema tõlgendus kasutajasõnastikku. nd-lõpulisõnastiku sõnavorme, millele sai kasutajasõnastiku abil anda nud-partitsiibi analüüsi, oli üle viiesaja.

Keerulisemaks tegi kasutajasõnastiku automaatse täiendamise asjaolu, et samas sõnavormis võib olla kasutatud mitut teisendust, nt sõnavormides *n2ind* või *viicix*.

Regulaarsed teisendused, mida rakendasime, olid korduvate tähtede asendamine (kolme asemel kaks või üks, kahe asemel üks ja konkreetset täheasendused (nt 2 asemel ä, x asemel ks), nagu on kirjeldatud eespool punktis 2.3. Sageli võimaldas selline teisendus tunda ära ka modifitseeritud versiooni korpuse-spetsiifilisest sõnast, mis oli juba varem kasutajasõnastikku lisatud.

Omaette grupi tundmatuks jäävate sõnavormide hulgas moodustavad pärisnimed, mida jututubades reeglina ja teistes allkorpustes juhuslikumalt kirjutatakse väikese algustähega. Üksjagu pärisnimesid tähistab isikuid, kes vestluses osalevad ja kelle poole pöördatakse. Eriti kehtib see jututubades, kus nimeline pöördumine on sageli ainus viis oma postituse adressaati näidata; 58% pärisnimedest ehk 3,8% kõigist tekstisõnadest selles korpuses on kasutajanimed. Õnneks on jututubade korpuses osalejate nimed (õigemini pseudonüümid) olemas ja märgendatud (26-27). Seega saab automaatselt teha osalejate nimesid sisaldava kasutajasõnastiku, kus nimedele antakse pärisnime analüüs. Juhul, kui kasutajanimi langeb kokku mõne tavalise sõnaga (26), tuleb pärisnime analüüs tavaanalüüsile lisada, et ei juhtuks nii, et sama sõna kasutus tavapärasel tähenduses jääks ilma õige analüüsita. Seejuures tuleb kasutajasõnastikku panna suurtähelised nimed ka väiketähelistena, sest jututubades on kombeks suurtähti vestluses mitte kasutada.

(26) <speaker> kaabakas </speaker> <p> suusatamine on mingi gei-sport </p>

(27) <speaker> Frode_Estil </speaker> <p> kaabakas milline ei ole gei sport ? </p>

Ülaltoodud näite puhul lisanduvad kasutajasõnastikku kolm kirjet:

Frode_Estil Frode_Estil+0 //_H_ sg n, //

frode_estil frode_estil+0 //_H_ sg n, //

kaabakas kaabakas+0 //_S_ sg n, // kaabakas+s //_S_ sg in, // kaabakas+0 //_H_ sg n, //

Automaatne kasutajasõnastiku täiendamine suurendas jututubade sõnastiku 30 000 sõnani, kuid kasutajanimed on sellest arvust välja jäetud. Ühe jututoa kasutajate nimed lisati ajutiselt kasutajasõnastikku sellesama jututoa töötlemise käigus ja neid järgmise jututoa analüüsil ei kasutatud. Põhjuseks oli see, et kasutajanimed on väga muutuv klass, mõned neist langevad kokku tavaliste sõnadega (nt. *keegi*) ja seega võiksid nad kaasa tuua mitmete analüüsides mõtetu kasvu seal, kus vastavate nimedega kasutajaid ei olegi.

Automaatselt loodud kasutajasõnastiku väärtus omaette, uute tekstide analüüsiks, ei ole arvatavasti kuigi suur: temas on palju sõnavorme, mis esinesid ainult selles korpuses, mille põhjal ta tehti, ja uues korpuses olevaid sõnavorme ta ei hõlma. Pigem võib väärtuseks pidada tema loomise meetodikat: regulaarsete teisenduste kasutamist sõnavormide modifitseerimisel ja sellele järgnevat sõnastikupõhist analüüsi, mille õnnestumise korral oletatakse, et teisenduse tulemuseks saadigi just algsele sõnavormile vastav kirjakeelne vorm. Uue korpuse analüüsil on mõtet rakendada sama meetodikat ja luua uus korpusepõhine kasutajasõnastik.

4. Teine katse: morfoloogiline analüüs eeltöötuse ja kasutajasõnastikuga

Teisel katsel analüüsisime uue meedia keelt rakendades eelmises osas kirjeldatud kasutajasõnastikku ja eeltöötust. Teostasime jälle morfoloogilise analüüsi ilma oletamise ja ühestamiseta. Tundmatuks jäanus sõnade hulk ja osakaal allkorpuste kaupa on esitatud tabelis 2. Tabeli veerg „1. katse T%“ näitab ilma kasutajasõnastikuta morfoloogilisel analüüsil tundmatuks jäänud tekstisõnade protsenti ja tabeli veerg „2.katse T%“ eeltöötuse ja kasutajasõnastikuga morfoloogilisel analüüsil tundmatuks jäänud tekstisõnade protsenti.

Tabel 2. Kasutajasõnastikuta (1. katse) ja eeltöötuse ning kasutajasõnastikuga (2. katse) morfoloogilisel analüüsil tundmatuks jäänud tekstisõnade %

	Maht	1.katse T %	2.katse T%
Jututoad	7016800	27,2	10,5
Foorumid	4981121	10,3	8,8
Kommentaariumid	1987304	5,6	4,8
Uudisgrupid	6850873	11,7	10,5

Paranemine on suurim jututubade korpuses, kus algne tulemus oli halvim. Tulemuse väikest paranemist foorumite ja uudisgruppide allkorpustes saab vähemalt osaliselt seletada ingliskeelsete sõnade rohkusega nendes tekstides.

Enamus nüüd veel tundmatuks jäävatest sõnavormidest on inglise, vähemal määral ka vene või mõne muu võõrkeele sõnad. Nüüd on ka jututubade tundmatute sõnade sagedusloendi tipus inglise keele sagedased sõnad – *the, is, to, in, my, it, what, or, of, am, go* jpt.

Muude tundmatuks jäävate sõnade analüüsil kooruvad välja veel mõned regulaarsed kirjaipildi muutused, täpsemalt täheasendused, millega kasutajasõnastiku koostamisel ja tekstide eeltöötusel polnud arvestatud. Vokaalidest on vastastikku asendatud e-d ja ä-d, rohkem esineb siiski ä kasutamist e asemel (*tära 'tere', eksolä, lähän, polä, okäi* jt) kui e kasutamist ä asemel (*jergmine, jelle, verk, reegiks*). Vähesel määral esineb ka ä kasutamist a asemel (*ärä*) ja y kasutamist ö asemel (*yelda* jt verbi *ütleva* vormid; *yysel*) ning y kasutamist i asemel (*yrw*).

Regulaarsetest kirjaipildi muutustest on veel suhteliselt sagedased ia kasutamine ea asemel verbide pidama (*pian, piab, piaks* jt) ja teadma (*tian, tiate* jt) vormides ning muutumatutes sõnades *pial ja piale, vahepial, piaeagu* ning *sialt*).

Keerulisem on lugu tähejärjendiga ää, mis võib asendada nii tähejärjendit ea (*hääd, pääle, pääst, päält, säält*), tähejärjendit äe (*pääv, nään, näab* jt verbi *nägema* vormid, *kääs* jt), tähejärjendit ähe verbi minema vormides (*ei lää, lääme, läävad* jt) kui ka tähejärjendit äi (*nääta*).

Sõna alguses puuduva h-ga sõnad on kasutajasõnastiku automaatse täiendamise tulemusena analüüsi saanud, nüüd jäävad tundmatuks veel mõned verbi *minema* vormid, kus h on ära jäetud sõna keskelt (*läed, lääme, ei läe, läeks* jt).

Suhteliselt sagedase ja regulaarse teisendusena saab välja tuua p kasutamise b asemel verbi oleviku ainsuse 3. isiku vormi tunnuse (*käip, näep, saap, tulep, akkap, tahap, vaatap, jääp* jpt) ja t kasutamise d asemel verbi oleviku ainsuse 3. isiku vormi tunnuse (*teet, olet,*

näet, võtat jt). Ka sõna algul esineb gbd asendamist kpt-ga (*tiivan, tüsgraafid, karaas, karderoob, kümnaasium, panaan, peib, pakter*), kuid need jäävad pigem nn juhuslike asenduste tasemele.

Esineb ka sõnavorme, milles kirjapildi muutused on kombineerunud ortograafiareeglitele mittevastava liitsõnamoodustusega, nt vormid *maitia*, ka *maidia, eitia, ääküll*.

5. Kokkuvõte

Ka internetis kasutatav eesti keel on eesti keel ning väärrib uurimist ning automaatset töötlemist. Käesolev artikkel keskendub aspektidele, mis eristavad uue meedia keelekasutust normeeritud ja toimetatud kirjakeelest – ortograafiale ja sõnavormide erinevusele kirjakeelsetest. Analüüsi erinevuste liike ja pakuti välja viis, kuidas olemasolevat tarkvara kohandada, et ta internetikeele töötlemisega hakkama saaks – luua (osaliselt automaatselt) morfoloogilisele analüsaatorile uus korpusepõhine kasutajasõnastik.

Interneti keelekasutust ei saa mingil juhul vaadelda ühtse allkeelena, samuti on ta ajas kiiresti muutuv nähtus. Seetõttu ei ole ühe korpuse põhjal loodud sõnastiku väärtus sellisena kuigi suur: temas on palju sõnavorme, mis esinesid ainult selles korpuses, mille põhjal ta tehti, ja uues korpuses olevaid sõnavorme ta ei hõlma. Pigem võib väärtuseks pidada tema loomise meetodikat: regulaarsete teisenduste kasutamist sõnavormide modifitseerimisel ja sellele järgnevat sõnastikupõhist analüüsi, mille õnnestumise korral oletatakse, et teisenduse tulemuseks saadigi just algele sõnavormile vastav kirjakeelne vorm. Uue korpuse analüüsil on mõtet rakendada sama meetodikat ja luua uus korpusepõhine sõnastik.

Uue meedia keelekasutus erineb normeeritud kirjakeelest nii oma leksika kui ka ortograafia poolest. Leksikaalsete omapärade hulka kuuluvad partiklid, emotikonid, allkeelsete spetsiifilised uudissõnad, lühendid ja toorlaenud ning nn kõnekeelsused. Leksikoni koostise seisukohalt on eripärane pärisnimede suur hulk tekstides.

Ortograafiliste eripärade hulka kuuluvad nn ortograafiamängud – ühe tähe või tähejärjendi asendamine teise või teistega, sõnaalgulise h ärajätmine ning suur- ja väiketähtede kasutamine mitte ortograafiareeglitest lähtuvalt, nt pärisnime või lause alguse tähistamiseks, vaid emotsioonide väljendamiseks.

Analüüsidest jututubade korpuses kasutatud sõnavara ja selle erinevust kirjakeelsega, jääb mulje, et valdav osa erinevustest on teadliku keelemängu tulemus.

Corpus-based Approach to the Automatic Morphological Analysis of Estonian Computer-mediated Communication

The article concentrates on aspects of Estonian that are different in the computer-mediated communication and the edited written language, and namely on orthography and the divergence of word forms. The article presents an analysis of these differences and proposes a way to custom an existing morphological analyser for the purpose of analysing the computer-mediated communication. The method entails the creation of a custom dictionary for the morphological analyser, largely in an automated manner.

When one analysis the word forms used in chat rooms, compared with those of edited written language, she gets a feeling that most of the differences result from conscious language play. The lexical traits of chat rooms include particles, emoticons, newly coined words, specific for the genre, acronyms, borrowings from foreign languages and colloquial words. There is a lot of play with orthography: substituting letters for similar sounding ones, omitting, lengthening and shortening letter sequences, and omitting capitalisation.

Kirjandus

EKK = Ereht, Mati, Tiiu Ereht ja Kristiina Ross 2007. Eesti keele käsiraamat. Kolmas, täiendatud trükk. Tallinn: Eesti keele Sihtasutus.

Hennoste, Tiit 2002. Suulise kõne uurimine ja sõnaliigi probleemid. – Teoreetiline keeleteadus Eestis. TÜ üldkeeleteaduse õppetooli toimetised 4. Tartu: TÜ kirjastus, lk 56-73.

Hennoste, Tiit, Liina Lindström, Olga Gerassimenko, Airi Jansons, Andriela Rääbis, Krista Strandson, Piret Toomet ja Riina Vellerind. Suuline kõne ja morfoloogiaanalüsaator. – Tähenäpüüdjä. pühendusteos professor Haldur Öimu 60. sünnipäevaks. Toim. Renate Pajusalu, Tiit Hennoste. Tartu Ülikooli üldkeeleteaduse õppetooli toimetised 3, 2002, lk 161 - 171.

Kaalep, Heiki-Jaan ja Kadri Muischnek 2002. Eesti kirjakeele sagedussõnastik. Tartu: TÜ kirjastus.

Kaalep, Heiki-Jaan ja Tarmo Vaino 2000. Teksti täielik morfoloogiline analüüs lingvisti töövahendite komplektis. – Arvutuslingvistikalt inimesele. TÜ üldkeeleteaduse õppetooli toimetised 1. Tartu: TÜ kirjastus, lk 73-100.

Kerge, Krista 2004. Veebikommentaariumi mitmetahuline maailm. – Tekstid ja taustad III Lingvistiline tekstianalüüs. Tartu: TÜ kirjastus, lk 51-73.

Lindström, Liina, Liisi Bakhoff, Mari-Liis Kalvik, Anneliis Klaus, Rutt Läänemets, Mari Mets, Ellen Niit, Karl Pajusalu, Pire Teras, Kristel Uihoaed, Ann Veismann ja Eva Velsker. Sõnaliigituse küsimusi eesti murrete korpuse põhjal. – E. Niit (toim.) Keele ehe. Tartu Ülikoolieesti keele õppetooli toimetised 30. Tartu 2006. lk 154-167.

Oja, Anni. Eesti keel internetis. – Keel ja arvuti. Tartu ülikooli üldkeeleteaduse õppetooli toimetised 6. Tartu 2006. lk 259-267.

Oja, Anni. Sissevaateid internetisuhtlusse. – Oma Keel nr 1 2010, lk 11-18.

Ooi, Vincent 2002. Aspects of computer-mediated communication for research in corpus linguistics. – P. Peters, P. Collins, A. Smith (toimetajad). New frontiers in corpus research. Amsterdam, Rodopi, lk 91-104.

Salla, Sigrid. Jututuba kui võrgusuhtlusvorm. – R. Kasik (toim) Tekstid ja taustad. Artikleid tekstianalüüsist. Tartu 2002. lk 128-156

Soodla, Karin. Morfoloogilisi, morfosüntaktilisi ja sõnamoodustuslikke erijooni eesti internetikeeles. Magistritöö (teadusmagistritöö) Tartu Ülikool, filosoofiateaduskond, eesti keele osakond. Tartu 2010 <http://hdl.handle.net/10062/15263>