

TARTU ÜLIKOOL
FILOSOOFIATEADUSKOND
Eesti ja üldkeeleteaduse instituut

Kaarel Veskis

PARALLEELKORPUSED ARVUTILINGVISTIKAS:
LEKSIKONIDE GENEREERIMINE JA KORPUSTE
VÕRDLEMINE

Magistritöö

Juhendaja dr. Heiki-Jaan Kaalep

Tartu 2007

SISUKORD

Sissejuhatus.....	- 5 -
1. Taust.....	- 7 -
1.1. Paralleelkorpused	- 7 -
1.2. Paralleelistamine	- 8 -
1.2.1. Probleemid	- 8 -
1.2.2. Meetodid	- 9 -
1.2.3. Tasandid	- 11 -
1.3. Võrdluskorpused	- 11 -
1.4. Paralleelkorpuste kasutamine.....	- 12 -
1.5. Paralleelkonkordantsid.....	- 13 -
1.6. Vorming	- 14 -
1.7. Väikesed ja vähemuskeeled	- 15 -
2. Korpusepõhised masintõlke- ja tõlkeabisüsteemid.....	- 17 -
2.1. Statistiline masintõlge	- 18 -
2.2. Puude pangad	- 19 -
2.3. Tõlkemälusüsteemid	- 20 -
2.4. Näitepõhine masintõlge.....	- 22 -
3. Paralleelkorpuste kasutamine lingvistilistes uurimustes ja keeleõppes	- 24 -
3.1. Lingvistilised ja tõlketeoreetilised uuringud.....	- 24 -
3.2. Paralleelkorpused keeleõppes	- 25 -
4. Leksikoni loomine paralleelkorpuse baasil	- 27 -

4.1. Automaatselt loodavate leksikonide vajalikkus	- 27 -
4.2. Paralleelkorpuste ja võrdluskorpuste põhjal leksikonide genereerimise meetodid	- 28 -
4.3. Grammatilise analüüsi osa sõnastike genereerimisel	- 30 -
4.4. Sõnastiku genereerimise praktilised võimalused	- 32 -
4.4.1. Eeldused ja eeltöötlus	- 32 -
4.4.2. Lihtsa leksikograafilise abivahendi kavand	- 35 -
4.4.3. Poolautomaatsed vahendid	- 36 -
4.4.4. Täisautomaatne leksikoni genereerimine tarkvarapaketi PWA abil	- 40 -
4.5. Kakskeelse leksikoni genereerimine PWA abil paralleelkorpuste põhjal	- 41 -
4.5.1. Uppsala Word Aligner	- 41 -
4.5.2. Linköping Word Aligner	- 47 -
4.5.3. Tulemuste hindamine: ARCADE ja PWA	- 48 -
4.5.4. Eesti-inglise paralleelkorpuste põhjal genereeritud leksikonide hindamisest	- 49 -
4.5.5. Võrdlus ESTERMiga	- 52 -

5. Paralleelkorpuste võrdlemine ja paralleelistuse kvaliteedi

hindamine	- 54 -
5.1. Korpuste kirjeldused	- 55 -
5.1.1. Maht	- 55 -
5.1.2. Paralleelistused	- 56 -
5.1.3. Paralleelistusvigadele viitavad tunnused	- 58 -
5.2. Eeltöö	- 60 -
5.3. Võrdlemine	- 61 -
5.3.1. Metoodika ja algoritm	- 61 -
5.3.2. JRC Vanilla versiooni ja TÜ korpuse võrdlus	- 63 -
5.3.3. 0-vastavuste protsent	- 65 -
5.3.4. JRC HunAligni ja TÜ korpuse võrdlus	- 66 -

5.3.5. JRC-Acquis' korpuse Vanilla ja HunAligni versioonide võrdlus.....	- 68 -
6. Kokkuvõte.....	- 69 -
THE ROLE OF PARALLEL CORPORA IN COMPUTATIONAL LINGUISTICS: COMPARISON OF PARALLEL CORPORA AND GENERATION OF BILINGUAL LEXICONS FROM PARALLEL CORPORA <i>Summary</i>.....	- 73 -
Kirjandus.....	- 76 -
Lisa 1. Fragment UWA-ga TÜ paralleelkorpusest ekstraheeritud leksikonist.....	- 86 -
Lisa 2. Fragment LWA-ga TÜ paralleelkorpusest ekstraheeritud leksikonist.....	- 101 -

Sissejuhatus

1990-ndate aastate algus tähistab mitmeid märkimisväärseid saavutusi paralleeltekstide vaheliste vastavuste automaatse tuvastamise osas (nt Brown jt 1991; Gale, Church 1993). Edasiste aastate jooksul on esile kerkinud suur hulk paralleelistamisega seotud probleeme, kuid ka palju huvitavaid lahendusi nendele probleemidele. Samuti on tekkinud teadlikkus mitmesugustest uutest võimalustest, mida paralleelkujul keelekorpused võivad tähendada erinevate loomuliku keelega seotud eluvaldkondade jaoks.

Magistritöö üheks eesmärgiks on olla esimeseks eestikeelseks sissejuhatuseks paralleelkorpuste temaatika tähtsamatesse aspektidesse ja ühtlasi anda ülevaade Tartu Ülikooli üldkeeleteaduse õppetoolis toimuva paralleelkorpuste-alase töö hetkeseisust. Teiseks eesmärgiks on anda panus Eestis toimuva leksikograafiatöö arendamisse, asetades erilist rõhku paralleelkorpustes peituva leksikaalse info automaatse esiletoomise võimaluste tutvustamisele praktiliste näidete kaudu. Töö praktilise osa tulemuseks olnud leksikonid võivad kasutust leida masintõlkerakendustes või olemasolevate leksikonide täiendamist hõlbustavate vahenditena.

Magistritöö jaguneb kuueks peatükiks. Esimesed kolm peatükki ja neljanda peatüki esimesed alaosad kujutavad endast ülevaatlikku sissejuhatus paralleelkorpuste temaatikasse. Neljanda peatüki osas 4.5 ja viiendas peatükis on ülekaalus praktilise eesti keelt puudutava paralleelkorpuste-alase töö kirjeldus.

Esimene peatükk annab ülevaate paralleelkorpuste koostamisest ja kasutamisest. Teine peatükk keskendub paralleelkorpuste peamisele rakendusvaldkonnale – korpusepõhisele masintõlkele ja tõlkeabisüsteemidele. Kolmas peatükk käsitleb paralleelkorpuste kasutusvõimalusi lingvistilises uurimistöös ja keeleõppes.

Neljas peatükk¹ tutvustab üht esialgu suhteliselt marginaalset, kuid samas perspektiivikat paralleelcorpuste rakendusvaldkonda – kakskeelsete leksikonide genereerimist. Peatükk kirjeldab ka esimest teadaolevat katset genereerida inglise-eesti leksikon automaatselt inglise-eesti paralleelcorpustest ja katse tulemuste analüüsi. Püüan siin lisaks osutada erinevatele alternatiivsetele võimalustele kasutada paralleelcorpuseid leksikograafiatöös ning kirjeldada eeldusi, millest saab lähtuda ühe keelena eesti keelt sisaldavate leksikonide genereerimiseks sobiva tarkvara loomisel tulevikus.

Viendas peatükis² on kirjeldatud paralleelcorpuste võrdlemise ja hindamise alast praktilist tööd, mis kujutab endast loomulikku jätku leksikonide genereerimise teemale ning annab teisalt mitmeid vastuseid küsimustele, mida tõstatasid eesti-inglise statistilise masintõlke (Fishel jt 2007) esimesed tulemused.

Samas on tegemist uudse lähenemisega paralleelcorpuste paralleelistuskvaliteedi hindamisele – erialasest kirjandusest ei ole teada analoogseid katseid hinnata paralleelistuse kvaliteeti paralleelcorpuste võrdlemise teel poolautomaatselt.

Tahaksin tänada oma juhendajat Heiki Kaalepit, kes töö valmimisele väga olulisel määral kaasa aitas.

¹ Magistritöö neljas peatükk põhineb Eesti Rakenduslingvistika Ühingu 2007. aasta aastaraamatus (Veskis 2007) ilmunud artiklil.

² Viies peatükk rajaneb koos Heiki-Jaan Kaalepiga kirjutatud ja 2007. a septembris Kopenhaagenis toimuvale XI masintõlke-alasele tippnõupidamisele (MT Summit) stendiettekandena avaldamiseks esitatud artiklil.

1. Taust

1.1. Paralleelkorpused

Paralleelkorpus on korpus, mis sisaldab mingit teksti originaalkeeles ja selle tõlget teise keelde või tõlkeid teistesse keeltesse. Paralleelkorpuse paralleeltekstid võivad olla ka mõne kolmanda korpuse mittekuuluva teksti tõlked. Paralleelkorpuste kasutamiseks on neid korpuseid vaja eelnevalt rohkem töödelda kui tavalisi ükskeelseid tekstikorpuseid – kahe paralleelse teksti märgendus peab olema omavahel seotud.

Tuntuimaks paralleelkorpuseks on peetud Kanada *Hansardit*³. See korpus koosneb Kanada parlamendidebattidest, mida avaldatakse riigi kahes ametlikus keeles – inglise ja prantsuse keeles. Väiksemate korpuste hulgast võib näiteks tuua 1 miljoni sõnalise inglise-sloveeni paralleelkorpuse⁴. Sellistest paralleelkorpustest, kus üheks keeleks on eesti keel, tuleks mainida europrojekti Multext-East⁵ raames valminud paralleelkorpust, mis sisaldab George Orwelli romaani „1984” kaheksas keeles. Selle lausetasandil paralleelistatud korpuse maht on 75 000 sõnet.

On koostatud ka seaduste ja õigusaktide tekste ning tõlkeid sisaldavad inglise-eesti ja eesti-inglise paralleelkorpus (TÜPK) ja väiksemahulisemate ettevõtmistena eesti-rootsi, eesti-norra ja vene-eesti paralleelkorpused. Nii keeletehnoloogia kui ka kontrastiivne lingvistika vajab aga tingimata vähemalt paarkümmend miljonit sõna sisaldavat suurt eestikeelse osalusega paralleelkorpust; tõsise masintõlke-alase töö jaoks läheks vaja 100

³ <http://www ldc.upenn.edu/Catalog/CatalogEntry.jsp?catalogId=LDC95T20>

⁴ <http://nl.ijs.si/elan>

⁵ <http://nl.ijs.si/ME/>

miljoni sõna suurust paralleelkorpust (Muischnek jt 2003:25). Praegu sellises mahus paralleelkorpust, milles oleks esindatud ka eesti keel, veel ei ole.

1.2. Paralleelistamine

1.2.1. Probleemid

Probleeme tekitab paralleelistamisel keelte erinev morfoloofiline tüpoloogia (aglutinatsioon vs fleksioon või analüütilisus) ja eesti ja indo-euroopa keelte lause-ehituslikud erinevused, näiteks eesti liitsõna väljendatakse inglise keeles tihti hoopis mitmesõnalise fraasiga.

Ometi on mujal maailmas isegi inglise ja aglutinatiivsete keelte sõnatasandil paralleelistamisel teatud edu saavutatud. Seejuures ei piisa tegelikult sõnatasandist, vaid aglutinatiivsete keelte eripärast johtuvalt tuleb inglise keele sõnadega või suuremate üksustega tihti vastavusse asetada ka aglutinatiivse (aglutineeriv-flekteeriva) keele morfeeme ja sõnaosi (Martin jt 2003).

Nimetatud lahknevus inglise ja eesti keele liitsõnade vahel esineb aga sarnaselt inglise-eesti keelepaariga ka inglise ning germaani keelte vahel. Seda lahknevust on arvestatud inglise-rootsi (Ahrenberg jt 2000a) sõnatasandi paralleelistaja loomisel – kuna tegemist on keelest sõltumatu paralleelistajaga, siis võib selle inglise-rootsi keelepaariga tehtud töö võtta üheks aluseks inglise-eesti joondamisprogrammi väljatöötamisel.

Korrektne paralleelistamine on paralleelkorpuse hilisema kasutatavuse seisukohalt kõige olulisem ja ka üksnes indo-euroopa keeli hõlmava korpuse puhul kõige töömahukam probleem. Paralleelistamist raskendab lisaks keeltevahelistele erinevustele asjaolu, et tihti sisaldavad paralleeltekstid “müra”, st ühes tekstis on midagi rohkem või vähem kui temale vastavas teises tekstis, mistõttu ei saa lauseid omavahel üksüheselt kokku viia.

1.2.2. Meetodid

Paralleelistamisel ehk joondamisel eristatakse statistilisi ja lingvistilisi meetodeid, kusjuures statistilisi meetodeid peetakse tõhusamaks suuremate korpuste ja lingvistilisi meetodeid väiksemate joondamisel (Oakes, McEnery 1998). Potentsiaalselt kõige edukamaks peetakse siiski statistiliste meetodite kombineerimist lingvistiliste meetoditega ja lisaressursside (sõnastikud) kasutamise. Statistiliste meetodite populaarsus on viimastel aastatel oluliselt tõusnud ka üldisemalt keeletehnoloogias ja arvutilingvistikas seoses infotehnoloogia kiire arengu ja järjest laiema levikuga.

Tartu Ülikooli inglise-eesti ja eesti-inglise paralleelkorpuse paralleelistamiseks kasutati Gale'i ja Churchi algoritmile toetuvat statistilist Vanilla paralleelistajat⁶.

Gale'i ja Churchi algoritmil põhineb ka programmi ParaConc⁷ poolautomaatne joondamisfunktsioon. Programmi ParaConc on Vanilla kõrval peetud üheks levinumaks, „vaikimisi” kasutatavaks paralleelistamisvahendiks. (Rosen 2005)

Gale'i ja Churchi algoritm eeldab, et mõlemad paralleelistatavad tekstid koosnevad võrdsest hulgast mingil kindlal moel eristatud terviktekstist väiksematest osadest, mis on juba algselt paralleelsed. Teiseks eeldab algoritm, et tõlketeksti laused on samas järjestuses kui lähteteksti laused. Kolmandaks eeldab algoritm, et tõlke pikkus sõltub originaali pikkusest. (TÜPK)

Katset on tehtud ka inglise- ja hiina keelsete tekstide joondamisega (McEnery jt 2000). Gale'i ja Churchi algoritmi rakendamine inglise keelsete ja ladina tähestikku translitereeritud (Pinyin) hiina keelsete tekstide joondamiseks kukkus läbi, kuna mõlema keele vastava lause tähemärkide arvu vahel puudub sellisel juhul piisav korrelatsioon. Rööpsus on aga märksa suurem, kui tähemärkide arvu asemel võetakse arvesse sõnade arvu, sellest on ka ühe võimaliku paralleelistamismeetodi väljatöötamisel lähtutud (McEnery jt 2000). (Mõningast edu on samuti saavutatud kakskeelse leksikoni ekstraheerimisel inglise-hiina paralleelkorpusest (Fung 2000)).

⁶ nl.ijs.si/telri/Vanilla

⁷ <http://www.athel.com/para.html>

Ülevaate Gale'i ja Churchi algoritmi kasutavast ja teistest tekstide erinevatel tasemetel paralleelistamise meetoditest ja nendega seonduvatest probleemidest annab näiteks Jörg Tiedemann (Tiedemann 1997).

Lausetasandil paralleelistamise meetodid saab üldjoontes jagada kas lausete pikkusest, sõnade vastavusest või sõnade sarnasusest lähtuvateks ning hübriidmeetoditeks. (Singh, Husain 2005)

Lausete pikkusest lähtuvad meetodid põhinevad oletusel, et lähte- ja sihtkeele lausete pikkused on omavahel seotud. Lausete pikkuse tunnuseks on kasutatud nii sõnade arvu (Brown jt 1991) kui ka (Gale, Church 1991) tähemärkide arvu. Mõlemad lausete pikkusest lähtuvad meetodid eeldavad, et enne automaatset lausetasandil paralleelistamist on tekstid juba lõigutasandil paralleelsed.

Sõnade vastavusest lähtuvad meetodid (nt Kay 1991) rajanevad eeldusel, et sõna ja tema tõlge teises keeles paiknevad paralleelsetes tekstides positsioonilistes vastavustes. Esialgselt sõnatasandi paralleelistusest tuleb iteratiivne algoritm tõenäolise lauseparalleelistuse, mida omakorda kasutatakse sõnatasandi täielikumaks joondamiseks.

Sõnavastavuse meetodeid arendati edasi IBMi mudel 1 raames⁸ statistilise masintõlke eesmärkidel. Hübriidmeetodiks võib nimetada näiteks Robert C. Moore'i algoritmi (Moore 2002), mis toimib kahes järgus. Esimene järk baseerub lausete pikkustel ja teine IBMi 1. mudelil.

Erinevaid paralleelistusmeetodeid võrreldes on jõutud järeldusele, et paralleelistuse kvaliteet sõltub peamiselt sisendtekstide mitmesugustest omadustest: tekstidevahelistest erinevustest näiteks tabelite, jooniste või viidete osas, tekstitüübist (kas tegu on sõnasõnalise tõlkega või loomingu tõlkega või millegi vahepealsega), keelte omavahelisest sugulusest või suguluse puudumisest, sisendtekstide mahust jne. Sõltuvalt sisendtekstide omadustest on tekstide paralleelistamiseks sobivam üks või teine paralleelistusmeetod. (Singh, Husain 2005)

⁸ IBM-i mudelite kohta vt lähemalt ptk 4.4.3.

1.2.3. Tasandid

Lausetasandil paralleelistamine kui üldiselt kõige olulisem etapp võib mõnede tekstitüüpide (DVD-subtiitrid, tarkvara kasutajaliideste tõlkimisel produtseeritavad .po-failid, nummerdatud lõikudega pühakirjatekstitid jms) puhul osutada ka üsnagi lihtsalt teostatavaks, mis selliste tekstitüüpide puhul võimaldab keskenduda spetsiifiliste rakendusvõimaluste uurimisele (vt nt Simões 2004).

Lause- ja sõnatasandil tekstide joendamisele (eeskätt nn „ankurpunktidega” seostuvale G&C algoritmi edasiarendusele) on pühendatud kõige rohkem uurimistööd ja kirjutisi (nt Hofland, Johansson 1998), vähem tööd on esialgu tehtud teiste tasandite (mitmesõnalised üksused, süntaktiline ning semantilis-pragmaatiline struktuur) paralleelistamismeetodite väljatöötamiseks.

1.3. Võrdluskorpused

Lisaks paraleelkorpustele on kasutusel ka **võrdluskorpuste** (*comparable corpora*) mõiste. Võrdluskorpused võivad lisaks erinevatele keeltele sisaldada ka samakeelseid tekste (erinevates keelevariantides), mida ühendab teemavaldkond, tekstitüüp, tekstide loomise periood või mõni muu omadus.

Võrdluskorpustel on mõned kasutuseelised paraleelkorpuste ees, kuigi võrreldavate tekstide segmente ei ole võimalik omavahel paraleelkorpustele sarnaselt seostada. Tuntuim võrreldav korpus on ICE – International Corpus of English⁹, mis sisaldab 20 inglise keele varianti kas inglisekeelsetest maadest nagu Suurbritannia või USA või riikidest, kus inglise keel on teiseks ametlikuks keeleks või kõrghariduse keeleks nagu näiteks Indias. Selle korpuse põhieesmärgiks on inglise keele erinevate regionaalsete variantide kontrastiivne uurimine. (Muischnek jt 2003:18)

⁹ <http://www.ucl.ac.uk/english-usage/ice>

1.4. Paralleelkorpuste kasutamine

Nii paralleel- kui võrdluskorpustel on aga lisaks erinevate keelte või keelevariantide kontrastiivuuringutele väga palju erinevaid kasutusvõimalusi nii teoorias kui praktikas, mitmeid neist võimalustest ollakse alles avastamas. Enim levinud paralleelkorpuste rakenduslikud eesmärgid saab jagada kolme suuremasse rühma (Borin 2002:14):

(1) kontrastiivsed ja tüpoloogilised grammatikat ja leksikograafiat hõlmavad lingvistilised uurimused (vt nt Ebeling 1998);

(2) paralleelkorpuste kasutamine masintõlkesüsteemides ja tõlkeabiprogrammides (nt Melby 2000); eraldi võib nimetada mitmesuguste toodete keelelist lokaliseerimist ning internatsionaliseerimist kui masintõlke kitsama fookusega allsuunda;

(3) paralleelkorpuste kasutamine keeleõppes ja -õpetuses (nt Botley jt 2000).

Sii võiks lisada veel paralleelkorpuste põhjal sõnastike koostamis- või täiendamisevõimaluse terminite ekstraheerimise läbi (nt Bowker, Pearson 2002: 171-174)

Veidi marginaalsemate, kuid samuti huvitavate praktiliste kasutusvõimalustena võiks esile tuua mitmekeelse info-otsingu (Davis 1998) ja sõnastike käsitlemise paralleelkorpustena või paralleelkorpuste osadena, et luua uusi ja täielikumaid tõlkevahendeid (Geisler 2002) või mõnel muul eesmärgil. Sõnatähenduste ühestamiseks saab kasutada semantiliselt märgendatud paralleelkorpust (Dien 2002). Samuti on võimalik paralleelkorpust kasutada keele lingvistiliseks analüüsiks või sünteesiks vajalike vahendite automaatseks produtseerimiseks (Kuhn 2004), erinevates keeltes tekstide automaatseks kategoriseerimiseks (Gliozzo, Strapparava 2005) jne.

Sageli on paralleelkorpustealastest uurimustest kasu praktiliste rakenduste väljatöötamisel ja vastupidi, samuti saab paralleelkorpuste erinevates kasutusvaldkondades rakendada ühtesid ja samu algoritme ja meetodeid. Nii näiteks on elektrooniliste sõnastike põhjal loodud sõnade foneetiliste transkriptsioonide paralleelkorpust analüüsides ehk hääldusvariante automaatselt võrreldes võimalik mõõta sugulaskeelte (Müller 2005) või ka

ühe keele dialektide (Nerbonne, Heeringa 1997) omavahelist sugulusastet, selle ja teiste sarnaste analüüside tulemusi saab kasutada keeleõppes. Erinevate keelte või keelevariantide selline võrdlemine lähtub masintõlke jaoks välja töötatud algoritmidest, millega otsitakse sarnasusi sõnade ortograafilisel tasandil. Paralleelkorpused võiksid ka hõlbustada terminoloogide tööd ülevaate saamisel mingi valdkonna olemasolevast sõnavarast ning uute terminite loomisel (Trosterud 2002:120).

Lynne Bowker ja Jennifer Pearson (2002: 94-95) nimetavad kolme põhilist inimgruppi, kes peaksid olema huvitatud paralleelkorpuste kasutamisest. Need kolm gruppi on (a) keeleõppurid/-õpetajad, (b) kõik tõlkimisega kokku puutuvad inimesed ja (c) keeletehnoloogid/arvutilingvistid.

Arvutilingvistide põhihuvi on suunatud paralleelkorpuste kasutamisele treeningmaterjalina, mille põhjal välja töötada ja arendada joondamistarkvara. Mida parem tarkvara, seda suuremaid paralleelkorpuseid on võimalik koostada. Suurte paralleelkorpuste tähtsaim rakendusvaldkond on masintõlkesüsteemid.

1.5. Paralleelkonkordantsid

Esmane paralleelkorpuste kasutusfunktsioon seisneb võimaluses leida sõnade või fraaside paralleelkonkordantse erinevates keeltes (st tekstis või tekstides esinevate sõnade loendeid koos oma vahetu kontekstiga ja paralleelkontekstiga). Selleks ei ole tingimata vaja spetsiaalset tõlkeabitarkvara, kuid tõlkeabiprogrammid võimaldavad ka konkordantside leidmist, nagu ka tekstide paralleelistamist.

Tõlkeabiprogrammides esitatakse paralleelkonkordantse tõlgitava lause (lauseosa) täpse või hägusa (*fuzzy*) sarnasuse põhjal tõlkemälukorpuses olevate lausete või lauseosadega. Tõlkevastete leidmiseks tõlkemälust võib kasutada nii statistilisi meetodeid kui ka lausestruktuurianalüüsi (Gaussier jt 2000).

On olemas ka peamiselt paralleelkonkordantside esitamisele suunatud tarkvara (Para-Conc) ja samuti on välja töötatud spetsiaalseid lingvistilise suunitlusega tekstitöötlusprogramme, mis sisaldavad funktsioone tekstide joondamiseks, paralleeltekstide töötlemiseks, uurimiseks ja vajaliku info ekstraheerimiseks paralleelkorpusest, näiteks TUSTEP (Tuebingen System of Text Processing Programs, Stahl 2002) või Uplug (Tiedemann 2002).

Paralleelkonkordantside leidmist muudab efektiivsemaks tekstide eelnev lausetasandist detailsem paralleelistus, mis võimaldab vastava funktsiooniga varustatud tarkvaral otsitavaid sõnu või väljendeid ülejäänud tekstist esile tõsta või ka lasta kasutajal paralleelistus tekkinud vigu parandada (Gaussier jt 2000).

Konkordantside võrdlemisest saavad kasulikku infot nii keeleõppijad ja -õpetajad kui ka tõlkijad. Paralleelkorpust saab nõnda kasutada kakskeelse sõnaraamatu asemel või olemasolevate sõnastike täiendusena, eriti idiomatika ja keelelise stiili osas. Erialased sõnastikud toovad harva näiteid sõnade reaalsest kasutamisest lausetes, küll aga annab paralleelkorpus võimaluse uurida sõna võimalikke kontekste erinevates keeltes. Tõlkimistöös on oluline rõhutada tõlkemälukorpusest leitavate paralleelkonkordantside rolli terminite ühtlustamisel ja tõlke järjepidevuse tagamisel.

1.6. Vorming

Maailmas ja ka Eestis on olnud üheks levinud keelekorpuste formaadiks TEI (Text Encoding Initiative), kuid nii paralleelkorpuste (Godwin-Jones 2001) kui ka tavaliste korpuste rakendatavuse seisukohalt on soovitatud hakata kasutama TEI asemel XML-formaati (*Extensible Markup Language*). XML-formaadis tekstide teisendamiseks teistsugustele kujudele on välja töötatud mitmeid spetsiaalseid meetodeid (XSLT jm) ning korpuse meta-andmestik võib XML-i puhul paikneda tekstist eraldi XML-failides. Nüüdseks on olemas XML-iga ühilduv TEI-formaat (Sperberg-McQueen, Lou Burnard

2004), mis võimaldab TEI-kujul korpuste töötlemisel kasutada olemasolevaid XML-i jaoks loodud vahendeid.

1.7. Väikesed ja vähemuskeeled

Mitmed autorid (nt Scannell 2003) on pidanud tähtsaks paralleelkorpuste rolli keel-etehnoloogia arendamisel iseäranis väikeste ja vähemuskeelte jaoks. Seda võiks silmas pidada ka eesti keelele mõeldes. Väikeste keelte puhul peaks teine korpuse keel kuuluma suurte globaalsete keelte hulka, nõnda saab suures keeles olemas olevate rakenduste eeskujul lihtsamini luua rakendusi väiksemate keelte jaoks. Väikeste ja/või vähemuskeeltega paralleelkorpuse puudutavatest projektidest võiks nimetada järgmisi:

- EMILLE projekt (McEnery jt 2000) – Inglismaal elavate India päritoluga vähemusrühvuste keelte korpused ja paralleelkorpus.
- Inglisekeelse tesauruse põhjal iiri keele jaoks statistiliste meetoditega loodud teaurus (Scannell 2003). Scannelli hilisem kirjutis¹⁰ käsitleb Saint Louis' ülikoolis inglise-iiri paralleelkorpuse põhjal loodud tarkvara, mis seab normeeritud iiri kirjakeele sõnad vastavusse iiri keele normeerimisele eelnenud perioodi erinevate murrete vastavate sõnadega. Selline programm hõlbustab info-otsingut ja leksikograafilist tööd juhtudel, kui üks ja sama sõna võib keelesiseselt esineda mitmetel erinevatel kujudel.
- OPUS (Tiedemann, Nygaard 2004) – üle 30 miljoni sõna paralleeltekste 60 keeles (sh eesti keel), mis on saadud avatud lähtekoodiga tarkvara kasutajaliidestest ja mujalt
- STRAND (Resnik, Smith 2003) – internetist automaatselt leitud paralleeltekstide kogu, nende seas 59 dokumenti paralleelselt baski ja inglise keeles

¹⁰ <http://borel.slu.edu/pub/ccgb.pdf>

- PTOLEMAIOS (Kuhn 2004) – see projekt puudutab väiksemaid keeli, mille jaoks ei ole veel välja töötatud vajalikku grammatikat ja leksikoni morfoloogilise, süntaktilise ja semantilise analüüsi ning sünteesi jaoks. (Kuhn 2004) väidab, et analüüsivahendid saab automaatselt tuletada üksnes paralleelcorpusest, mis ei ole suures osas eelnevalt lingvistiliselt märgendatud.

2. Korpusepõhised masintõlke- ja tõlkeabisüsteemid.

Masintõlke (MT) idee on juba väga vana, kuid MT statistiliste meetodite uurimist alustati 1980. aastatel. Samal perioodil algas ka tõlkeabiprogrammide levik. Praeguseks on jõutud olukorrani, kus masintõlkealane uurimine ja arendustöö ongi põhiliselt keskendunud korpustest saadava andmestiku kasutamisele masintõlkes.

Korpusepõhises masintõlkes võib üldjoontes eristada kaht põhisuunda: statistilist masintõlget (Statistical Machine Translation) ja näitepõhist masintõlget (Example-Based Machine Translation).

Statistilises MT-s saab lause- ja sõnatasandil joondatud paralleelkorpusi kasutada tõlke originaaliläheduse hindamiseks. Paralleelkorpusest saab leida, mitu korda teatud sõna, fraas või struktuur kujutub üheks või teiseks võimalikuks tõlkeks, et hinnata lause tõlkevaste adekvaatsuse tõenäosust. Parima tõlke leidmiseks võib seejärel kasutada intellektitehnikast tuntud heuristilise otsingu algoritme. Lisaks originaalilähedusele on aktsepteeritava tõlke saamiseks vaja hinnata ka tõlke soravust, mida saab teha näiteks tõenäosusliku generatiivse grammatika abil. (Koit 2003)

Paralleelkorpusi on seni kasutatud põhiliselt lähedaste sugulaskeelte vahelises MT-s, kuna korpuse paralleelistamine on sugulaskeelte puhul hõlpsam. Mittesugulaskeelte puhul peab selleks, et neist masintõlkes kasu oleks, paralleeltekstid süntaktiliselt märgendama. St, vaja on paralleelseid nn puude pankasid (vt ptk 2.2).

2.1. Statistiline masintõlge

Statistiline masintõlge rajaneb nn müranivooga kanali (*noisy channel*) meetodil, mis võimaldab kasutada mitmeid informatsiooniteooriast, side, kommunikatsiooni, raadio, kõnetuvastuse jm valdkondadest pärit algoritme. Näiteks tõlkides inglise keelest eesti keelde tuleb leida selline eestikeelne lause, mille puhul tõenäosus, et see eestikeelne lause on inglisekeelse lause tõlge, on suurim. Selle tõenäosuse välja-arvutamiseks Bayesi valemi abil on tarvis teada ka tõenäosust, et mingi lause üldse eesti keeles võib esineda. Need tõenäosused saaks leida piisavalt suure paralleelkorpuse alusel. (Muischnek jt 2003:53-54)

IBM-i uurimiskeskuses valmis 1990. aastate algul statistikapõhine MT süsteem Candide inglise keelest prantsuse keelde tõlkimiseks. Statistilised andmed selle süsteemi jaoks saadi kakskeelsest Hansardi korpusest, mis sisaldab parlamendidebattide üleskirjutusi. Süsteemi edasiarendamise võimalusena nähakse lingvistiliste ja statistiliste meetodite kombineerimist: kasutada traditsioonilisi reeglipõhiseid meetodeid lausete morfoloogiliseks ja süntaktiliseks analüüsiks ja genereerimiseks ning statistilisi meetodeid sõnatähenduste ühestamiseks ja sõnavalikuks. (Koit 2003)

MT-süsteemis kasutatav paralleelkorpus ei pea alati olema lausete kujul täistekst, vaid võib olla näiteks ka ainult noomenifraase ja nende tõlkeid sisaldav paralleelkorpus. Philipp Koehn (2003) näitab, et kui integreerida tänapäevastesse statistilise masintõlke süsteemidesse eraldi noomenifraase tõlkiv allsüsteem, siis on lõpptulemus senistest saavutustest parem.

Hiljuti on katsed tehtud ka eesti-inglise statistilise masintõlkega (vt Fishel jt 2007). Seejuures kasutati samu paralleelkorpuseid, mille võrdlus on esitatud käesoleva töö viiendas peatükis – TÜ paralleelkorpus ja JRC-Acquis' paralleelkorpuse inglise-eesti osa. Fraasitabelite koostamiseks ja dekodeerimiseks kasutati Moses¹¹ – statistilist, keelest

¹¹ <http://www.statmt.org/moses/>

sõltumatut, fraasipõhist masintõlkesüsteemi. Sõnatasandi vastavused saadi GIZA++ (Och, Ney 2000) abil¹².

2.2. Puude pangad

Prahas asuvas Karli ülikoolis väljatöötatud tšehhi-inglise masintõlkesüsteem (Cmejrek jt 2003), kasutab lisaks inglisekeelsele korpusele ja elektroonilistele sõnastikele ka tšehhi-inglise paralleelset, lause lingvistilist tähendust esitavat puude panka (puude pangaks nimetatakse süntaktiliselt anoteeritud lausete kogumit, kus märgendus on käsitsi üle kontrollitud; Karli ülikooli projektis kasutati nn sõltuvuspuid – Dependency Based Machine Translation). See puude pank hõlmab inimestest tõlkijate poolt spetsiaalselt projekti tarbeks inglise keelest tšehhi keelde tõlgitud 11 000 lausepaari. Tõlkimisel püüti säilitada originaallausete struktuuri, et hiljem saada võimalikult häid automaattõlkeid.

Eesti keele osalusel on koostatud nn Sofie puude pank projekti Nordic Treebank Network raames. See korpus hõlmab Jostein Gaarderi romaani “Sofie maailm” kahe esimese peatüki tõlkeid kuude põhja-euroopa keelde, mis on süntaktiliselt märgendatud. Sofie puude panga märgendamisskeemiks valiti märgendamisskeem VISL¹³, mis kombineerib fraasistruktuuri- ja sõltuvuspuu head omadused. Puude kombineerimise standardiks valiti TIGER XML formaat, mis võimaldab paralleelpuudepanga ükskeelsete osade töötlemiseks kasutada redigeerimisvahendit *Annotate* ning päringu- ja visualiseerimisvahendit *TigerSearch* (Nivre jt 2005).

Puude pankade vajalikkust nii keeletehnoloogias kui ka lingvistilises uurimistöös on esile toonud mitmed autorid, vt näiteks (Volk, Samuelsson 2004) või (Abeillé 2003). Nii puude pangad kui ka paralleelkorpused on arvutuslingvistikas viimasel ajal saanud väga populaarseks uurimisaineks, ent nende kahe valdkonna kombinatsiooni, paralleelsete puude pankade vastu on huvi tuntud märksa vähem.

¹² Vt lähemalt GIZA++ kohta käesoleva töö ptk 4.4.3 ja samuti (Muischnek 2006)

¹³ <http://visl.sdu.dk/>

Siiski on paralleelpuudepankade tähtsus väga suur nii masintõlkesüsteemide arendamise kui paralleelistamisprogrammide töö parandamise, samuti komparatiivse lingvistika seisukohalt. Paralleelpuudepankasid saab tavalistest paralleelkorpustest paremini rakendada ka masintõlkeprogrammi töö tulemuse automaatseks hindamiseks, kuna lisaks paralleelvastete leidmisele saab tõlgete hindamisel kasutada ka puude pankadest saadavat morfosüntaktilist teavet. Seni on paralleelpuudepankasid loodud aga väga vähe ja olemasolevatel paralleelpuudepankadel on mitmeid puudusi – enamasti hõlmavad nad ainult spetsiifilise ainevaldkonna tekste ja pole fraasi- ning sõnatasandil joondatud. (Volk, Samuelsson 2004; Uchimoto jt 2004)

2.3. Tõlkemälusüsteemid

Tõlkemälusüsteemide osakaal professionaalsete tõlkijate poolt kasutatavate abivahendite hulgas on viimastel aastatel kasvanud, kuna infotehnoloogia arengud annavad tõlkemälu kasutamiseks järjest paremaid võimalusi. Traditsiooniline automaattõlge (kus väljundkeele lause produtseeritakse automaatselt keelereeglite ja leksikoni abil) aga õigustab end siiani vaid teatud üksikute juhtumite korral (sisend on kas eeltoimetatud, kujutab endast mingit piiratud ainevaldkonda või on väljund üksnes teksti sisu mõistmiseks vajalik toortõlge). Reeglipõhised MT-süsteemid ei oska olla loovad, arvestada piisavalt konteksti ega lahendada korrektselt semantilisi mitmesusi. Arvatakse, et ideaalses MT-süsteemis oleksid omavahel kombineeritud lingvistiline analüüs empiirilise andmestiku (korpused) kasutamisega.

Tõlkemäluprogrammides kasutamiseks on võimalik paralleelkorpusest lihtsal teel moodustada tõlkemälu, kuhu programmi kasutaja saab tõlketöö käigus omapoolseid täiendusi lisada. Arvatakse, et tõlkemälu tõlkimisel tarvitavad ja töö käigus täiendavad ettevõtted hakkavad tulevikus oma tõlketöö läbi tekkinud paralleelkorpust käsitlema kommertsiaalse produktina – potentsiaalseteks ostjateks oleksid sama tegevusvaldkonda

jagavad ettevõtted, kelle tõlkeprogrammides vastav tõlkemälu veel puudub (Bowker, Pearson 2002: 96).

Üks tuntumaid tõlkeabiprogrammide tootjaid SDL International pakub tarkvarapaketi Trados Freelance osana ka paralleeltekstide joondamisprogrammi Winalign, mis tekstide joondamisel arvestab nii tekstide struktuurilisi kui ka sisulisi omadusi ning võimaldab kasutajal joondamistulemusi korrigeerida. Samasuguseid paralleelistamisvahendeid on lansseerinud ka teised konkureerivad firmad.

Kommertseesmärkidel toodetud joondamisprogrammide tööprintsüübid sarnanevad vägagi arvutilingvistide poolt lingvistiliste uuringute tarbeks välja töötatud joondamisprogrammide tööpõhimõtetele, kuid kommertsjoondajad on kasutajasõbralikumad, sisaldades kasutajaliideseid automaatse paralleelistamise tulemuste käsitsi kinnitamiseks või parandamiseks (Bowker, Pearson 2002: 102).

Joondamisalgoritmide arendamise tulemusel loodetakse tõlkeabiprogramme tõhustada selles suunas, et oleks hõlpsam tõlgitavatele keeleüksustele ainult osaliselt sarnanevaid tõlkemäluvasteid tõlketöös kasutada (Gaussier jt 2000). Samuti nähakse tõlkeabiprogrammide tulevikku paralleelkonkordantside genereerimise protsessi kombineerimises terminoloogiliste andmebaaside kasutusega ja tõlkemälu lausete osade paralleelistamisega (Gaussier jt 2000).

Martin Volk (2005) on veendunud, et tuleviku tõlkeabiprogrammid põhinevad internetis leiduvatel keeleressurssidel. Üht võimalust internetist automaatselt paralleeltekstide leidmiseks on kirjeldanud Philip Resnik (1999). Nõnda leitud paralleeltekstide laused saaks pärast automaatset statistiliste meetoditega hindamist (kus saab võtta arvesse ka näiteks mõlemas keeles HTML-dokumentide struktuuri sarnasusi) lisada tõlkeprogrammi mällu.

Teksti tõlkevaste leidmiseks internetist saaks kasutada ka mõnd automaatset MT-süsteemi. On võimalik, et suured paralleelkorpused ei leiagi otseselt tuleviku tõlkeprogrammides rakendust ja selle asemel integreeritakse tõlkesüsteemidesse otsisüsteemid, mis otsivad internetist konkreetsele tõlgitavale tekstile sarnaseid tekste ning seejärel nende tekstide tõlkeid. Kui selline tõlkesüsteem leiab ühele lausele mitu erinevat tõlkevarianti, siis

võib statistilise MT eeskujul automaatselt hinnata tõlkevariantide kvaliteeti statistiliste keelemudelite abil. Kõigepealt otsiks see süsteem terviklausete täpseid ja nn hägusaid vasteid (*fuzzy matches*), jätkates osalause ja fraasidega. Kui suuremaid üksuseid ei õnnestu leida, siis võib tõlkesüsteem lause tõlkimiseks kasutada muid võimalusi (nt lingvistilisel analüüsil põhinevat sihtlause automaatset genereerimist) või abistada tõlkijat sõnatasandil, otsides sõnadele vasteid onlain-sõnastikest. (Volk 2005)

Internetist paralleeltekstide leidmisest keerulisemaks ja palju tähtsamaks küsimuseks osutub tegelikult kogutud ja lause- või fraasitasandil joondatud tekstide tõlkekvaliteedi automaatne hindamine. Tõlkeprogrammi mälu peaks hõlmama ainult häid ja adekvaatseid tõlkeid. (Volk 2005)

Järgnevalt iseloomustan lühidalt paralleelkorpustel põhinevaid näitepõhiseid tõlkesüsteeme, mis kujutavad endast tegelikult tõlkemäluprogrammide edasiarendusi.

2.4. Näitepõhine masintõlge

Näitepõhine masintõlge ehk analoogtõlge hõlmab tihti ka reeglipäraseid ning statistilisi masintõlkemeetodeid. Näitepõhise hübriidsüsteemi ja tõlkemäluprogrammi tööprintsip ei erine siiski üksteisest seni, kuni tõlkesüsteemi poolt kasutatavas paralleelkorpuses leidub enam-vähem täpne vaste tõlgitavale lausele ning selle lause tõlge teise keelde. Kui sobivat lauset ei leita, siis otsib näitepõhine süsteem tõlke genereerimiseks väiksemaid lauseüksusi, mida töödeldakse statistiliste ja/või lingvistiliste reeglite alusel, aga primaarseks jäävad lauseosade tõlkevastete leidmisel endiselt tõlkenäited. (Carl, Way 2003: xix)

Mõned olemasolevad näitepõhised süsteemid kasutavad paralleelkorpusi vahetult tõlkeprotsessi käigus, kuid tõlkimiseks vajalik teave võidakse tõlkesüsteemi poolt korpusest ammutada ka eelnevalt spetsiaalse õppemooduli poolt. Näitepõhise tõlkesüsteemi poolt nõnda automaatselt genereeritud tõlkemallid sarnanevad reeglipõhises masintõlkes

kasutatavate lingvistide poolt loodud ülekandereeglitega, kuid on siiski põhimõtteliselt nendest erinevad (Carl, Way 2003: xx).

Tõlkeprotsessi käigus paralleelsete tõlkenäidete poole pöörduvad süsteemid on üldiselt seda efektiivsemad ja parema tõlkekvaliteediga, mida suurem on tõlkenäidete korpus; mida suurem on võimalus, et tõlke saamiseks ei ole vaja kasutada grammatilisi kirjeldusi ja keerulisi reegleid, seda suurema eelise saavutavad vahetult tõlkenäiteid kasutavad süsteemid ülejäänud tõlkesüsteemide ees (Carl, Way 2003: xxvi). Loomulikult aitaks tõlke kvaliteedile kaasa kasutatava korpuse sõnatasandil paralleelistatus.

Eiichihiro Sumita (2003) kavandatud näitepõhine MT-süsteem kasutab lausetasandil joondatud kakskeelset korpust, kakskeelset tõlkesõnastikku ja mõlema keele tesaurust. Tesauruse abil mõõdetakse kõigepealt sisendlause ja paralleelkorpuse lähtekeele lausete sõnade semantilist kaugust, mille alusel leitakse sisendlausele semantiliselt kõige lähedasemad näitelauseid. Kõige lähedasema näitelause ja selle teise keele vaste kattuva osa vahel luuakse jooksvalt tõlkemall. Kui genereeritakse mitu tõlkemalli (kuna mitu tõlkenäidet olid sisendist samal semantilisel kaugusel), siis valitakse kas kõige sagedasem tõlkemall või kasutatakse tõlkemalli valimiseks sõnasageduste summeerimist või juhuslikku valikuprintsiipi. Sisendlause sõnad, mida tõlkemall ei hõlma, tõlgitakse automaatselt sõnastiku abil ja nõnda saadakse väljundlause teises keeles.

Selline küllalt lihtne tõlkesüsteem on ebaefektiivne pikemate lausete korral, kuna sarnase lause leidmise tõenäosus näitelauseste korpusest on väiksem. Katsetades seda meetodit tõlkimaks jaapani keelest inglise keelde, saadi ometi 80 % korral lausetest rahuldav tulemus. Ebaefektiivsust pikkade lausete korral saab vähendada, kui lasta enne tõlkeprotsessi sisend- ja näitelauseid automaatselt osadeks jagada. (Sumita 2003)

3. Paralleelkorpuste kasutamine lingvistilistes uurimustes ja keeleõppes

3.1. Lingvistilised ja tõlketeoreetilised uuringud

Paralleelkorpused kujutavad endast kasulikku uurimisainest kontrastiivsele lingvistikale ning tõlketeooriale, pakkudes neile distsipliinidele omavahelisi kokkupuutepunkte ja võimaldades lingvistikal ja tõlketeoorial seeläbi senisest ulatuslikumalt üksteist täiendada.

Paralleelkorpuste abil on võimalik allutada tõlkijate loomingulised lähenemised erinevate keelekonstruktsioonide tõlkimisel statistilisele analüüsile, mille tulemused võivad pakkuda ka lingvistidele huvitavat infot vastavate keelte ülesehituse kohta. Seejuures tuleb esmalt piiritleda sagedasemad nendest juhtudest, mil ühe keele lause või lauseosa on teises keeles edasi antud võrreldes lähtekeelega erineval viisil. Sellised juhtumid võivad seisneda tüüpilistes väljajäätudes, osalauseste alistustüübi muutustes jne.

Seejärel tuleb määratleda teatav *interlingua*, milles esitatakse tõlgitavate keeleüksuste või mõistete sarnasused ja erinevused vastavas keelepaaris, kusjuures tõlkeüksuste maht võib olenevalt keeltest ja uurimiseesmärkidest suuresti varieeruda (Siin kasutatav *interlingua* mõiste on analoogne reeglipõhises masintõlkes mitmete keelepaaride vahelist tõlkimist hõlbustama pidava tehisliku üleminekukeelele.) Mitte-ekvivalentsete tõlkeüksuste kaardistamisel võime lähtuda mitmetest olemasolevatest teoreetilistest raamistikest. Teine võimalus on koondada tõlkeüksuste võrdlusandmed kontrastiivsesse mitmekeelsele *wordnet*ile sarnanevasse andmebaasi. (Salkie 2002)

Keele ja tõlkimisprotsessi olemuse valgustamiseks sobivad üksikutest keelepaaridest paremini mitmeid erinevaid keeli hõlmavad paralleel- või võrreldavad korpused. Mitme-

keelsete korpuste puhul saab huvitavaid järeldusi teha näiteks ka erinevate lähtekeeltega tõlgete võrdlemisest ühe ja sama keele siseselt.

Samas võime ka ainult ükskeelset korpust uurides saada konkreetseid vastused mitmetele tõlkimisega (nagu ka keeleõpetamisega) seotud küsimustele – näiteks samakeelsete tõlgitud ja originaaltekstide võrdlemisel selgub, mis määral tõlkijad suudavad vältida võõrapäraste lausekonstruktsioonide ja sõnade põhjendamatu ülevõtmist tõlgitavast keelest.

Selliste uuringute tulemuste kinnitamiseks on aga vajalik uurimisse kaasata ka vastavad tekstid selles keeles, millest tõlgiti, ning soovitatavalt ka teisi keeli, millesse samad tekstid veel tõlgitud on – nõnda saame teada, kas originaali ja tõlke süntaktilised või leksikaalsed erinevused on mingis osas võrreldes teiste keeltega märkimisväärsed. (Johansson 2002)

Nii näitis Martin Gellerstam (1996: 59) rootsi verbi *tillbringa* ülemäärast tarvitust rootsikeelsetes tekstides, mis olid inglise keelest tõlgitud, ja oletas, et tegu on inglise verbi *spend* liiga sagedase otsetõlkega (sõna *spend* tähendussfäär inglise keeles on laiem kui *tillbringa* oma rootsi keeles).

Stig Johansson (2002) kontrollis seda oletust inglise-norra paralleelkorpuse peal, kuhu on lisatud inglisekeelsete tekstide tõlkeid saksa, hollandi ja portugali keeles. Johansson leidis, et ühelt poolt inglise ja teiselt poolt norra ja saksa keel kujutavad aja möödumist erinevalt: inglise keel tõlgendab aja möödumist üldjuhul kui aja kulutamist (*consume*) või möödasaatmist, norra ja saksa keeles eelistatakse aga rääkida sündmuste kestmisest (*duration of an event*). Tõlkimisel kaldutakse aga tõepoolest kasutama inglise keelega analoogseid keelekonstruktsioone, mida tähistab ka *tillbringa* kasutamine tõlgetes. (Johansson 2002)

3.2. Paralleelkorpused keeleõppes

Nagu juba mainitud, paralleelkorpust saab kasutada kakskeelse sõnaraamatu asemel või olemasolevate sõnastike täiendusena, kui sealt tekstiotsingu abil esile tuua vajalikud konkordantsid. Tegelikult on sõnaseletuste ja grammatilise infoga varustatud paralleel-

tekstid aga üks traditsioonilisi keeleõppe vorme näiteks antiikkeelte puhul, arvuti abi paralleelkorpustest keeleõppimiseks vajaliku info ammutamisel võib pidada selle traditsioonilise keeleõppe vormi moodsaks täienduseks. Et traditsiooniline paralleeltekstidest õppimine on populaarne ka tänapäeval, sellest annavad tunnistust mitmete nimekate kirjastajate (Penguin, Harvard, Random House, Reclam, Mercier jt) poolt välja antavad paralleeltekstid (Nerbonne 2000: 3). Paralleeltekste on seostatud nn “tõlkemeetodiga” keeleõppes ning sellel meetodil on olnud nii apoloogeete kui ka kritiseerijaid; tõlkemeetodit on kritiseeritud “kommunikatiivse pädevuse” arendamise vajalikkuse seisukohast (Widdowson 1990: 117ff).

Ometi ei saa eitada, et keele (mitte ainult kõnekeele) võimalikult täielikuks omandamiseks ning tõlkijate koolitamisel pakuvad paralleeltekstid emakeelena teist keelt kõnelevale keeleõppijale asendamatu abimaterjali. Paralleeltekstid võimaldavad nii koos korpuseanalüüsitarkvaraga kui ka ilma selleta keeleõppijale ligipääsu autentsele keelekasutusele koos võõrkeelsete lausete tõlgetega emakeelde ja see võimalus on sageli olnud õppimist hoogustavaks teguriks. Kui aga tegemist on juba elektroonilisel kujul paralleelkorpusega, siis osutub kvaliteetse korpuse kasutamiseks mõeldud tarkvara olemasolu korpuse abil õppimisel hädavajalikuks – ilma spetsiaalse arvutiprogrammiga tabab selliseid õppeprojekte tõenäoliselt ebaedu (Nerbonne 2000: 6).

Väidetavalt (Nerbonne 2000) osutuvad paralleelkorpused keeleõppijale kõige kasulikumaks leksikaalsete sõltuvuste võrdlemise osas ning sõnade või väljendite erikeelsete konkordantside esiletoomine peaks kujutama endast soovitatavalt üht osa multifunktsionaalsest keeleõppetarkvarast (nagu GLOSSER – rahvusvaheline projekt, milles oli osaline ka Eesti). Korpuse lekseemid – mitte sõnavormid – peaksid olema indekseeritud, et otsinguga leitaks ka sama sõna teistsuguse tüvekujuga vormid. Korpus peaks olema lausetasandil paralleelistatud. (Nerbonne 2000)

4. Leksikoni loomine paralleelcorpuse baasil

4.1. Automaatselt loodavate leksikonide vajalikkus

Tuleviku seisukohalt on uute lähenemiste leidmine kas korpustest või muudest allikatest arvutis loetavate sõnastike loomiseks keeletehnoloogias üks kahest põhisuunast, millele kogu maailmas tähelepanu pööratakse (Muischnek jt 2003). Paralleelcorpustest automaatselt genereeritud sõnastikud leiavad rakendust nii inim- kui masintõlkes, keeleõppes ja ka mujal, näiteks sellise spetsiifilise arvutitehnoloogilise ülesande hõlbustajana, nagu seda on semantiline ühestamine (Ide jt 2002).

Esiolgu saab rääkida ainult poolautomaatselt mitmekeelsete erialasõnastike koostamisest paralleelcorpuste baasil, sest ilma inimkorrektori parandusteta saab sellisel viisil koostada üksnes leksikograafidele, terminoloogidele või tõlkesüsteemidele abiks olevat toormaterjali. Erialased tekstid sobivad leksikoni ekstraheerimiseks paremini näiteks kirjanduslikest tekstidest seetõttu, et erialased terminid tähistavad enamasti kindlapiirilisi mõisteid, millele leidub kindel ja järjepidev vaste paralleelteksti sõnade hulgas (Bowker, Pearson 2002: 171–174, 220; Fung 2000). Tõlke järjepidevusest sõltub suuresti ka leksikoni genereeriva programmi võime luua korrektseid seoseid sõnade vahel.

Isegi kui keele erialase valdkonna jaoks on juba olemas kakskeelne sõnastik, siis on üldjuhul paralleelcorpuse põhjal võimalik automaatselt genereerida olemasolevast tunduvat mahukam sõnastik. Tänapäeva uut tüüpi corpusepõhistes õppijasõnastikes esitatakse lisaks tavapärasele leksikaalsele ja grammatilisele infole ka teavet tähenduste piirangute, kollokatsioonide, grammatiliste mallide, stiili, registri ja kasutussageduse kohta ning luuakse seosed sõna struktuuri, kasutamise ja tähenduse vahel (Kitsnik 2006: 96). Corpusest automaatselt ekstraheeritud sõnastikud on üheks sellise teabe allikaks, võimaldades

muuhulgas ka parandada ja täiendada olemasolevaid sõnastikke nii uute kirjete osas kui ka näiteks muuta sõnatähenduste hierarhiat genereeritud sõnastikus sisalduva korpusest ammutatud sagedusinfo põhjal.

Olgugi et kaasaegne keeletehnoloogia võimaldab automaatselt genereerida vaid üsna suure veaprotsendiga sõnastikke isegi erialakeele puhul, on siiski tekstiressursside olemasolul mahukate ja osaliselt vigaste sõnastike loomine enamasti vajalik mitmetel põhjustel.

Näitepõhise tõlkesüsteemi jaoks on genereeritava leksikoni juures veaprotsendist olulisem kirjete arv, kuna tõlgitava üksuse vastetest korpuses peab tõlke õnnestumiseks vaid ühel olema korrektne paralleelistus (Brown 1997: 6). Suurem maht, mis tähendab saagise eelistamist täpsusele¹⁴, on ka terminoloogi, tõlkija või sõnaraamatu koostaja seisukohalt parem lahendus. Lihtsam, kui otsida tekstist ise programmi poolt leidmata jäänud tõlkevasteid, on programmi poolt välja pakutud valesid märksõnakandidaate eraldada korrektsetest tõlgetest või kustutada.

4.2. Paralleelkorpuste ja võrdluskorpuste põhjal leksikonide genereerimise meetodid

Leksikoni genereerimine seisneb paralleelkorpuses üksteise tõlgeteks olevate sõnade või väljendite tuvastamises ja ekstraheerimises. Seda saab aga teha mitmel erineval moel.

Erinevad lähenemised tõlkevastele automaatsele ekstraheerimisele jagunevad kahte põhikategooriasse: nn „hüpoteesi kontrollimise” ehk heuristilised meetodid (nt Smadja jt 1996) ja estimateerivad meetodid (nt Hiemstra 1997).

¹⁴ Saagise all on siin mõeldud leksikoni jaoks ekstraheeritud korrektsete seoste arvu suhet kõigi paralleelkorpuses sisalduvate ja omavahel tõlkelises seoses olevate sõnade või väljendite arvuga. Paralleelistuse saagise mõõtmiseks kasutatakse ka nn kuldstandardeid. Sellisel juhul kujutab saagis endast korrektsete seoste arvu suhet seoste arvuga korrektselt paralleelostatud alamkorpuses. Täpsus on korrektsete tõlkevastavuste hulk leksikonis, võrrelduna kogu leksikoni mahuga.

„**Hüpoteesi kontrollimine**” tähendab tõlkevastete kandidaatide loendi genereerimist. Kandidaadid allutatakse statistilisele analüüsile, mis peab näitama, kas tegemist on tegelike tõlkevastetega või mitte.

IBM-i teadlaste ideedest (Brown jt 1990)¹⁵ lähtuvast statistilise masintõlke paradigmast inspireeritud **estimeerivad meetodid** põhinevad tõenäosusliku bitekst-mudeli loomisel, mis võimaldab tõlkevasteid hinnata mitte ainult eraldi, vaid ka rühmadesse jaotatuna.

Mõlemal lähenemisel on oma plussid ja miinused. Allpool vaatluse alla tulevatest leksikoni genereerimise vahenditest võib PWA-d pidada heuristiliste meetodite ja Giza++ estimeerivate meetodite paradigmasse kuuluvaks. (Tufiş, Barbu 2001: 156, Tiedemann 2003)

Varasemad lähenemised kakskeelse leksikoni ekstraheerimisele (K-vec algoritm, DK-vec algoritm jmt) põhinesid nn „ankurpunktidest” lähtuvatel lausetasandil paralleelistamise meetoditel, mida kombineeriti leksikaalse koosinemuse analüüsiga (nt Church jt 1991). Heuristilised meetodid sõnade joondamisel lähtuvadki üldiselt ideest, et tuleb leida sõnapaar, mis esineb koos märgatavalt sagedamini kui seda võiks lubada juhus (*t-score*, Dice'i koefitsient jm).

Lisaks koosinemuse analüüsile on üldkasutatav ka sõnesarnasuse mõõtmine leidmaks omavahel etümoloogiliselt seotud sõnu. Sõnesarnasuse mõõdupuu näiteks on LCSR (*longest common sub-sequence ratio*), mis tähendab kahe sõne pikima ühise tähejärgnevuse ja sõnepaari pikema sõne pikkuse suhet. Sõnesarnasuse mõõtmine eeldab, et mõlemal keelel on sarnane tähestik ja et etümoloogilises suguluses olevate sõnade kirja piltide vahel eksisteerib arvestatav sarnasus. Tõlkepaaride ekstraheerimiseks sobivaid sõnesarnasuse mõõtmise meetodeid on võrrelnud Lars Borin (1998).

Sõnade joondamise abivahendina saab kasutada ka olemasolevaid elektroonilisi kakskeelseid sõnastikke. Seejuures sõltub selliste lisaressursside mõju ulatus sõnaparalleelistuse tulemustele suuresti konkreetse sõnastiku temaatilisest sobivusest paralleelcorpuse temaatikaga ning samuti sõnastiku mahust. Mõnede tekstide puhul võib joondamisel abi

¹⁵ IBM-i mudelite kohta vt lähemalt ptk 4.4.3.

olla ka tekstide vorminduse võrdlemisest. Lisaks saab korpusevälise ressursina kasutada ka mitmesugust keelespetsiifilist informatsiooni sõnajärje reeglite, süntaktiliste suhete jms kohta, kuid selliste andmete kombineerimises statistiliste meetoditega ei ole esialgu märkimisväärset edu saavutatud. (Tiedemann 2003)

Kuna võrdluskorpusi on lihtsam koostada kui paralleelkorpusi, siis on püütud välja töötada ka abivahendeid sõnavastavuste leidmiseks joondamata võrdluskorpustes (nt Fung 2000).

Kui on tegemist sugulaskeeltega, siis saab sõnadevahelise vastavuse osaliselt tuletada sõnatüvede sarnasusest, mittesugulaskeelte puhul tuleb vastavuste leidmiseks kasutada statistilisi meetodeid ja/või infot sõnade konteksti ja paiknemise ning nende süntaktilise rolli kohta.

Kasutades info-otsingule sarnaseid meetodeid, võib võrdluskorpuste abiga leida vasted paralleelkorpuste põhjal või traditsioonilisel meetodil loodud leksikonides puuduvatele sõnadele. Nõnda võib osutada võimalikuks näiteks olemasolevate elektrooniliste sõnastike pidev automaatne täiendamine uute, varem mitte eksisteerinud sõnadega või tuletistega. (Bowker, Pearson 2002: 171–174, 220; Fung 2000).

Tüüpiliselt seisneb sõnavastavuste leidmine võrdluskorpustest lähte- ja sihtkeele kõigi sõnade sagedasemate kollokatsioonide põhjal moodustatud kontekstvektorite võrldemises, kusjuures kontekstvektorite tõlkimiseks kasutatakse olemasolevaid kakskeelseid sõnastikke (Déjean jt 2002: 1). Paremate tulemuste saavutamiseks on kontekstvektor-tõlkimist vaja kombineerida keeleliste lisaressursside (mitmekeelsed tesaurused vms) abi kasutavate meetoditega.

4.3. Grammatilise analüüsi osa sõnastike genereerimisel

Eric Gaussier jt (2000: 254–255) toovad välja kolm faasi, milles sõnade või sõnaühendite tõlkepaaride ekstraheerimine paralleelkorpusest üldjuhul seisneb: tõlgitavate üksuste tuvastamine ning filtreerimine igas keeles eraldi, millele järgneb seoste leidmine tuvastatud

üksuste vahel statistiliste algoritmide abil ning lõpuks leksikoni genereerimine omavahel seostatud sõnade loendi põhjal. Vajalike üksuste tuvastamine võib toimuda ka dünaamiliselt joondamisprotsessi käigus.

Kui kasutada sõnaühendite tuvastamiseks lingvistilisi meetodeid, siis on siinjuures eelduseks, et kõigi keelte (mõlema keele) lausete morfosüntaktiline analüüs toimub piisavalt heal tasemel ja sarnasel moel, mis alati ei ole võimalik. Samuti ei ole väga lihtne kindlaks teha esimeses faasis leitud grammatiliste mallide keeltevahelisi vastavusi, eriti juhul, kui soovitakse leksikoniga hõlmata ka mitmesõnalisi sõnaühendeid.

Viimast probleemi on püütud lahendada erinevatel viisidel. Üheks väljapääsuks on igasuguste grammatiliste korrelatsioonide leidmine keelepaaride vahel enne terminitele või soovitatavatele üksustele vastavate grammatiliste mallide tuvastamist. Teine võimalus on kõigepealt leida grammatilise analüüsi abil kõik üksused ainult esimeses keeles olevates tekstides ja seejärel leida nende üksuste tõlked teises keeles. Kui aga eesmärgiks on luua kakskeelne ja kahe-suunaline oskussõnastik, siis tuleks esmalt tuvastada terminid mõlemas keeles ja seejärel need terminid paralleelistada.

Grammatilise analüüsi abil leitud tõlgitavate keeleüksuste piire saab täpsustada ka paralleeltekstide põhjal loodud statistiliste tõlkemudelite abil – nõnda osutub reaalseks näiteks ka kõigi püsiühendite ekstraheerimine paralleeltekstidest (Melamed 2001).

Võimalikud on ka keerulisemad algoritmid, mille korral sõnade või sõnaühendite grammatilise struktuuri vaatlemise abil parandataks joondamisel tehtud vigu, ning vastupidi, joondamismalle hinnates täpsustataks vajalike üksuste tuvastamist (Gaussier jt 2000). Allpool pöoran aga põhitähelepanu grammatilist analüüsi minimaalselt kasutatavatele rakendustele, mille puhul nii sõnade või sõnaühendite tuvastamine kui ka paralleelistamine toimub suuremas osas statistiliste meetoditega.

Sõnastiku automaatsel genereerimisel paralleelkorpusdest tuleb vahet teha eeldefineeritud terminite ekstraheerimise ja laiema tähendusega sõnastikugenerereerimise vahel (mida võidakse siiski rakendada ka erialaste tekstide põhjal). Arvutilingvistikas määratletakse tehnilisi termineid sageli kitsalt teatud kindlatele morfosüntaktilistele tunnustele vastavate

noomenifraasidena (Blank 2000: 240). Selliste kindlate omadustega noomenifraaside automaatne tuvastamine nõuab korpuselt kindlasti süntaktilist märgendust ja leksikoni genereerimine peab sellisel juhul toetuma konkreetse keele grammatilisele kirjeldusele. Seega moodustab keelest sõltumatute meetodite puhul selline kitsatähenduslik terminoloogia vaid osa korпустest genereeritavast leksikonist, kuna ekstraheeritavatele üksustele ei saa morfosüntaktilise info puudumisel kehtestada sõnaliigist vms lähtuvaid piiranguid.

Esiialgu puudub eesti keele jaoks tarkvara, mis teostaks paralleelistamist ja sõnastiku genereerimist lähtudes eesti keele grammatilisest struktuurist ja selle struktuuri vastavusest mõne teise keele struktuuriga. Keelest sõltumatud meetodid tähendavad muuhulgas ka seda, et esialgsed tekstist ekstraheeritud üksused, mille omavahelise joondamise kaudu saadakse lõpuks sõnastik, peavad olema leitud põhiliselt statistiliste meetoditega. Mitmeid uurimusi, millest saab lähtuda kakskeelse leksikoni genereerimisel, on aga tehtud ka ükskeelsest korpusel statistiliste vahenditega mitmesõnaliste üksuste tuvastamiseks (Tiedemann 2003).

Kui lähtuda leksikoni võetavate üksuste joondamisel keele grammatikast, siis tekib järgmine probleem: kuigi sagedasti esinevad sõnühendite grammatiliste mallide vastavused erinevate keelte vahel on tuvastatavad, esineb siiski suhteliselt palju mitmesusi ja kõrvalekaldeid reeglitest.

4.4. Sõnastiku genereerimise praktilised võimalused

4.4.1. Eeldused ja eeltöötlus

Mitmed kommertstarkvara tootjad (Xerox, Ahead Software, SensoLogic, SDL International) pakuvad mõne oma tarkvarapaketi osana ka terminisõnastiku automaatse ekstraheerimise võimalust etteantud tekstide põhjal, kuid nendel programmidel puudub esialgu eesti keele tugi ning kommertstarkvara poolt kasutatavad meetodid ei võimalda töö kohaldamist teiste keelte tarbeks. Seetõttu on esimeseks loogiliseks sammuks teel eestikeelse osalusega sõnastiku automaatse genereerimise poole keelest sõltumatute

meetodite katsetamine vabavara abil või eesti keelt toetava leksikonigenererimistarkvara loomine.

Mida tuleks silmas pidada, kui on kavas välja töötada eriotstarbelist leksikograafiatööle orienteeritud tarkvara, mis genereerib paralleelcorpuse põhjal muuhulgas loodavasse sõnastikku sobivaid sõna- või väljendipaare?

Enamik leksikoni ekstraheerimise alastest töödest põhineb kindlatel eeldustel, mida tuleks arvesse võtta eriti juhul, kui tarkvara loomisel soovitakse alustada lihtsamatest algoritmidest. Ükski nendest postulaatidest ei kehti tegelikkuses sajabrotsendiliselt, kuid erandid ei põhjusta nii suurt langust tulemuste kvaliteedis, et eelduste rakendamine poleks põhjendatud. Need eeldused on:

a) mitmetähenduslikku leksikaalset üksust kasutatakse ühe ja sama teksti siseselt ainult ühes kindlas tähenduses;

b) tõlkeüksuste paari kuuluvate sõnade sõnaliigid peavad omavahel sobima, st näiteks verbile ühes keeles võib vastata üksnes verb või mõni teine sõnaliik, mis on tunnustatud võimeliseks täitma tõlkes verbi funktsiooni – ka see reegel ei ole muidugi tegelikkuses absoluutne;

c) tõlkevastete kandidaatide seas on tõenäolisemad tõlkevasted need, millesse kuuluvate sõnade suhteline asend lauses on üksteisele lähedasem;

d) tõlkepaari ühe poole leksikaalsele üksusele vastab maksimaalselt üks leksikaalset üksust tõlkepaari teisel poolel.

Esimene eeldus vastab paraku seda vähem tõele, mida vähemkasutatava sõna või sõnaühendiga on tegu ja ka suur hulk oskussõnu tõlgitakse teise keelde mitmel erineval moel isegi sama teksti siseselt. Ingeborg Blank (2000) leidis prantsuse-saksa näidiskorpuse varal tehtud katse abil, et 5–15 % terminitest on sellised, millel on teises keeles rohkem kui üks tõlkevaste (Blank 2000: 246–247).

Ühe termini tõlkevasted võivad olla aga erineva grammatilise struktuuriga. Nii näiteks esineb sõna *sihthiikmesriik* vastena TÜ inglise-eesti paralleelcorpuses kahel korral *the Member State of destination*, aga ühel korral ka *the destination Member State*; sõna

lähetuskoht esineb inglise keeles kaheksal korral kui *place of dispatch*, kahel korral kui *place of destination* ning kahel korral vastab noomenifraasile hoopis umbisikuline verbivorm: *(the products) are dispatched from*. Selliste mittevastavuste võimalikkust tuleb arvesse võtta nii terminite automaatsel piiritlemisel kui ka joondamisfaasis.

Ka viimast, üks-ühele vastavuse eeldust on leksikonide genereerimisel küll laialdaselt kasutatud, kuid ka see eeldus tekitab siiski suhteliselt palju ebakorrektsed tõlkeid, kui ühe keele liitsõnale vastab teises keeles mitmesõnaline väljend. Niisiis on see nn „1:1-kaardistuse hüpotees” inglise-eesti keelepaari korral küsitav: eesti keele liitsõnale vastab inglise (või ka näiteks prantsuse) keeles tavaliselt mitmesõnaline üksus. I. Blank (2000: 247) toob näite selle kohta, kuidas saksa keele liitsõnaline termin (*Einspruchsbeschwerdeverfahren*) esineb vastavates prantsusekeelsetes tekstides järjepidevalt kujul *procédure de recours engagée à l’encontre d’une décision rendue sur opposition* (kompleksne noomenifraas).

Seda probleemi on täheldanud Pim van der Eijk (1993), kelle uurimus põhineb inglise-hollandi korpusel, Lars Ahrenberg jt (1998), kes tegelesid inglise-rootsi korpusel ja samuti I. Blank (2000), tuginedes saksa-inglise-prantsuse korpusel, jt. Keelespetsiifilise eel- (automaatne segmenteerimine) ja järeltöötuse (osaliste tõlgete filtreerimine) abil on 1:1-kaardistusest tulenevad probleemid vähemalt osaliselt ületatavad (Tufiş, Barbu 2001: 157). Parema tulemuse saamiseks tuleks siiski vähemalt eesti keele puhul üks-ühele paralleelsetusega ühendada mitmesõnaliste üksuste tuvastamine ja joondamine.

Paralleelkorpuse eeltöötlus võiks lisaks leksikaalsete üksuste segmenteerimisele eesti keele puhul hõlmata ka Kadri Muischneki (2006) poolt käsitletud inglise-eesti masintõlke kvaliteedi parandamiseks sobivaid meetodeid – ühendverbide restruktureerimist ning liitsõnade osadeks jaotamist. Selline eeltöötlus eeldab aga omakorda morfoloogilist analüüsi, mis tähendab küll eemaldumist esialgsest keeltevahelise portatiivsuse printsiibist, kuid võimaldab tõlkevasteid ekstraheerida tunduvalt lihtsamalt.

4.4.2. Lihtsa leksikograafilise abivahendi kavand

Morfoloogiline analüüs lubaks tarkvara kavandamisel lähtuda kõigepealt ainult ühest sõnaliigist, näiteks verbidest. Sellisel juhul peaks programm esmalt morfoloogilise analüsaatori abil kõik valitud sõnaliigi esindajad sisendkorpuses tuvastama ja lemmatiseerima ja korpuse paralleelistama lausetasandil. Paralleelistamiseks võib kasutada mõnd Gale'i ja Churchi algoritmi modifikatsiooni (nt Davis jt 1995), mille puhul on võimalik enne joendamist laused filtreerida, teostades joendamise ainult punktuatsiooni, pärisnimede vms põhjal.

Seejärel võiks programm luua nimekirja kõigist võimalikest valitud sõnaliiki kuuluvatest keeltevahelistest sõnapaaridest, mis ei ületa oma joondamisüksuse piire. Selline loend võib endast juba kujutada arvestatavat abimaterjali leksikograafide, kuid tulemuse parandamiseks peaks programm lisaks võrdlema kõigi korpuse valitud sõnaliiki kuuluvate keeltevaheliste sõnapaaride elementide omavahelist ortograafilist sarnasust ja koosinemise tõenäosust (väljendatuna näiteks seosetugevuse üldtuntud mõõdu, Dice'i koefitsiendi või mõne selle variandina). Kuna tegemist on sõnapaaridega, siis saab seejuures rakendada kollokatsioonide analüüsimisel kasutatavaid meetodeid (Brew, McKelvie 1996: 48).

Igale sõnapaarile omistaks programm nende näitajate alusel arvulise märgendi, mis iseloomustab tõenäosust, et tegemist on vastastikuste tõlgetega. Kasutajale esitataks edasiseks töötluks üksnes need sõnapaarid, millega vastavusse seatud arv ületab kasutaja seatud lävendi. Nõnda saab tekstist automaatselt esile tuua väidetavalt 30 % kõigist tekstis leiduvatest korrektsetest tõlgetest, täpsusega 90 % (Brew, McKelvie 1996: 51). Kui tõenäosuse arvutamisel lähtuda üksnes ortograafilisest sarnasusest, siis on võimalik esile tuua ka potentsiaalsed eksitavad valepaarid (nn *faux amis*) – kirjapildilt üksteisele sarnanevad, kuid tähenduselt erinevad sõnad. (Brew, McKelvie 1996)

4.4.3. Poolautomaatsed vahendid

Kui tarkvara väljatöötamine ei ole mingil põhjusel võimalik või otstarbekas, siis on sõnastiku genereerimiseks võimalik kasutada ka vabavara.¹⁶ Kõige lihtsam viis sõnade käsitsi joondamiseks või ka paralleelteksti põhjal sõnastiku toormaterjali loomiseks on kasutada selleks mõnd graafilist joondamisvahendit, mis võimaldab biteksti üksteisega vastavuses olevad sõnad hõlpsasti omavahel ühendada ja genereerida automaatselt omavahel seostatud sõnade loendi.

Sellise joondamisvahendi näide on Rebecca Hwa ja Nitin Madnani poolt välja töötatud Java-programm¹⁷, mis genereerib loendi kasutaja poolt arvutihiire abil seostatud sõnadest. Katse selle programmi abil paralleelistada sõnatasandil TÜ inglise-eesti paralleelkorpusest juhuslikult valitud lause andis järgneva tulemuse (1). Toon kõigepealt ära kasutatud lause inglise ja eesti keeles ja seejärel selle lause paralleelistusele vastava väljundi.

(1)

a) The representatives agree on the desirability of acceptance of the WTO Agreement by all participants in the Uruguay Round of Multilateral Trade Negotiations (hereinafter referred to as "participants") with a view to its entry into force by 1 January 1995, or as early as possible thereafter.

b) Esindajad lepivad kokku, et on soovitatav, et kõik mitmepoolsete kaubanduslääbirääkimiste Uruguay voorus osalejad (edaspidi "osalejad") tunnustaksid WTO lepingut, et jõustada see 1. jaanuaril 1995 või niipea kui võimalik pärast seda.

1 1 (The, Esindajad)

2 1 (representatives, Esindajad)

¹⁶ Vabavarana saadaval olevatest sõnatasandil paralleelistamist võimaldavatest tarkvarapakettidest ülevaate saamiseks vt <http://www.cse.unt.edu/~rada/wa/#softwareWA> (21.08.2006)

¹⁷ Vt <http://www.umiacs.umd.edu/~nmadnani/alignment/forclip.htm> (21.08.2006)

3 2 (*agree, lepivad*)
3 3 (*agree, kokku,*)
5 6 (*the, soovitav,*)
6 6 (*desirability, soovitav,*)
8 16 (*acceptance, tunnustaksid*)
18 11 (*Uruguay, Uruguay*)
16 12 (*in, voorus*)
17 12 (*the, voorus*)
19 12 (*Round, voorus*)
15 13 (*participants, osalejad*)
14 8 (*all, kõik*)
7 7 (*of, et*)
11 17 (*WTO, WTO*)
10 18 (*the, lepingut,*)
12 18 (*Agreement, lepingut,*)
21 9 (*Multilateral, mitmepoolsete*)
22 10 (*Trade, kaubanduslääbirääkimiste*)
23 10 (*Negotiations, kaubanduslääbirääkimiste*)
24 14 (*hereinafter, (edaspidi)*)
28 15 (*"participants", "osalejad"*)
36 20 (*force, jõustada*)
35 20 (*into, jõustada*)
33 21 (*its, see*)
38 22 (*I, I.*)

39 23 (*January, jaanuaril*)

40 24 (*1995., 1995*)

41 25 (*or, või*)

43 26 (*early, niipea*)

45 28 (*possible, võimalik*)

46 29 (*thereafter., pärast*)

46 30 (*thereafter., seda.*)

Selline väljund on saadud püüdes seostada omavahel kõik leksikaalsed üksused, mis üldse mingit vastastikust seost omavad. Inglise ja eesti keele morfoloogilistest erinevustest tingituna ei sobi selline tulemus eriti hästi leksikograafiatöök, kuid loomulikult saab leksikograafist kasutaja juba väljundi genereerimisele eelnevalt märgistada ainult potentsiaalselt vajalikud seosed, jättes välja korduvad sõnad, numbrid, mitmesõnalised üksused jms. Sellisel juhul võiks väljund (2) välja näha umbes niisugune:

(2)

2 1 (*representatives, Esindajad*)

6 6 (*desirability, soovitav,*)

8 16 (*acceptance, tunnustaksid*)

19 12 (*Round, voorus*)

15 13 (*participants, osalejad*)

14 8 (*all, kõik*)

12 18 (*Agreement, lepingut,*)

21 9 (*Multilateral, mitmepoolsete*)

23 10 (*Negotiations, kaubanduslääbirääkimiste*)

24 14 (*hereinafter, (edaspidi)*)

39 23 (*January, jaanuaril*)

41 25 (*or, või*)

43 26 (*early, niipea*)

45 28 (*possible, võimalik*)

46 29 (*thereafter., pärast*)

Sellise „käsitsi” sõnaparalleelistamise eelis on potentsiaalselt maksimaalne täpsus, kuid miinusteks muidugi täisautomaatsetele lähenemistele vastanduvalt märksa suuremad vajadused inimtööjõu osas ja esialgu ka võimaluse puudumine mitmesõnaliste üksuste joendamiseks.¹⁸

Järgmine samm täisautomaatse sõnastikugenerereerimise suunas on selline poolauto-
maatne protsess, mille puhul tarkvara poolt teostatud sõnaparalleelistuse tulemused
vaadatakse inimkasutaja poolt üle ja vajadusel parandatakse. See on võimalik näiteks
kasutades Chris Callison-Burchi poolt loodud graafilist abivahendit¹⁹, mille sisendiks on
vabavarana saadaval oleva tarkvarapaketi Giza++ (Och, Ney 2000) poolt sõnatasandil
paralleelistatud paralleeltekst. Kasutajale kuvatakse paralleelistus maatrikstabelina, kus on
hõlpsasti võimalik parandusi teha.

Giza++²⁰ tööprintsüübid hõlmavad sarnaselt suurema osaga statistilise masintõlke
rakendustest IBM-i uurijate mudeleid, mida tutvustasid Peter F. Brown jt (1993).²¹

¹⁸ On küll olemas ka samalaadne joondamisvahend, mis võimaldab lisaks sõnadele ka fraaside joondamist, kuid selle kasutamine leksikograafilistel eesmärkidel nõuaks väljundi lisatöötlust – vt

<http://www.isi.edu/~hdaume/HandAlign/> (03.08.2006)

¹⁹ Vt <http://demo.linearb.co.uk:8080/sandbox/start.jsp> (21.08.2006)

²⁰ Vt <http://www.fjoch.com/GIZA++.html> (21.08.2006)

²¹ Näide Giza++ sõnaparalleelistustest koos Giza++ sisendiks olnud TÜ korpuse alamosaga on allalaaditav aadressilt www.teataja.ee/leksikonid.zip.

IBM-i mudel 1 leiab sõnadevahelised vastavused lausetasandil joondatud bitekstist sõnade koosinemise alusel, alustades ühtlustatud tõlkevastetõenäosustest. Mudel 2 lisab sellele lihtsale tõlkemudelile positsioonilised parameetrid ning mudel 3 nn viljakusparameetrid. Viljakusparameetritega tuuakse esile mõnede sõnade kalduvus olla tõkelises ühenduses tõenäolisemalt ühe- või mitmesõnalise vastega, mis sisaldab teatud arvu sõnu. Mudel 4 hõlmab meetodeid mitmesõnaliste üksuste tuvastamiseks biteksti põhjal genereeritud sõnaklasside võrdlemise abil lauses. Mudeliga 5 on püütud parandada eelmiste mudelite töö käigus esile tulnud vigasid. (Brown jt 1993)

4.4.4. Täisautomaatne leksikoni genereerimine tarkvarapaketi PWA abil

Rootsi teadlaste projekti Plug raames välja töötatud kahe joondamisrakenduse – Linköping Word Aligner (LWA) ja Uppsala Word Aligner (UWA) – loomisel ja arendamisel on muuhulgas arvesse võetud eespool kirjeldatud sõnade üksühese vastavuse eeldusega seotud probleeme. Mõlemad süsteemid kasutavad tekstide sõnatasandil joondamiseks võrdlemisi vähe keelespetsiifilist infot ning on seetõttu üpris lihtsalt rakendatavad ka eestikeelse osalusega sõnastike koostamiseks; mõlemad süsteemid on koostatud programmeerimiskeeles Perl ning on internetist tasuta allalaaditavad.

LWA ja UWA on integreeritud paralleelkorpuste töötlemiseks loodud laiema funktsionaalsusega tarkvaraplatvormi nimega Uplug (Tiedemann 2002). Projekti Plug käigus loodud süsteemide peamiseks rakendusvõimalusteks on peetud (Sågvall Hein 2002) erinevate masintõlkeliikide täiendamist tõlkeinfo ja leksikonidega ning samuti inimtõlkijate abistamist leksikonide loomise läbi. Vaatlen järgnevalt mõlemat joondamisprogrammi ja nende kasutusvõimalusi lähemalt.

UWA sisendiks on lause- või fraasitasandil joondatud bitekst. Operatsioonisüsteemi Windows jaoks praegu saadaval olev PWA (Plug Word Aligner – tarkvarapakett, mis ühendab endas UWA ja LWA) versioon on eelseadistatud rootsi-, inglise- ja saksakeelse sisendteksti jaoks, kuid seadistusi on võimalik kohaldada ka teistele keeltele.

Sisendtekst jagatakse programmi poolt kas juba olemasoleva (kasutaja poolt lisatud) lihtsavormilise sõnastiku või teksti enese põhjal ühe- või mitmesõnalisteks üksusteks. Sõnastiku puudumisel arvestatakse selle etapi juures tähejärjendite sagedusi ja pikkust, samuti sõnatüüpe ja punktuatsiooni.

Järgnevalt püütakse mõlema keele vastavad üksused omavahel kokku viia. See protsess algab nn „kindlate juhtumite” (nt kui korpuse paralleelistatud segment kujutab endast vaid ühesõnalist elementi) eristamisest. Seejärel hindab süsteem sõnede omavahelist sarnasust ja suhtelist paiknemist tekstis (arvestades konteksti) ja märgib tõlkevastete kandidaatidena ära eelnevalt seatud lävendid ületanud sõnepaarid. Eraldi üritatakse seostada vähesagedasi sõnesid. Viimaks toimub tõlkevastete automaatne filtreerimine, millele võib järgneda tulemuste „käsitsi” korrigeerimine.

4.5. Kakskeelse leksikoni genereerimine PWA abil paralleelkorpuste põhjal

4.5.1. Uppsala Word Aligner

Järgnevalt kirjeldan katset luua UWA abil lähtematerjal erialasõnastiku koostamiseks või täiendamiseks, kasutades selleks mahukat paralleelkorpust (730 880 paralleelset ühest või mitmest lausest või (ala)pealkirjast koosnevat lõiku, 24 169 586 sõnet).

TÜ inglise-eesti paralleelkorpuse (TÜPK) korpuse ühendasin leksikoni koostamise otstarbeks JRC-Acquis' mitmekeelse paralleelkorpuse²² inglise-eesti alamkorpusega, kuna mõlemad korpused sisaldavad Euroopa Liidu seadusandlusega seotud tekste. Mõlemad korpused on samuti paralleelistatud Vanilla paralleelistaja²³ abil.

²² Vt <http://langtech.jrc.it/JRC-Acquis.html> (21.08.2006)

²³ Vt <http://nl.ijs.si/telri/Vanilla/> (21.08.2006)

JRC-Acquis' korpus on alternatiivina paralleelistatud ka HunAligniga²⁴. Kuigi HunAlign osutus võrdluses Vanillaga (vt osa 5.3.5) oluliselt täpsemaks paralleelistusvahendiks, oli TÜ korpusega leksikoni genereerimise eesmärgil ühendatavaks alamkorpuseks siiski Vanillaga paralleelistatud JRC-Acquis' korpuse inglise-eesti osa, kuna paralleelistusmeetodite võrdluse tulemused ei olnud leksikoni genereerimisel veel teada.

TÜ korpus oli algselt kujul, kus eesti ja inglisekeelsed üksused paiknevad vaheldumisi ja on üksteisest eristatud keele nimetust sisaldavate märgenditega. JRC-Acquis' korpus oli algselt mitmesugust märgendust sisaldaval TEI-kujul, jaotatuna paljudesse failidesse. Et muuta need korpused PWA-le „arusaadavaks”, tuli need korpused teisendada sellisele kujule, mille korral mõlema keele kõik üksused on ümber tõstetud kahte eraldi faili ja on tähistatud omavahel vastavuses olevate numbriliste tähistega. Selleks kasutasin osaliselt Camelia Ignat' poolt loodud Perli-programmi²⁵.

Katse²⁶ käigus genereeritud inglise-eesti õiguskeele leksikon sisaldab 130 865 märksõna (vrd EKI inglise-eesti elektrooniline sõnastik – u. 86 000 märksõna, Eesti Õiguskeele Keskuse terminibaas ESTERM 57 829 märksõna) ja koos tõlkevastetega 482 571 sõnet. Sisendina kasutatud koondkorpus hõlmas 730 880 paralleelset ühest või mitmest lausest või (ala)pealkirjast koosnevat üksusepaari. Lisaks sisendkorpusele hõlmas programmi kasutajapoolne sisend ka väiksemahulist morfoloogiainfot faili eesti keele tüüpiliste sõnalõppude ja ebareeglipäraste verbide kohta ning ka lähtesõnastikku, milleks otsustasin valida EKI inglise-eesti sõnastiku²⁷. EKI sõnastik tuli UWA jaoks teisendada sobivale kujule, mis tähendas põhiliselt mitmesugust lühendamist – mittesobivate kirjete, sulgudes asuvate märkuste, liigsete vastete jms automaatset kustutamist UNIXI vahenditega.

²⁴ Vt <http://mokk.bme.hu/resources/hunalign>

²⁵ Vt <http://wt.jrc.it/lt/Acquis/JRC-Acquis.2.2/alignments/index.html> (21.08.2006)

²⁶ Nii UWA kui LWA abil genereeritud leksikonid on allalaaditavad aadressil www.teataja.ee/leksikonid.zip (29.08.2006). Lisaks koondkorpusest ekstraheeritud leksikonidele leiab sealt ka üksnes TÜ korpuse põhjal UWA-ga genereeritud sõnastiku ning näite Giza++ sõnaparalleelistustest koos Giza++ sisendiks olnud TÜ korpuse alamosaga.

²⁷ Vt <http://www.eki.ee/dict/inglise/> (21.08.2006)

Järgnevalt genereeritud eesti-inglise leksikoni sisendkorpuseks olid üksnes TÜ paralleelkorpuse inglise-eesti ja eesti-inglise osad.²⁸ See leksikon sisaldab kokku 97 300 märksõna ja koos tõlkevastetega 318 025 sõnet.

Esitan juhusliku parandusteta fragmendi (3) UWA poolt TÜ inglise-eesti paralleelkorpuse ja JRC-Acquis' mitmekeelse paralleelkorpuse inglise-eesti alamkorpuse põhjal genereeritud leksikonist, mis ei ole esitatud UWA graafilise liidese vahendusel, vaid tavalise tekstifailikatkendina:

(3)

```
{reserve officer}
{
  IX:reservohvitseri
}
{reserve officer candidates}
{
  IX:reservohvitserikandidaat
}
{reserve officer courses}
{
  IX:reservohvitserikursusel
}
{reserve positions and}
```

²⁸ Fragment üksnes TÜ paralleelkorpuse põhjal genereeritud eesti-inglise leksikonist on ära toodud käesoleva töö esimeses lisas.

{
2X:reservipositsioonid ja
 }
 {*reserve power system*}
 {
1X:reservelektrisüsteem
 }
 {*reserve ratio*}
 {
3X:reservibaasi
2X:reservimäär
 }
 {*reserve service*}
 {
1X:reservteenistus

Nagu näitest (3) paistab, on väljundsõnastikus ka eesti keele puhul üsna edukalt lähenenud juhtumid, mille puhul eestikeelsele sõnale vastab inglise keeles mitmesõnaline üksus. Küll aga hakkavad silma rohked keeltevahelistest morfoloogilistest erinevustest tingitud ebatäpsused²⁹, mis eksisteerivad samamoodi eesti-inglise keelepaari jaoks nagu ka rootsi ja inglise keelte jaoks (UWA on loodud arvestades rootsi keele eripärasid, kuid mitmeid väljundi vigasid ei ole suudetud ka rootsi keele puhul esialgu parandada). Samuti

²⁹ Seda probleemi võib tõlgendada ka andmehõreduse efektist lähtudes, nagu tehti näiteks inglise-eesti statistilise masintõlke tulemuste hindamisel (vt Fishel jt 2007).

on probleemiks ka inglise-eesti statistilise masintõlke puhul mainitud (Fishel jt 2007) inglise ja eesti keele erinev sõnajärjekord.

Kuna UWA paneb tõlkevastavuste leidmisel suurt rõhku sõnesarnasusele, siis tulenevad sellest inglise ja eesti keele mittesuguluse tõttu paljud vead ja osalised vastavused väljundis, näiteks:

(4)

{selling}

{

1X:Lepingus

1X:Selleks

1X:eelkõige

2X:lepingu

3X:müümine

10X:selle

1X:selles

1X:sellest

1X:selline

2X:sellise

1X:sellisel

1X:selliselt

2X:selliste

}

Anna Sågvall Hein (2002: 74) väidab, et erinevast morfoloogiast tulenevad ebatäpsused lahendaks rootsi keele puhul olemasoleva ülekandepõhise tõlkesüsteemi (Multra) liitmine UWA-ga. Eesti keele jaoks sellist tõlkesüsteemi kahjuks esialgu ei leidu. Vähene keelespetsiifiline morfoloogiainfo ei võimalda UWA-l teostada automaatset lemmatiseerimist, kuid oletan, et tõenäoliselt parandaks väljundsõnastiku korrektsust siiski sõnastiku genereerimisele eelnev sisendkorpuse lemmatiseerimine muude vahenditega.

Tulemust parandaks kindlasti JRC korpuse (sisendkorpuse) automaatse paralleelsete manuaalne korrigeerimine, samuti süsteemi poolt kasutatavate statistiliste meetmete osaline laiendamine sisendkorpusele internetile andmehõreduse efekti kahandamise eesmärgil. St üks süsteemi alam-moodul võiks olla ühendatud mõne interneti (mitmekeelse) otsimootoriga, millele võib esitada samu päringuid kui sisendkorpusele.

PWA-d puudutavatest publikatsioonidest ei selgu, kas süsteem kasutab rootsi keele puhul morfoloogilist infot ka selleks, et parema joondamistulemuse eesmärgil liitsõnad eelnevalt koostisosadeks liigendada (nagu seda on mitmel puhul tehtud saksa-inglise sõnaparalleelisel – nt (Déjean jt 2002: 6)).

Genereeritud leksikoni saaks täiustada ka mitmesuguse automaatse järeltöötusega, mis võiks tegelikult olla ekstraheerimistarkvarasse integreeritud. Näiteks oleks soovitatav automaatselt kustutada kõik üksnes mittealfabeetilisi sümboleid sisaldavad märksõnad, kui soovitud tulemus kujutab endast loomuliku keele leksikoni.

Kasutaja poolt lisatud morfoloogiainfo põhjal saaks sõnastiku genereerinud moodul kustutada ka ühe ja sama märksõna erinevaid muutelõppe sisaldavatest, kuid muus osas sarnastest tõlkevastetest ebavajalikud. Samuti ei oleks vaja kohelda erinevate leksikaalsete üksustena identseid sõnu, mille algustähed on lausesisesest positsioonist tingituna erineva suurusega.

A. Sågvall Hein näeb süsteemi arendusvõimalustena lisaks selle täiendamisele reegli- põhiste meetoditega veel sõnastiku märksõnade lemmatiseerimist ning viidete lisamist genereeritud sõnastikukirjetelt kontekstidele korpuses (Sågvall Hein 2002: 74).

4.5.2. Linköping Word Aligner

LWA loomisel on toetunud Pascale Fungi ja Kenneth Churchi (1994) ning Dan Melamedi (1997) joondamisalastele uurimustele. LWA hõlmab nagu UWA-gi põhiliselt statistilisi meetodeid. Algoritm on iteratiivne – tõlkevasted genereeritakse biteksti põhjal, siis kustutatakse genereeritud sõnapaarid bitekstist ja korratakse tsüklit. Statistilistel tõenäosustel põhinevat baasalgoritmi täiendavad neli lisamoodulit ning programmi kasutajaliides võimaldab seadistada mitmeid parameetreid (lisatestid, lävendid, tsüklite arv).

Esimene moodul alustab tööd sõnade jagamisega etteantud info põhjal mitmesugustesse kategooriatesse ja alamkategooriatesse: relevantseteks ja irrelevantseteks (irrelevantsetena on määratletud näiteks inglise keele abiverb *do*, millel enamasti puudub üksühene vaste teises keeles), suletud ja avatud sõnaklassi sõnadeks – avatud sõnaklassi kuuluvaid sõnu saab edaspidi joondada vaid teiste sama sõnaklassi sõnadega ja suletud sõnaklassi sõnu saab vastavalt joondada ainult suletud sõnaklassi sõnadega. Suletud sõnaklassi sõnad jagatakse järgnevalt veel alamkategooriatesse. Kategooriate põhjal toimub iteratiivne paralleelistamine.

Morfoloogiamoodul tunneb vastava keele sufiksiloendi järgi ära ühe sõna erinevad vormid. Kui leksikaalsete üksuste paar (X, Y) on mõne eelneva tsükli käigus tunnistatud tõlkevastavuses olevaks, siis otsib moodul teisi kandidaatpaare, mille esimene element on X ja teine element Z , nõnda et leiduvad ka sõned W, F ja G , mille puhul $Y = WF$ ja $Z = WG$ ning F ja G sisalduvad sufiksiloendi ühes ja samas paradigmas. Kui leitakse mitu erinevat üksust, mida saab tähistada sümboliga Z ja mille sufiksud kuuluvad erinevatesse paradigmadesse, siis valitakse neist suurima sagedusega paradigma.

Järgmine moodul tegeleb mitmesõnaliste üksuste paralleelistamisega, kasutades selleks samasuguseid võtteid kui ühesõnaliste üksuste puhulgi. Mitmesõnaliste üksuste kogum koosneb süsteemi poolt automaatselt leitud üksustest ning eelnevalt lisatud keelespetsiifilistest kollokatsioonidest. Lõpuks parandatakse paralleelistust veel vastavalt

joondatavate segmentide suhtelisele asendile tekstides. LWA miinuseks on UWA-ga genereeritavatele leksikonidele sarnase sagedusinfo puudumine väljundist.

Sisendkorpuse mahupiirangu tõttu kasutasin LWA-ga katselise sõnastikumaterjali genereerimiseks TÜ paralleelkorpuse 10 000 paralleelsest lõigust koosnevat alamkorpust. Toon siin ära fragmendi³⁰ (5) tulemuseks saadud leksikonist, mis sisaldas 8516 kirjet:

(5)

<i>kinnitades</i>	<i>confirming</i>
<i>kinnitades</i>	<i>reaffirming</i>
<i>kinnitatud</i>	<i>affixed</i>
<i>kinnitava</i>	<i>assurance</i>
<i>kinnitavad</i>	<i>reaffirm</i>
<i>kirikute</i>	<i>churches</i>
<i>kirjalik</i>	<i>written</i>
<i>kirjalike</i>	<i>written</i>
<i>kirjalikke</i>	<i>written</i>
<i>kirjaliku</i>	<i>written</i>
<i>kirjalikult</i>	<i>writing</i>
<i>kirjalikust</i>	<i>written</i>

4.5.3. Tulemuste hindamine: ARCADE ja PWA

Joondamisprogrammide võrdleva hindamise standardiseerimisprojekti ARCADE (Véronis, Langlais 2000) raames on kindlaks tehtud, et parimate sõnadevahelise paralleel-

³⁰ Pikem fragment on esitatud magistritöö lisas 2.

listamise süsteemide täpsus ja saagis ulatuvad umbes 75 protsendini (Véronis, Langlais 2000: 386). Seejuures erineb tulemus sõnaliigiti, ulatudes 94 protsendini adjektiivide puhul ja ainult 60–70 protsendini verbide puhul (samas).

Nii hea tulemuse saavutab joondamissüsteem aga üksnes rohke keeletespitsiifilise info kasutamiseiga bitekstide analüüsil. Kuna UWA ja LWA loomisel on järgitud keeltevahelise portatiivsuse eesmärki, siis on nende programmide töötulemuste hindamisel saadav täpsus ja saagis oluliselt väiksem. Projekti ARCADE juhendite järgi arvatud rootsi-inglise näidiskorpusest ekstraheeritud leksikoni täpsus on UWA poolt teostatud joondamise puhul üksnes 42,2 % ja saagis 37 %; LWA puhul on samad näitajad vastavalt 51 % ja 41,3 %.

Mitmesõnaliste üksuste korral on saagise arvutamine tunduvalt keerulisem kui sellise vastavuse korral, mille puhul ühele sõnele vastab alati ainult üks sõne. On väga raske hinnata terves korpuses sisalduvate mitmesõnaliste üksuste koguhulka ja seetõttu peaks mitmesõnalisi üksusi kontrolli kaasates teostama arvutuse väga väikese tekstinäite põhjal.

ARCADE-projektis esitatud meetodid ei võimalda hinnata seda, kui hästi süsteem jagab biteksti lähtepoole üksused korrektseteks ühe- või mitmesõnalisteks üksusteks, vaid keskenduvad süsteemi poolt pakutud sihtkeele sõnade võrdlemisele etaloniga ehk nn „kuldstandardiga” (Ahrenberg jt 2000b: 4). Seetõttu on PWA autorid välja arendanud uue kontrollmeetodi (samas), mis erinevalt ARCADE-st võtab arvesse ka mitmesõnalisteks üksusteks jaotamist nii lähte- kui sihtkeele siseselt ja samuti erinevaid juhtumeid, mille korral mitmesõnaliste üksuste joondamist saab lugeda korrektseks vaid osaliselt. Sama rootsi-inglise näidiskorpuse põhjal süsteemiga UWA genereeritud leksikoni täpsus on selle meetodi kohaselt 71,8 % ja saagis 37,4 % (LWA-ga genereeritud leksikoni puhul on need näitajad vastavalt 71,9 % ja 42,6 %).

4.5.4. Eesti-inglise paralleelkorpuste põhjal genereeritud leksikonide hindamisest

PWA tarkvarapakett sisaldab muuhulgas moodulit joondamistulemuste automaatseks hindamiseks nii ARCADE kui ka PWA meetodil. Mooduli kasutamine eeldab PWA nn

kuldstandardi formaadis võrdlusfaili olemasolu vastava keelepaari jaoks, kus korrektsed vastavused on registreeritud koos mitmesuguse juurdekuuluva infoga ning näitelausetega. Eesti-inglise keelepaari tarbeks selline kuldstandard esialgu puudub. PWA kuldstandardifailide loomiseks on kasutatud samuti tarkvaraplatvormi Uplug kuulunud programmi Plug Link Annotator (PLA), kuid kahjuks ei ühildu PLA enam aastal 2006 kasutatavate operatsioonisüsteemide ja muu vajaliku tarkvaraga (Java) ega ole seetõttu kasutuskõlblik. Seega tuleb standardfaili loomiseks leida mõni alternatiivne poolautomaatne võimalus või teostada kontroll käsitsi.

Ka PWA hindamismeetodi tulemused varieeruvad üsna palju sõltuvalt tekstide žanrist ja kuldstandardi koostamise kriteeriumitest – st sõltuvalt sellest, millise sagedusega ja funktsiooniga sõnu ja millistes proportsioonides valitakse nende näidete hulka, mille alusel paralleelistustulemusi hinnatakse (Ahrenberg jt 2000b: 6).

Rootsi-inglise korpusest pärineva kolme erineva tekstikategooria UWA-süsteemi joondamistulemuste hindamine PWA-meetodil andis kõige kõrgema tulemuse (täpsus: 81,26 %; saagis: 64, 47 %) tehnilise sisuga tekstide puhul ja kõige madalama tulemuse (täpsus: 69,04 %; saagis: 41,44 %) poliitiliste tekstide puhul (samas). Vahepealse tulemuse saavutas ilukirjanduse alamkorpus (samas). Võib eeldada, et nende kolme tekstikategooria hulgast asetub TÜ eesti-inglise õigusakte ja seadusi sisaldav paralleelkorpus žanriliselt kõige lähemale poliitika-alastele tekstidele ning oodatav tulemus leksikoni täpsuse ja saagise hindamisel peaks olema võrreldav tulemusega, mille nimetatud katse käigus andis see tekstikategooria.

See oletus pidas tulemuste kontrollimisel vähemalt osaliselt paika: UWA-ga TÜ paralleelkorpuse ja JRC-Acquis' korpuse põhjal genereeritud inglise-eesti leksikonist juhuslikult valitud 50 kirje kontrollimisel sain selle alamosa täpsuseks 60 %. UWA-ga ainult TÜ paralleelkorpuse põhjal genereeritud eesti-inglise leksikonist juhuslikult valitud 50 kirje kontrollimisel sain selle alamosa täpsuseks 87 %. Samasugune kontroll LWA-ga

TÜ eesti-inglise paralleelkorpuse põhjal genereeritud leksikoni 50-kirjelise alamosa peal andis täpsuseks 61 %.³¹

Paremad tulemused eesti-inglise leksikoni kontrollimisel on seletatavad TÜ paralleelkorpuse parema paralleeljustusega (võrreldes JRC-Acquis' korpuse Vanilla versiooniga) ja sellega, et sisendkorpuse suure mahu tõttu ei suutnud programm inglise-eesti leksikoni genereerimisel läbi viia viimast, automaatse filtreerimise etappi.

Saagise hindamine ilma kuldstandardfailita oleks korpuse suure mahu tõttu keerulisem ülesanne ja väikese alamosa põhjal saadud hinnangud ei pruugi olla adekvaatsed, seetõttu saagise arvutamisest esialgu loobusin.

Uplugi tarkvara abil on proovitud genereerida ka kreeka-inglise leksikoni (Charitakis 2007), mille manuaalsel hindamisel otsustati hindamisprotsessi mitte kaasata tõlkepaare, mille esinemissagedus oli alla kolme. Ülejäänud tõlkepaarid jagati viide gruppi sõltuvalt nende esinemissagedusest ja kõiki gruppe vaadeldi hindamisel eraldi. Leiti, et tõlkepaaride esinemissagedus ja tõlke korrektsus on genereeritud leksikoni puhul otseses proportsionaalses vastavuses – suurema esinemissagedusega tõlkepaaride korrektsuse protsent on suurem. (Charitakis 2007)

Sama tendentsi kinnitab ka Uplugiga genereeritud inglise-eesti leksikonide vaatlus. Kontrollisin kümmet juhuslikku tõlkepaari UWA-ga TÜ paralleelkorpuse ja JRC-Acquis' korpuse põhjal genereeritud inglise-eesti leksikonis, mille esinemissagedus oli 3. Nendest 5 (50 %) olid ebakorrektsed. Seejärel kontrollisin kümmet juhuslikku tõlkepaari esinemissagedusega 20. Nende seas leidis ainult üks tõlkepaar, mida võis lugeda ebakorrektses.

Projekti Plug meetoditele väga sarnast lähenemist leksikoni genereerimisele esindavad Dan Tufiş ja Ana-Maria Barbu (2001), kes kasutasid leksikoni koostamiseks europrojekti

³¹ Täpsuse arvutamisel lugesin leksikoni kirje korrektseks (1 punkt), kui tõlkevastete seas leidis vähemalt üks täielikult märksõnale sisuliselt vastav tõlkevaste ja osaliselt korrektseks (0,5 punkti), kui tõlkevastete seas leidis vähemalt üks osaliselt märksõnale sisuliselt vastav tõlkevaste. Täpsuseks lugesin kirjete arvu (50) jagatud summeeritud punktide arvuga.

Multext-East³² raames valminud paralleelkorpust, mis sisaldab George Orwelli romaani „1984” kaheksas, sh eesti keeles.

D. Tufiş ja A.-M. Barbu kirjeldavad muuhulgas ka väikese eesti-inglise sõnastiku ekstraheerimist ning tulemuste kontrollimist. Kontroll teostati „kuldstandardi” puudumisel käsitsi ja eesti-inglise sõnastiku hindamisel saadi selle täpsuseks 96,2 % ja saagiseks 57,9 %. Paremaid tulemusi võrreldes Plug-projektiga võib seletada mitmesuguse eeltöötusega, mida oli rakendatud MULTEXT-East korpusele ja millele toetus tõlkevastekandidaatide loendi ekstraheerimine: lisaks lausetasandil paralleelistusele ka leksikaalsete üksuste segmenteerimine, morfoloogiline ühestamine ja lemmatiseerimine.

Sõnatasandil paralleelistamist saab teostada ka statistilise masintõlke rakenduste abil, eespool mainisin tarkvarapaketti Giza++. TÜ inglise-eesti korpuse paralleelistamise tulemuste kohta Giza++ abil vt K. Muischnek (2006).

4.5.5. Võrdlus ESTERMiga

Eesti Õigustõlke Keskus on koostanud tõlgitud Euroopa Liidu ja Eesti Vabariigi õigusaktide põhjal terminibaasi ESTERM³³, mis umbes neljandiku ulatuses sisaldab ka fraase. Lisaks õiguskeelele sisaldab ESTERM sarnaselt paralleelkorpuste põhjal genereeritud leksikonidega termineid ka paljudest teistest valdkondadest, mille käsitlemine seadusandlikes tekstides on vajalik.

Üks reaalsemaid inglise-eesti paralleelkorpuste põhjal genereeritud leksikonide kasutusvõimalusi näib olevat ESTERMI terminibaasi täiendamine uute terminitega leksikonidest saadud info põhjal. Kontrollimaks seda hüpoteesi, esitasin ESTERMI terminibaasile päringu kahekümne viie UWA-ga TÜ paralleelkorpuse ja JRC-Acquis' korpuse põhjal genereeritud leksikonist vabalt valitud ja leksikonis korrektset tõlget omava termini kohta.

³² Vt <http://nl.ijs.si/ME/> (21.08.2006)

³³ <http://mt.legaltext.ee/esterm/> (29.05.2007)

Kaheteistkümmel juhul kahekümne viiest esines leksikonist valitud märksõna koos oma tõlkevastega täpselt samasugusel kujul ka ESTERMIS. Üheteistkümmel juhul puudus leksikoni kirje täielikult ESTERMist. ESTERMist puuduvad kirjed olid järgmised:

- 1) *rewinder* ('üंबरkerija')
- 2) *active biocides* ('aktiivbiotsiidid')
- 3) *post-release* ('keskkonda viimise järgne')
- 4) *allocation method* ('eraldamismeetod')
- 5) *pyraclostrobin* ('püraklostrobiin')
- 6) *pyrasolidon* ('pürasolidoon')
- 7) *textile committee* ('tekstiilikomitee')
- 8) *textile pulp* ('tekstiilimass')
- 9) *piston rod* ('kolvivarras')
- 10) *efficiency indicator* ('toimearv')
- 11) *terminator* ('terminaator')

Ühel juhul ei sisaldanud ESTERM täpsel kujul inglisekeelset märksõna ('*bioaccumulated*'), kuid sisaldas selle eestikeelset vastet ('*bioakumuleeruv*' inglisekeelse tõlkega '*bioaccumulative*'). Samuti ei olnud ESTERMis märksõna '*semiconductor product*' ('*pooljuhttoode*'), kuigi hädusa otsingu meetodil on võimalik terminibaasist leida pikem termin '*topography of semiconductor products*' ('*pooljuhttoodete topoloogia*'), mille koostisosaks on leksikonist valitud termin.

Kuna päringutest rohkem kui 50 % korral puudus ESTERMis leksikonis esinenud termin täpselt samasugusel kujul, siis leidis kinnitust oletus, et on võimalik ESTERMi terminibaasi täiendamine uute terminitega leksikonidest saadud info põhjal.

5. Paralleelkorpuste võrdlemine ja paralleelistuse kvaliteedi hindamine³⁴

Kõige olulisem kvaliteedikriteerium on paralleelkorpuse puhul see, et paralleelüksuse koosseisus olev tõlge vastaks samas üksuses olevale originaalile. Juhul, kui paralleelkorpus luuakse (pool)automaatselt, on sellise nõude 100%-line täitmine pea võimatu. Ka paralleelistuste korrektsuste automaatne hindamine on väga raske, sest tegemist on tähenduste omavahelise vastavuse hindamisega.

Seetõttu jäetakse suurte paralleelkorpuste puhul hindamise küsimus sageli üldse tõstmata: nt OPUS (Tiedemann, Nygaard 2004), Europarl (Koehn 2002), tšehhi-inglise paralleelkorpus (Bojar, Žabokrtský 2006). Mõne korpuse dokumentatsioon sisaldab küll hoiatust, et paralleelistus on tehtud automaatselt või öeldakse lausa, et paralleelistustes on vigu. Ainus viis korpuse paralleelistuse kvaliteeti hinnata on teha seda käsitsi, väiksema osavalimi peal, mis on aga äärmiselt töömahukas, vt nt. (Samy jt 2006; Singh, Husain 2005).³⁵

Et peatükis 5.5. kirjeldatud leksikoni genereerimise ja masintõlkealaste eksperimentide tulemusi parandada, tuleks olemasolevad eesti keele osalusega paralleelkorpused ühendada nõnda, et loodavast koondkorpusest välja jätta osalised kattuvused ja võimalikult palju vigaselt paralleelistatud osi. Enne korpuste ühendamist on vaja ühtlasi hinnata korpuste

³⁴ Viies (ja osaliselt ka kuues) peatükk rajaneb koos Heiki-Jaan Kaalepiga kirjutatud ja 2007. a septembris Kopenhaagenis toimuvale XI masintõlke-alasele tippnõupidamisele (MT Summit) stendiettekandena avaldamiseks esitatud artiklil. Suurem osa siin kirjeldatud praktilisest tööst on minu poolt läbi viidud. Heiki-Jaan Kaalepilt pärinevad mitmed ideed ja soovitusel, millest praktilise töö teostamisel lähtusin. Samuti on viiendas peatükis osaliselt kasutatud nimetatud artiklis esitatud töö kirjeldust H.-J. Kaalepi sõnastuses.

³⁵ Paralleelistuskvaliteedi hindamist puudutavad küsimused kattuvad osaliselt korpustest automaatselt genereeritud mitmekeelsete leksikonide kvaliteedi hindamise probleematikaga, millest oli lähemalt juttu peatükis 5.5.4.

paralleeliste kvaliteeti, et teha kindlaks, millist korpust kattuvate osade kaasamisel eelistada. Järgnevalt kirjeldan katset hinnata paralleelkorpusete kvaliteeti poolautomaatselt võrdlemismeetodil. Erialasest kirjandusest ei ole teada uurimistöid, mille käigus oleks varem paralleelkorpusi nende kvaliteedi hindamise eesmärgil võrreldud.

5.1. Korpusete kirjeldused

5.1.1. Maht

Kaks teadaolevalt kõige mahukamat inglise-eesti paralleelkorpuset on Tartu Ülikooli inglise-eesti ja eesti-inglise õigustekstide paralleelkorpus (TÜPK) ja Ispras loodud JRC-Acquis'³⁶ paralleelkorpusete inglise-eesti alamkorpus.

TÜ korpusetes on kaks osa, mida eristab lähte-keel. Esiteks eesti seadused ja nende tõlked inglise keelde, 150 000 paralleelüksust (lauset või loendi elementi) 400 tekstis. Eesti keeles 1,7 miljonit sõnet, inglise keeles 2,9 miljonit sõnet. Teiseks Euroopa Liidu õigusaktid ja nende eestikeelsed tõlked, 280 000 paralleelüksust (lauset või loendi elementi) 4000 tekstis. Eesti keeles 3,3 miljonit sõnet, inglise keeles 4,9 miljonit sõnet.

JRC-Acquis' korpusete paralleelstatud versioonis on eesti keeles 4,6 miljonit sõnet, inglise keeles 6,8 miljonit sõnet. JRC-Acquis' korpus hõlmab samuti Euroopa Liidu seadusandlike tekste kokku 21 Euroopa keeles, kusjuures keskmiselt on keele kohta 8,8 miljonit sõnet 7600 tekstis. JRC-Acquis' korpusete paralleelstatud versioonis on eesti keeles 4,6 miljonit sõnet, inglise keeles 6,8 miljonit sõnet. Korpusete paralleelstatumata versioonis on korpusete suurused eesti ja inglise keele osas vastavalt 7,2 miljonit ja inglise keeles 9,9 miljonit sõnet. Seega sisaldab paralleelstatumata korpusete inglise-eesti alamosa oluliselt rohkem sõnu kui paralleelstatud versioon.

JRC-Acquis' korpusete paralleelstatumata osa sisaldab 236 796 (eesti osa) + 193 260 (inglise osa) paralleelüksust (kokku 430 056 paralleelüksust). TÜ korpusete paralleelstatumata

³⁶ (<http://langtech.jrc.it/JRC-Acquis.html>)

inglise-eesti osa (ilma eesti-inglise osata) hõlmab 176 921 (eesti osa) + 746 219 (inglise osa) paralleelüksust (kokku 923 140 lõiku). Paralleelüksuste suurem arv TÜ korpuses on ilmselt tingitud korpuse paralleelistamata toorversiooni liigendusvigadest (laused olid jaotunud mitme lõigu vahel).

JRC-Acquis' korpuse kõigi keelepaaride paralleelistamiseks saab kasutada kahte paralleelistajat: Vanilla ja HunAlign³⁷. JRC-Acquis' korpuses on iga keelepaari kohta keskmiselt 269 148 paralleelüksust. Inglise-eesti keelepaari kohta on Vanilla paralleelistuses 295 189 paralleelset üksust ja HunAligni paralleelistuses 301 647 paralleelset üksust - keskmisest keelepaarist Ispra korpuses on inglise-eesti osa mõlemal juhul suurem.

paralleelistus	TÜ korpus		JRC-Acquis' korpus		
	lausetasandil paralleelistatud	lõigutasandil paralleelistatud ³⁸	paralleelistamata kujul	JRC Hunalign	JRC Vanilla
Sõnede hulk (milj)	12,7	8,2	16	11,5	11,5
paralleelüksuseid	436 043	230 580	430 056	301 647	295 189

Tabel 1. Korpuste suurused

JRC Hunalign – JRC-Acquis' korpuse inglise-eesti alamosa paralleelistatuna Hunaligniga
 JRC Vanilla – JRC-Acquis' korpuse inglise-eesti alamosa paralleelistatuna Vanillaga

5.1.2. Paralleelistused

Mõlema korpuse paralleelistamisel on kasutatud keelest sõltumatuid meetodeid. TÜ korpuse paralleelistamisel kasutati Vanillat³⁹. Seejuures järgiti sama strateegiat mida Europarli korpuse loomisel (Koehn 2002): kui paralleelistatavate tekstiosade formaalsed struktuurid olid liiga erinevad, siis jäeti vastavad osad korpusest hoopis välja.

³⁷ <http://mkk.bme.hu/resources/hunalign>

³⁸ Ilma eesti-inglise osata

³⁹ <http://nl.ijs.si/telri/Vanilla/>

Paralleelistamine toimus kolmes etapis: esiteks peatükkide ja lisade paralleelistamine, seejärel lõikude ja lõpuks lausete. Seejuures kasutati igal etapil kontrollimiseks (mitte paralleelistamiseks) kõige lihtsamaid ankurpunkte: peatükkide, paragrahvide, loendite nummerdust.

Juhul, kui paralleelsetes tekstides paistis olevat erinev arv paragrahve, artikleid või nummerdatud loendite elemente; või kui vastavate järjekorranumbritega elemendid ei olnud kohakuti, siis vastavaid tekste paralleelkorpusesse ei võetud. Põhjuseks oli oletus, et selliste tekstide formaalne struktuur on liiga erinev, et neid antud lihtsa meetodiga töödeldes oleks tulemused usaldusväärsed.

TÜ korpuses võivad eesti- ja ingliskeelsed laused olla omavahel 1-1, 1-2 või 2-1 vastavuses. Muud vastavused on korpusest välja jäetud. (TÜPK)

HunAligni töö jaotub kolme faasi, millest esimene on paralleeltekstide segmentide esialgne üks-ühele paralleelistus, teine faas paralleeltekstide põhjal automaatne leksikoni moodustamine ja kolmas faas paralleelistamine leksikoni abil. (Steinberger jt 2006)

Ei ole päris selge, kas JRC-Acquis' korpuse HunAligni paralleelistusfail on saadud kõiki kolme faasi läbides, kuid tõenäoliselt on siiski kasutatud ka automaatselt genereeritud leksikoni abi, kuna paralleelistusfailis sisalduvate vastavuste seas on ka mitmeid selliseid, kus ühele lõigumärgenditega tähistatud segmendile vastab teises keeles mitu segmenti, näiteks:

```
<link type="2-1" xtargets="7 8;6">
```

```
<s1><p>THE COUNCIL OF THE EUROPEAN</p><p>COMMUNITIES,</p></s1>
```

```
<s2>EUROOPA ÜHENDUSTE NÕUKOGU,</s2>
```

```
</link>
```

Erinevalt TÜ korpusest ei luba JRC-Acquis' korpuse dokumentatsioon väita, et sealt oleks paralleelistusraskuste tõttu midagi puudu. Esinevad ka 0-1 vastavused.

Et TÜ korpuse avalik versioon on paralleelistatud lausetasandil, JRC-Acquis' aga lõigutasandil, siis ei saa võrdluses kasutada TÜ korpuse avalikku versiooni. Selle asemel kasutasin korpuse autoritelt saadud TÜ korpuse versiooni, millel oli paralleelistamise viimane etapp tegemata ja milles olid seega paralleelistatud ainult lõigud.

5.1.3. Paralleelistusvigadele viitavad tunnused

Püüdsin leida korpuste paralleelistusvigu esile toovaid tunnuseid esmalt juhuslikke faile läbi vaadates ja üks teisega võrreldes. Selle tegevuse käigus selgusid järgmised üldisemad tunnused.

TÜ korpuses on kokku 40 086 juhtumit, mil ühe lõiguga on paralleelistatud täpselt sama sisuga teine lõik (vahe on ainult märgendites: <inglise> vs <eesti>). Näiteks:

<inglise> *La cantidad exportada no debe superar* </inglise>

<eesti> *La cantidad exportada no debe superar* </eesti>

Enamasti on sellistel juhtudel tegemist mõne kolmanda keele tsitaatidega (mitte inglise ega eesti keeles), mis tähendab, et selliste vastavuste rohkus ei viita mitte vigadele, vaid pigem paralleelistuse korrektsusele. JRC-Acquis' korpuse Vanilla versioonis on selliseid juhtumeid vaid 1260 (HunAligni paralleelistuses 1443), mis võivad olla tingitud JRC-Acquis' korpuse viletsamast paralleelistusest.

Nende lõigupaaride otsing tõi aga välja, et kohati on lõigupaaride kordumine TÜ korpuses tulenenud ka eesti keele paralleelistamisest eesti keelega, seda vähemalt ühel juhul terve faili ulatuses.

JRC-Acquis' korpuse inglise-eesti osa HunAligniga tehtud paralleelistus sisaldab 15 532 juhtumit, mille puhul eestikeelse jaotusega <s2> on paralleelistatud täiesti tühi (inglisekeelset osa tähistav) segment <s1>. Vanillaga tehtud paralleelistuses on <s1> tühi 10 721 juhul. Märgendile '<s2>' vastav jaotus on tühi HunAligni paralleelistuses 2788

juhul, Vanilla paralleelistuses 1475 juhul. TÜ lõppkorpuses analoogseid tühje segmente ei leidu, aga võrdluses kasutatavas lõigutasandil paralleelistatud korpuses on kokku 1453 tühja segmenti.

Korpuse loomist käsitlevas artiklis (Steinberger jt 2006) väidetakse, et kahe paralleelistaja võrdlevat hindamist pole veel teostatud. Vanilla joondamistulemusi on küll võrreldud teiste joondamismeetodite abil saadud tulemustega (Rosen 2005). Alexandr Roseni poolt kirjeldatava eksperimendi tulemusel oli Vanillaga paralleelistamisel võrreldes teiste meetoditega kõrgem saagis (kõigi kasutatud tekstide puhul) ning täpsus tekstide puhul, mis ei sisaldanud „müra” (nt väljajätud, tekstide struktuurilised erinevused või eeltöötlusel tehtud vead). (samas)

Tühjade jaotuste tunduvalt suurem hulk HunAligni paralleelistuses tundub pealiskaudsel vaatlusel olevat argument Vanilla kasuks. Siiski ei saa järeldada tühjade segmentide suuremast arvust paralleelistuse suuremat korrektsust.

JRC-Acquis' korpuse Vanilla paralleelistuse kahjuks räägib asjaolu, et JRC korpuse tegijad ei võtnud korpuse paralleelistamisel Vanillaga arvesse nn "kindlaid" (*reliable*) seoseid paralleeltekstide vahel (näiteks nummerdatud lõigud). HunAligni algoritm aga vaatleb eraldi ka paralleelteksti segmentides sisalduvaid numbrilisi sarnasusi. (Steinberger jt 2006)

Korpuste kattuvate alamosade inglise ja eesti osade mahtu (baitides) failiti võrreldes selgus, et suhteliselt on JRC Vanilla paralleelistuses kahe keele vahelised mahulised erinevused failiti palju suuremad kui TÜ korpuses. See tähendab, et tõenäoliselt on TÜ korpuse paralleelistus korrektsem.

5.2. Eeltöö

Esimene ülesanne kahe korpuse võrdlemisel oli kindlaks teha, millised tekstid esinevad mõlemas korpuses, sest nende võrdlemine võimaldab korpuse hindamist automatiseerida ning annab kõige rohkem infot korpuste tegemise protseduuri mõjust tulemusele.

Selliste tekstide leidmine ei ole üldjuhul lihtne ülesanne, sest tekstide samasuse üle otsustamine on keeruline. Samasisulised tekstid võivad olla erinevalt kujundatud, mõned osad võivad ühes tekstis puududa või olla teistsuguses järjestuses. Näiteks seadusandluse puhul võib kuupäev olla teksti alguses või lõpus, mõni lisa võib hoopis puududa jms.

Teiselt poolt esineb ka juhtumeid, mil näivalt väga sarnased tekstid on tegelikult erinevad, nt. juhul, kui uus seadusandlik akt kordab peaaegu sõna-sõnalt mõnda varasemat, või sisaldab sellest suuremahulisi tsitaate. Seega dokumentide automaatne võrdlemine, ka juhul, kui kasutatakse ligikaudset võrdlemist (Levenšteini distantssi), ei pruugi anda nende samasuse kohta adekvaatseid hinnanguid.

Õnneks on JRC-Acquis' korpuse failides kirjas CELEXi kood – Euroopa Liidu seadusandlike dokumentide identifitseerimiseks kasutatav kood. Kõigil EL tõlgetel on sama CELEXi kood mis originaalil. Nii JRC-Acquis' kui ka TÜ korpuse paralleelistatud kujust sai automaatselt teha sellise, kus iga paralleeltekst on omaette fail, kusjuures failis on kirjas ka dokumendi CELEXi kood⁴⁰. Seega oli võimalik kokku viia JRC-Acquis' korpuse ja TÜ korpuse failid, millel on sama CELEXi kood ehk mis peaksid sisaldama sama sisuga ja ühesuguse paralleelistusega tekste.⁴¹ Selliseid failipaare oli u 2000 ja korpuste edasisel võrdlemisel kasutasingi just neid.

Need failipaarid viisin efektiivsema võrdluse huvides omavahel võimalikult sarnasele kujule. JRC-Acquis' korpuse alamosa teisendasin TÜ korpusega sarnasele SGML-kujule,

⁴⁰ Kasutatud skript Unixi-laadses Windowsi-keskkonnas Cygwin: *csplit vanilla_jrc_originaal.xml*
*/*div.type=.*/*{*}*

⁴¹ Lisaks on suur osa mõlema korpuse tekste tähistatud ka teksti **teemavaldkonna** EUROVOC-koodiga. Nende koodide abil tekstide klassifitseerimisest võib abi olla näiteks eestikeelsete terminoloogiasõnastike täiendamisel.

ühtlustasin mõlema alamosa märgendusskeemi. Et failide võrdlust hõlbustada, nimetasin kõik failid automaatselt nendele vastava CELEXi koodiga.⁴²

5.3. Võrdlemine

5.3.1. Metoodika ja algoritm

Võrdlesin omavahel paarikaupa kolme paralleelistatud versiooni eelnimetatud 2000 failist: JRC-Acquis' Vanilla ja HunAligni versioone (edaspidi *JRC Vanilla* ja *JRC HunAlign*) ning TÜ versiooni. Kõigi võrdlemiste metoodika oli ühesugune.

Pärast JRC-Acquis' ja TÜ korpuste märgenduse, täpitähtede kodeeringu ja tühimärkide (tühik, reavahetus, tabulatsioon) ühtlustamist selgus, et tekstid sisaldavad ikka sisu mõttes mitteolulisi erinevusi, nt. mittestandardset viisil esitatud täpitähed, erinevate sulgude ja kirjavahemärkide kasutamine, erinevused numbrite ja arvude esitusviisis, ebaühtlane suur-tähe kasutus või erinevad viisid tähistada tekstidest välja jäetud osi.

Samuti sisaldas mõni fail ühes korpuses rohkem teksti kui sama Celexi koodiga tähistatud fail teises korpuses. Mõnikord oli see tingitud korpuste eri allikatest pärinevate lähtefailide erinevusest, kuid seda võis tingida ka TÜ korpuse paralleelistamismeetod, mis teatud juhtudel raskesti paralleelistatavad osad lihtsalt eemaldas.

Samuti selgus, et korpused sisaldavad lahknevusi ka eestikeelsete sünonüümiliste sõnade ja väljendite kasutamise osas, mis viitab sellele, et tegemist on lihtsalt erinevate tõlgetega või teksti toimetamise erinevate etappidega. Ka need juhtumid kuuluvad sisu mõttes mitteoluliste erinevuste hulka, kui eesmärgiks on paralleelistuse õigsuse kontroll.

⁴² Selleks kasutatud skript:

```
for file in `ls`  
do  
mv "$file" `cat $file | grep '<div type=.body' | sed 's/<div type=.body. *n=.'" | sed 's/^" select=\\"en et\>/'`  
done
```

Tavaline sõnade võrdlemine ei erista juhtumeid, kus lõigud on üksteisest ainult formaalselt, mitte sisuliselt erinevad (tänu kirjavahemärkide või suurtähtede kasutusele) nendest juhtumistest, mil tegemist on hoopis erinevate lõikudega. Paralleelistamise korrektsuse kontrollimisel oleks selline eristamine aga oluline. Lahenduseks on Levešteni distantse kasutamine ja teatud mittevastavuse protsendi lubamine, mille puhul võib lõike lugeda veel ühesugusteks.

Eesmärgiks oli esiteks välja selgitada, millises ulatuses on failid mõlemas korpuses ühtemoodi paralleelistatud. Ühtemoodi paralleelistused on tõenäoliselt korrektsed, samas kui erinevatest paralleelistustest on vähemalt üks vale. Teine etapp oleks erinevuste analüüs.

Ülaltoodu valguses tuli eri korpuste tekstide võrdlemiseks kasutada järgmist algoritmi:

1. Võrdle tekstide ingliskeelseid originaalosi ja leia ühisosa - kokkulangevate lõikude arv A. (Ei olnud ühtegi teksti, mille originaalversioonid mõlemas korpuses oleksid 100% kokku langenud. Kokkulangevate lõikude protsent oli vahemikus 0...99).

0-1 vastavusi võrdlemisel ignoreerisin, kuid hindamisel lugesin nad mittekorreksete vastavuste hulka.

Tekstide võrdlemiseks kasutasin Unixi utiliiti *diff*, mis teatavate parameetritega kasutades annab väljundina loendi täht-tähelt kokkulangevatest ridadest ja ridadest, mis mingis ulatuses erinevad. Neid erinevaid ridu võrdlesin omakorda Levenšteni distantse⁴³ kasutades, et eristada tegelikke paralleelistuserinevusi juhuslikest vormistuserinevustest.

2. Võrdle tekstide paralleelversioone, kus inglise lõik ja tema tõlge moodustavad ühe terviku, ja leia kokkulangevate vastavuste arv B.

Lähtetekstide ühisosa A leidmisel oli mittekokkulangevate märkide lubatud protsent 2 ja ühesuguste vastavuste B leidmisel 1 (kui lubatud protsent on A ja B leidmisel sama, siis võib juhtuda, et inglise originaalis on mõne rea erinevus failiti suurem kui lubatud protsent,

⁴³ Levenšteni distantse ehk teisenduskaugus on vähim teisendusoperatsioonide arv, mis tuleb teha selleks, et muuta üks sõne teiseks.

tõlkega koos aga väiksem, ja tekib loomuvastane olukord, et õigesti paralleelistatud lõike on rohkem kui ühesuguseid lõike lähtetekstis ehk $B > A$)

3. Leia paralleelistuse sarnasus $C=B/A$, s.t. ühesuguselt paralleelistatud lõikude ja ühesuguste originaal-lõikude suhe. C peab olema 0 ja 1 vahel (incl.).

C ei näita, kummas failis on paralleelistus parem ja kummas halvem, ta näitab ainult kahe faili kattuva osa erinevust.

C võtsin indikaatoriks, mille järgi valida, milliseid faile käsitsi läbi vaadata. Oletasin, et mingid esialgu tundmatud tegurid on paralleelistamise algoritme mõjutanud. Neid tegureid võis olla rohkem kui üks ja nende mõju eri meetoditele võis olla erineva tugevusega, seetõttu kontrollisin tekste, mis jäid C erinevatesse vahemikesse.

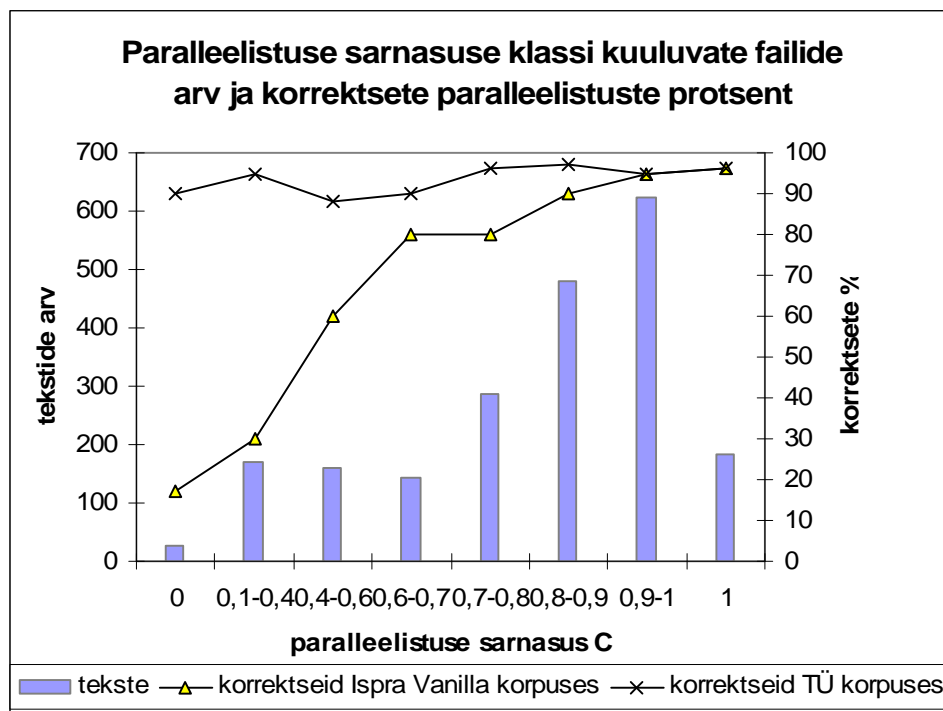
5.3.2. JRC Vanilla versiooni ja TÜ korpuse võrdlus

Jagasin failid nende paralleelistuse sarnasus C suuruse järgi 8 gruppi (vt joonist 1). Vaatasin läbi vähemalt 5% igast C vahemikku jäävatest failidest, et saada vastus küsimusele: milline on TÜ ja JRC-Acquis' korpuses korrektselt paralleelistatud lõikude protsent? Lisaks C -le on oluliseks indikaatoriks identsete ridade protsent failipaariti, mis näitab, kui suurel määral C -d antud faili puhul usaldada saab. Identsete ridade protsenti võtsin arvesse üldistuste tegemisel nende failide kohta, mis jäid läbi vaatamata.

Koostasin läbivaadatud failide põhjal tabeli korrektsete vastavuste protsendi hinnanguliste väärtustega. Iga faili korrektsete vastavuste protsendi algväärtuseks seadsin vastava failigrupi hinnangulise keskmise, välja arvatud läbivaadatud failid, mille puhul märkisin tabelisse täpse protsendi. Kuna täheldasin, et tühjad segmendid tähistavad paralleelistusvigu, siis korrigeerisin seda hinnangulist algväärtust automaatselt väiksemaks sõltuvalt tühjade segmentide protsendist antud failis. Kontrollisin hinnanguliste väärtuste keskmist failigrupiti läbivaadatud failide vastavate keskmistega, erinevuste korral korrigeerisin hinnangulisi väärtuseid vajalikus suunas.

Joonis 1 iseloomustab korrektsete vastavuste protsenti mõlemas korpuses, sõltuvalt failide paralleelistuste sarnasusest C . Näeme, et TÜ ja JRC Vanilla paralleelistuse erinevuse

korral on TÜ korpuse paralleelistus korrektsem. JRC-Acquis' korpuse alaosa korrektsete paralleelistuste protsent on üsnagi suures korrelatsioonis C väärtusega, TÜ korpuse korrektsete paralleelistuste protsenti mõjutab C vähem. Antud joonisel on esitatud keskmised statistilised väärtused, üksikute failide puhul võib siiski juhtuda, et JRC Vanilla paralleelistus on korrektsem.



Joonis 1.

Tuginedes käsitsi läbi viidud võrdlemisele võib väita, et TÜ korpuses tervikuna on korrektsete paralleelistuste protsent ca 95 ja JRC-Acquis' korpuse Vanilla paralleelistuses ca 84.

5.3.3. 0-vastavuste protsent

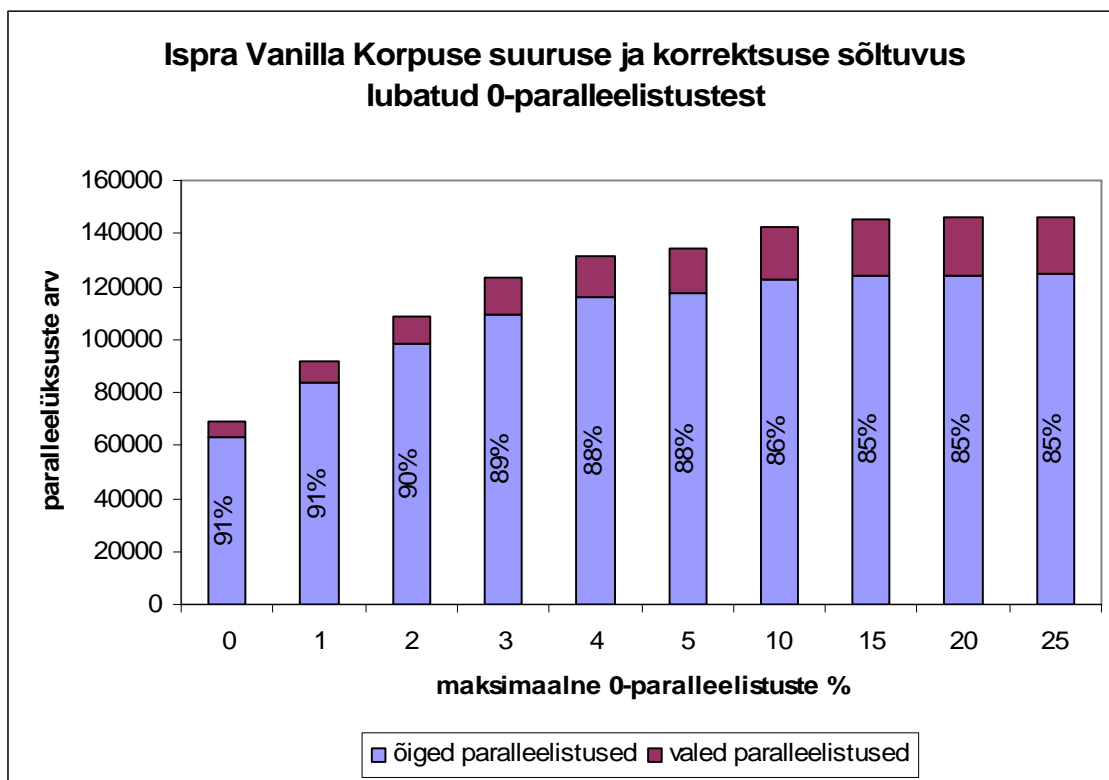
Kas on mingit viisi kindlaks teha teksti paralleelistamise korrektsust ilma teda teise korpuse tekstiga võrdlemata? Sellise tunnuse leidmine võimaldaks korpuste ühendamisel, mil vastav tekst teises korpuses puudub, vigased tekstid välja jätta.

Lähtudes paralleelistusalgoritmide tööpõhimõttest võiks oletada, et selliseks tunnuseks sobib 0-vastavuspiirkondade protsent: mida rohkem esineb juhtumeid, mil paralleelistusalgoritm pole originaalile või tõlkele vastavust leidnud, seda tõenäolisem on, et ka leitud vastavused on valed.

TÜ korpuse andmete põhjal arvatuna oli kovariatsioon tühjade ridade protsendi ja korrektsete vastavuste hinnangulise protsendi vahel 0,33 ja korrelatsioon 0,02. S.t. seos puudub. Ilmselt on siin põhjuseks TÜ korpuse paralleelistamisel kasutatud algoritm, mis enne väiksemate üksuste paralleelistamise juurde asumist eemaldas eelmisel etapil loodud 0-vastavused.

JRC-Acquis' korpuses oli kovariatsioon -53,82 ja korrelatsioon -0,42. S.t. eksisteerib keskmise tugevusega seos. Seega JRC-Acquis' korpuse puhul võib 0-vastavuste protsenti kasutada paralleelistuse kvaliteedi indikaatorina.

Joonis 2 kirjeldab JRC Vanilla korpuse (õigemini selle TÜ korpusega võrreldud osa) põhjal loodavate alamkorpuste suurusi ja korrektsusi, kui seada erineva ulatusega piiranguid 0-vastavusi sisaldavate tekstide korpusesse võtmisel. 25%-st leebemad piirangud ei avalda korrektsusele ega korpuse mahule erilist mõju, sest vastavate failide arv on väike. Usutavasti kehtivad samad proportsioonid ka JRC-Acquis' korpuse selle osa kohta, mida pole TÜ korpusega võrreldud.



Joonis 2.

5.3.4. JRC HunAligni ja TÜ korpuse võrdlus

JRC HunAligni korpuse failide paralleelsuse sarnasuse määr C TÜ korpuse vastavate failidega oli keskmiselt suurem kui JRC Vanilla ja TÜ korpuse vahel.

Juhuslike failide läbivaatamisel selgus, et kui HunAlign on leidnud vastavuse 0-1, siis antud korpuses on viga nii selles vastavuses kui ka eelmises või järgmises. 95% käsitsi kontrollitud failidest oli HunAligni paralleelistus väljaspool nimetatud 0-1 veapiirkonda täiesti korrektne. Paraku on 0-1 vastavuste hulk küllaltki suur, nii et keskmine hinnanguline korrektsete paralleel-lõikude protsent antud korpuses on 94, mis on praktiliselt võrdne TÜ korpuse 95-ga.

Erinevalt JRC Vanilla versioonist ei olnud HunAligni versioonis märgatavat seost 0-vastavuste protsendi ja sama faili teiste vastavuste korrektsuse vahel.

TÜ paralleelkorpuse veebilehel⁴⁴ korpuse koostamise protseduuri kirjelduses on muuhulgas öeldud:

„Mõnikord juhtus, et paralleelistatavate üksuste arv paralleelsetes tekstides on liiga erinev (s.t. ühes tekstis on paralleelistatavaid üksusi enam kui kaks korda rohkem kui teises), et tulemus võiks olla mõistlik. Näiteks kui ühes tekstis on ühes artiklis 1 lõik ja teises 5, siis on ilmselt tegemist kas elementide (nt tabelite või viidete) ärajäämisega ühest tekstist või kujundusliku erinevusega (nt ühes tekstis on loetelu esitatud ühel real, teises aga mitmel). Sellisel juhul jäeti antud artikli lõigud paralleeltekstist lihtsalt välja, kuid teksti ülejäänud osad võeti siiski sisse.” (TÜPK)

Selliseid üks-mitmele vastavusi, mis TÜ korpusest välja jäeti, esineb ka JRC Hunaligni versioonis. Kas need lõigud tuleks samuti TÜ korpusega liidetavast JRC-Acquis' korpuse osast välja jätta?

Sellele küsimusele vastamiseks vaatasin läbi JRC Hunaligni versioonis kümme juhuslikku 1-3 (st ühele inglisekeelsele lõigule on asetatud vastavusse kolm eestikeelset lõiku) vastavust. Kokku oli selliseid vastavusi 1565. Läbivaadatud vastavustest kaheksa olid sisuliselt korrektsed vastavused ja kaks osaliselt ebakorrektsed⁴⁵. Kümnest juhuslikust 1-4 vastavusest (kokku 187) olid neli täielikult või osaliselt ebakorrektsed. Kümnest juhuslikust 1-5 vastavusest (kokku 152) olid kõik kümme osaliselt ebakorrektsed. Mitmetel juhtudel on inglisekeelse lõiguga paralleelistatud rohkem kui 5 eestikeelset lõiku.

Kümnest juhuslikust 3-1 vastavusest (kokku 180) olid kõik kümme osaliselt ebakorrektsed. Kümne juhusliku 4-1 vastavuse (kokku 42) seas leidis üks korrektne vastavus, ülejäänud vastavused olid siiski ebakorrektsed. Kümnest juhuslikust 2-2 vastavusest (kokku 22) olid kolm korrektsed, ülejäänud ebakorrektsed. Võib oletada, et ka 1-2 vastavuste veaprotsent JRC korpuses on mõnevõrra suurem võrreldes 1-1 vastavustega.

⁴⁴ <http://www.cl.ut.ee/korpused/paralleel/index.php?lang=et>

⁴⁵ Ebakorrekseteks nimetan paralleelistusi, mille üks pool sisaldab originaallõiku või -lõike, millele paralleelistuse teises pooles ei leidu sisulist vastet.

Korrelatsiooni üks-mitmele (1-2 ja 2-1) vastavuste sageduse ja veaprotsendi vahel on tuvastatud ka araabia-hispaania-inglise paralleelcorpuse puhul (Samy jt 2006).

5.3.5. JRC-Acquis' corpuse Vanilla ja HunAligni versioonide võrdlus

JRC-Acquis' corpuse kahe erineval meetodil paralleelistatud failide võrdlemisel kasutasin varemkirjeldatud metoodikat: algul leidsin paralleelistuse sarnasuse määra C ja seejärel võrdlesin valikuliselt faile igast C vahemikust. Erinevus oli ainult selles, et lähtetekstide originaalid langesid alati 100% kokku.

Veerand võrreldavast 2000 failipaarist olid täiesti ühesugused. Ülejäänutest oli HunAligni paralleelistus alati korrektsem kui Vanilla abil tehtu. See kinnitab JRC HunAligni versiooni ja TÜ corpuse võrdlusel tehtud hinnangut JRC HunAligni versiooni ja TÜ corpuse korrektsusetaseme sarnasuse kohta.

6. Kokkuvõte

Lausete paralleelistamise meetodite võrdlemise suhtes ollakse seisukohal, et praegu pole võimalik välja tuua ühte ja parimat meetodit (Rosen 2005; Singh and Husain 2005), mis ületaks teisi mistahes korpusel. Samuti ei ole õigustatud paralleelistusmeetodi ühe korpuse peal saadud täpsuse ja saagise ülekandmine teise korpuse iseloomustamiseks, isegi kui seal on kasutatud sama meetodit (Singh, Husain 2005).

Viendas peatükis kirjeldatud võrdlus näitas veelkord, et ka juhul, kui sama korpuse peal kasutatakse sama meetodit (Vanilla), võivad tulemused olla radikaalselt erinevad, kui korpused on erinevalt normaliseeritud või meetodit rakendatakse veidi erineval moel. Samuti näitas võrdlus, et eri kohtades loodud, sama sisu, kuid erinevatest allikatest pärinevate ja erinevate meetoditega paralleelistatud korpusi on võimalik kasutada korpuste endi kvaliteedi hindamisel.

Võrdlemise tulemusel saab väita järgmist:

- JRC-Acquis' korpuse HunAligni versioon on palju korrektsem kui Vanilla versioon.
- JRC-Acquis' HunAligni versioonis 0-vastavus +/- üks paralleelne rida on suure tõenäosusega ebakorrektsed ja need võib eemaldada.
- Suure tõenäosusega on JRC -Acquis' HunAligni versioonis ebakorrektsed ka vastavused, kus ühe ingliskeelse lõiguga on paralleelistatud rohkem kui neli lõiku, ja need võib eemaldada. Ülejäänud üks-mitmele vastavuste veaprotsent vajab täpsemat määratlemist.
- Võrreldud korpused erinesid mitte ainult oma vastavuste ja tõlkeversioonide poolest, vaid ka lähtetekstide osas – polnud ühtegi ingliskeelset teksti, mis oleks formaalselt 100% kokku langenud oma ingliskeelse vastega teises korpuses.

- Tekstide automaatseks võrdlemiseks kasutasin vastavuste sarnasuse mõõtu C. C arvestamine võimaldas vähendada käsitsi läbivaatamise vajadust.
- Ankurpunktide arvestamine, olgu see siis paralleelistamisel (Hunalign) või vahepealse tulemuse filtreerimisel (TÜ versioon Vanillast) on isegi formaalselt lihtsate tekstide puhul oluline – vigade arv väheneb 10%, võrreldes ankrupunkte mittearvestava meetodiga (JRC-Acquis Vanilla).
- Leitud hinnangud korpustele võivad seletada samade korpuste peal tehtud (Fishel jt 2007) statistilise masintõlke eksperimentide tulemuste erinevusi.
- Leitud tunnused (tühjade ridade protsent failis JRC-Acquis Vanilla puhul, tühjade vastavuste olemasolu JRC-Acquis Hunaligni puhul) võimaldavad korpustele lisamisel eelistada usaldusväärsemaid vastavusi.
- Kirjeldatud võrdlusmeetodid peaksid olema rakendatavad ka muude korpuste puhul, mille jaoks eksisteerib sarnaseid, osalise kattuvusega paralleelkorpusi, nt sloveeni- inglise Euroopa seadusandluse korpused SVEZ-IJS ja JRC-Acquis' korpus (Erjavec 2006).

Viiendas peatükis kirjeldatud tööst selgub, et nii leksikoni genereerimiseks kui ka muudeks rakendusteks oleks otstarbekam JRC-Acquis' korpuse Vanilla versiooni asemel kasutada HunAligni versiooni ühendatuna TÜ paralleelkorpusega. Kuigi Fishel jt (2007) leidsid, et ka Vanilla versiooni ühendamine TÜ paralleelkorpusega aitab eesti- inglise statistilise masintõlke tulemusi mõnevõrra parandada, on siiski väga tõenäoline, et Vanilla versiooni asendamine HunAligni versiooniga oleks ka masintõlke puhul vajalik.

Lisaks TÜ paralleelkorpusele ja JRC-Acquis' paralleelkorpusele on koostatud veel üks tekstimahukas mitmekeelne paralleelkorpus – OPUS (Tiedemann, Nygaard 2004) – mis samuti sisaldab eestikeelset osa. Muuhulgas sisaldab OPUS 115 000 sõne ulatuses eestikeelset Euroopa Liidu seaduseandlikku teksti (ja eestikeelse tekstiga paralleelistatud inglisekeelset teksti), mille ühendamine TÜ paralleelkorpusega oleks samuti soovitatav. Kõnealuse OPUSi osa ühendamisel TÜ paralleelkorpusega on tõenäoliselt vähemalt osa-

liselt võimalik lähtuda eelnevas peatükis kirjeldatud kahe korpuse ühendamisele eelnenud tööst korpuste võrdlemisel ja hindamisel.

Arvesse võttes kahe paralleelistusmeetodi võrdlustulemusi (vt osa 5.3.5) võib oletada, et OPUSE osa ühendamisel TÜ paralleelkorpusega ei tohiks kasutada korpuse kodulehel⁴⁶ saadaval olevaid Vanillaga paralleelistatud inglise-eesti tekste. Selle asemel tuleks kasutada paralleelistamata tekste ja paralleelistus teostada HunAligni või mõne teise meetodi abil.⁴⁷

Lähemas tulevikus oleks vaja niisiis TÜ paralleelkorpus ühendada JRC-Acquis' korpuse HunAligni versiooniga ning OPUSE inglise-eesti seadusandlikke tekste hõlmava osaga. Probleeme tekitab seejuures samasuguste CELEXi koodidega failide olemasolu mõlemas korpuses. Tuleb välja töötada automaatne meetod, mille abil iga lõigupaari eraldi "vaadeldes" toimub teatud tunnuste abil kattuvate lõikude tuvastamine ja selliste lõigupaaride puhul kas ühe või teise korpuse lõigu eelistamine. Reeglina tuleks eelistada JRC HunAligni versiooni, kuna pärast 0-vastavuste ja vigaste üks-mitmele vastavuste eemaldamist on HunAligni paralleelistus suure tõenäosusega TÜ paralleelkorpuse paralleeljustusest korrektsem.

Eesti keelt hõlmavate mitmekeelsete sõnastike koostamine ja täiustamine on üks peamisi teid tulevikus eesti keele väljatõrjumise takistamiseks mitmetest eluvaldkondadest. Seetõttu osutub väga oluliseks ka magistritöö neljandas peatükis tutvustatud töö jätkamine praktilisel tasandil leksikograafide ja keeletehnoloogide poolt.

Esialgset tulemusi paralleelkorpustest leksikonide ekstraheerimise vallas jätavad veel kõvasti soovida – kontrollimisel osutus kõige suurema sisulise täpsusega genereeritud leksikoniks TÜ paralleelkorpuse põhjal genereeritud eesti-inglise leksikon, mille juhuslikult valitud 50 kirjeline fragmendi kontrollimine andis selle osa täpsuseks 87 %; teiste leksikonide täpsus jäi kontrollimisel eesti-inglise leksikonile suuresti alla. Siiski peaks kõik mitmekeelsete leksikonide koostamisega tegelevad inimesed mõtlema selle peale, kuidas

⁴⁶ <http://logos.uio.no/opus/EUconst.html>

⁴⁷ Samamoodi talitati OPUSE tekstide integreerimisel tšehhi-inglise korpuse (vt <http://ufal.mff.cuni.cz/czeng/>)

paremini ammutada automaatsete vahenditega leksikone rikastavat infot nii tava- kui paralleelkorpusdest.

Leksikoni genereerimise katset PWA-ga tuleks korrata loodava ühendparalleelkorpuse peal, kuna parem paralleelstuse kvaliteet tagab kindlasti korrektsema leksikoni. Võrreldes erinevaid leksikonigenererimismeetodeid tuleks leida inglise-eesti keelepaarile sobivaim ning sellest lähtuda spetsiaalselt sellele keelepaarile mõeldud leksikonigenererimistarkvara loomisel.

THE ROLE OF PARALLEL CORPORA IN COMPUTATIONAL LINGUISTICS: COMPARISON OF PARALLEL CORPORA AND GENERATION OF BILINGUAL LEXICONS FROM PARALLEL CORPORA

Summary

In addition to contrastive studies of languages or language variants, parallel/comparative corpora have many other uses both in theory and practice, while the potential of some of such uses is still awaiting discovery. One of the most interesting trends involves dictionary compilation or revision by means of extracting translation equivalents.

The present master's thesis attempts a survey of what has been done, with a view to some possible practical applications to Estonian in the future. For example, a simple lexicographic device has been outlined to enable the lexicographer to generate a list of translation equivalents by using a parallel text. Also, a report of attempts is presented to generate source material for an English-Estonian Estonian-English technical dictionary using not only the parallel corpus but also some free software, which needs little additional language resources beside the corpus material.

For the time being multilingual technical dictionaries can be compiled from parallel corpora only semi-automatically, because without intervention on the part of a human proofreader the method would yield but raw material to help lexicographers, terminologists or translation systems.

There is no software that could perform alignment and dictionary generation on the basis of the Estonian grammatical structure in correlation with the structure of some other language. Although the quality of the lexicon to be generated would certainly be improved

by preliminary morphological analysis of the parallel corpus, my present attention has been focused on language independent approaches to dictionary extraction. The word aligners UWA and LWA developed by Swedish researchers within the Plug project (Tiedemann 2002) use relatively little language-specific information, which makes them easily applicable in automatic generation of dictionaries containing Estonian material.

The master's thesis describes an attempt to develop source material for a technical dictionary by means of UWA and LWA, drawing on the English-Estonian parallel corpus of the University of Tartu and the English-Estonian subsection of the JRC-Aquis multilingual parallel corpus. One of the dictionaries so generated contains 130 865 entries and 482 571 word forms. The precision of a random sample of 50 entries turned out to be 60%. In addition the thesis provides a survey of the working principles of the used programmes, and some suggestions on how to improve UWA results with a view to an analogous device to be possibly developed for Estonian.

The availability of partially overlapping parallel corpora for a language pair opens up opportunities for automatically comparing, evaluating and improving them. I compare and evaluate the alignment quality of two English-Estonian parallel corpora that have been created independently, but contain overlapping texts. I describe how to determine the overlapping parts and find their alignment similarities that allow us to economize on manual evaluation effort. A feature is also suggested that could be used instead of comparing and manual checking to predict the alignment correctness.

While it has been common practice to compare and evaluate alignment quality when describing alignment methods, to my knowledge this is the first time when the alignments of completed parallel corpora are compared and evaluated.

It appears that corpora that have been created independently, containing essentially the same texts from independent sources, and which have been aligned with different methods, can be used for evaluating the alignment quality of the corpora themselves.

It appeared that the corpora compared were different not only in their translation versions and alignments, but in their source text parts as well. There was not a single text,

the English part of which had completely coincided with that of the corresponding one from the other corpus.

When comparing the corpora, I used the alignment similarity measure *C* and this allowed me to economize on manual evaluation.

The percentage of correctly aligned paragraphs ranged from 84% in JRC-Acquis Vanilla version to 94-95% in JRC-Acquis HunAlign version and the UT corpus.

The different levels of alignment quality may explain some of the differences in results by (Fishel et al, 2007), observed in statistical MT experiments conducted on different corpora.

The method used for determining *C* can be used also for selecting the most trustworthy alignments from the combination of two corpora.

The features that predict the quality of alignments – the proportion of 0-alignments in case of JRC Vanilla, and the mere existence of 0-alignments in case of JRC HunAlign – allow one to select the aligned units that are more trustworthy, even in the absence of a comparable text from another corpus.

The evaluation results of JRC-Acquis corpus might be transferable to other language pairs of the corpus; this assertion needs further investigation, however.

Kirjandus

- Abeillé, Anne* (ed.) 2003. Building and Using Parsed Corpora. Series: Text, Speech and Language Technology, vol. 20. Dordrecht: Kluwer, 440 p.
- Ahrenberg, Lars; Andersson, Mikael; Merkel, Magnus* 1998. A simple hybrid aligner for generating lexical correspondences in parallel texts. – Proceedings of the Thirty-Sixth Annual Meeting of the Association for Computational Linguistics and Seventeenth International Conference on Computational Linguistics. Montréal, Canada, 10–14 August 1998, 29–35.
- Ahrenberg, Lars; Andersson, Mikael; Merkel, Magnus* 2000a. A knowledge-lite approach to word alignment. –Parallel Text Processing: Alignment and Use of Parallel Corpora. Véronis, J. (ed.). Dordrecht: Kluwer, 97–116.
- Ahrenberg, Lars; Merkel, Magnus; Sågvall Hein, A.; Tiedemann, Jörg* 2000b. Evaluation of Word Alignment Systems. – Proceedings of LREC 2000, Athens/Greece, Volume III, 1255-1261.
- Blank, Ingeborg* 2000. Terminology extraction from parallel technical texts. – Parallel Text Processing: Alignment and Use of Parallel Corpora. Véronis, J. (ed.). Dordrecht: Kluwer, 237–252.
- Bojar, O.; Žabokrtský, Z.* 2006. CzEng: Czech-English Parallel Corpus, Release version 0.5. Prague Bulletin of Mathematical Linguistics, 86. (in print);
<http://ufal.mff.cuni.cz/~zabokrtsky/papers/pbml-06-OB-ZZ.pdf> (24.05.2007).
- Borin, Lars* 1998. Linguistics isn't always the answer: Word comparison in computational linguistics. – Proceedings of the 11th Nordic Conference on Computational Linguistics (NODALIDA), University of Copenhagen, Denmark, 140–151.
- Borin, Lars* 2002. ...and never the twain shall meet? – Parallel corpora, parallel worlds. Language and Computers: Studies in Practical Linguistics nr. 43. Borin, Lars (red.). Amsterdam: Rodopi, 47–59.
- Botley, Simon Philip; McEnery, Anthony Mark; Wilson, Andrew* (ed.) 2000. Multilingual Corpora in Teaching and Research. Amsterdam: Rodopi, 214 p.

- Bowker, Lynne; Pearson, Jennifer* 2002. Working with specialized language: A practical guide to using corpora. London/New York: Routledge, 242 p.
- Brew, C.; McKelvie, D.* 1996. Word-pair extraction for lexicography. – Proceedings of the Second International Conference on New Methods in Language Processing. K. Oazer; H. Somers, (ed.). Ankara: Bilkent University, 45–55; <http://citeseer.ist.psu.edu/brew96wordpair.html> (21.08.2006).
- Brown, Peter F.; Cocke, John; Della Pietra, Stephen A.; Della Pietra, Vincent J.; Jelinek, Frederick; Lafferty, John D.; Mercer, Robert L.; Roossin, Paul S.* 1990. A statistical approach to machine translation. – Computational Linguistics, 16(2), Cambridge: MIT Press, 79–85.
- Brown, Peter F.; Lai, Jennifer; Mercer, Robert L.* 1991. Aligning sentences in parallel corpora. – Proceedings of the 29th Annual Meeting of the Association for Computational Linguistics, Berkeley, CA, 169–176.
- Brown, Peter F.; Della Pietra, Stephen A.; Della Pietra, Vincent J.; Mercer, Robert L.* 1993. The mathematics of statistical machine translation: Parameter estimation. – Computational Linguistics, 19(2), Cambridge: MIT Press, 263–311.
- Brown, Ralf D.* 1997. Automated dictionary extraction for „knowledge-free” example-based translation. – Proceedings of the Seventh International Conference on Theoretical and Methodological Issues in Machine Translation (T MI-97); <http://citeseer.ist.psu.edu/brown97automated.html> (21.08.2006).
- Carl, Michael; Way, Andy* (ed.) 2003. Recent advantages in example-based machine translation. Dordrecht: Kluwer Academic Publishers, 520 p.; <http://www.cnts.ua.ac.be/~walter/papers/2004/d04.pdf> (26.05.2007).
- Charitakis, Konstantinos* 2007. Using parallel corpora to create a Greek-English dictionary with Uplug. – NODALIDA 2007 Conference Proceedings (to appear). Joakim Nivre, Heiki-Jaan Kaalep, Kadri Muischnek and Mare Koit (eds.), 212–215.
- Church, Kenneth W.; Gale, William A.; Hanks, Patrick; Hindle, Donald* 1991. Using statistics in lexical analysis. – Exploiting On-Line Resources to Build a Lexicon. Zernik, Uri (ed.). Hillsdale, NJ: Lawrence Erlbaum Associates, 115–164.
- Cmejrek, Martin; Curin, Jan; Havelka, Jiri* 2003. Treebanks in machine translation. – Proc. Of the 2nd Workshop on Treebanks and Linguistic Theories, Växjö, 209–212.

- Davis, Mark; Dunning, Ted; Ogden, Bill* 1995. Aligning noisy corpora. – Proceedings of the Seventh Conference of the European Chapter of the Association for Computational Linguistics, Belfield, Dublin: University College Dublin, 67–74.
- Davis, M. W.* 1998. On the effective use of large parallel corpora in crosslanguage text retrieval. –, Cross-Language Information Retrieval. Grefenstette, G. (ed.). Boston: Kluwer Academic Publishers, 11-22.
- Déjean, Hervé; Gaussier, Eric; Sadat, Fatia* 2002. Bilingual terminology extraction: an approach based on a multilingual thesaurus applicable to comparable corpora. – Proc. of COLING, Tapei, Taiwan, 218-224; <http://www.xrce.xerox.com/Publications/Attachments/2002-025/dejean.pdf> (21.08.2006).
- Dien, D.* 2002. Building a training corpus for word sense disambiguation in the English-to-Vietnamese Machine Translation. – Proceedings of Workshop on Machine Translation in Asia, COLING-02, Taiwan, 9/2002, 26-32; <http://acl.ldc.upenn.edu/W/W02/W02-1607.pdf> (26.05.2007).
- Ebeling, Jarle* 1998. Contrastive linguistics, translation, and parallel corpora. – *Meta*, 43:4, 602–615.
- van der Eijk, Pim* 1993. Automating the Acquisition of Bilingual Terminology. – Proceedings of the 6th Conference of the European Chapter of the Association for Computational Linguistics, Utrecht, 113–119.
- Erjavec, Tomaž* 2006. The English-Slovene ACQUIS corpus. – Proceedings of the 5th Intl. Conf. on Language Resources and Evaluations, LREC 2006. Genoa, Italy, 2138-2141; <http://nl.ijs.si/svez/bib/svez-lrec06.pdf> (26.05.2007).
- Fishel, Mark; Kaalep, Heiki-Jaan; Muischnek, Kadri* 2007. Estonian-English statistical machine translation: the first results. – Proceedings of NODALIDA. Tartu (to appear).
- Fung, Pascale* 2000. A statistical view on bilingual lexicon extraction – from parallel corpora to non-parallel corpora. – Parallel Text Processing: Alignment and Use of Parallel Corpora. Véronis, J. (ed.). Dordrecht: Kluwer, 219–236; <http://citeseer.ist.psu.edu/fung98statistical.html> (21.08.2006).
- Fung, Pascale; Church, Kenneth W.* 1994. K-vec: a new approach for aligning parallel texts. – Proceedings of the 15th International Conference on Computational Linguistics. Kyoto, Japan, 1096–1102.
- Gale, William; Church, Kenneth.* 1993. A program for aligning sentences in bilingual corpora. – Computational Linguistics 19(1), Cambridge: MIT Press, 75–102.

- Gaussier, E.; Hull, D.; Ait-Mokhtar, S.* 2000. Term alignment in use: machine-aided human translation, – Parallel Text Processing: Alignment and Use of Parallel Corpora. Véronis, J. (ed.). Dordrecht: Kluwer, 253–274.
- Geisler, Christer* 2002. Reversing a Swedish-English dictionary for the Internet. – Parallel corpora, parallel worlds. Language and Computers: Studies in Practical Linguistics nr. 43. Borin, Lars (red.). Amsterdam: Rodopi, 122–133.
- Gellerstam, Martin* 1994. Translations as a source for cross-linguistic studies. – Aijmer, Karin; Altenberg, Bengt; Johansson, Mats (ed.): Languages in contrast: papers from a symposium on text-based cross-linguistic studies, Lund: Chartwell-Bratt, 53-62.
- Genereeritud leksikonid; <http://www.teataja.ee/leksikonid.zip> (29.08.2006).
- GIZA++: Training of statistical translation models; <http://www.fjoch.com/GIZA++.html> (21.08.2006).
- Gliozzo, A.; Strapparava, C.* 2005. Cross language text categorization by acquiring multilingual domain models from comparable corpora. – Proc. of the ACL Workshop on Building and Using Parallel Texts (in conjunction of ACL-05), University of Michigan, Ann Arbor, 9–16.
- Godwin-Jones, Bob* 2001. Emerging technologies. Tools and trends in corpora use for teaching and learning. – Language Learning & Technology . Vol. 5, No. 3, September 2001, 7–12.
- Hiemstra, Djoerd* 1997. Deriving a bilingual lexicon for cross-language information retrieval. – M. Heemskerk; M. Diepenhorst (eds.), Proceedings of the fourth Groningen International Information Technology Conference for Students, 21–26.
- Hofland, Knut; Johansson, Stig* 1998. The Translation Corpus Aligner: a program for automatic alignment of parallel texts. – Stig Johansson, Signe Oksefjell (ed.), Corpora and Cross-Linguistic Research: Theory, Method and Case Studies. Amsterdam: Rodopi. 87–100.
- Hwa, Rebecca; Madnani, Nitin* 2004. The UMIACS word alignment interface; <http://www.umiacs.umd.edu/~nmadnani/alignment/foreclip.htm> (21.08.2006).
- Johansson, Stig* 2002. Towards a multilingual corpus for contrastive analysis and translation studies. – Parallel corpora, parallel worlds. Language and Computers: Studies in Practical Linguistics nr. 43. Borin, Lars (red.). Amsterdam: Rodopi, 47–59.
- Ide, Nancy; Erjavec, Tomaz; Tufiş, Dan* 2002. Sense discrimination with parallel corpora. – Proceedings of the SIGLEX Workshop on Word Sense Disambiguation: Recent Successes and

- Future Directions. ACL2002, Philadelphia, 56–60;
http://www.racai.ro/~tufis/Selected_Papers/sense-discrimination.pdf (21.08.2006).
- Index of /telri/Vanilla; <http://nl.ijs.si/telri/Vanilla/> (21.08.2006).
- Index of JRC-Acquis/alignments; <http://wt.jrc.it/lt/Acquis/JRC-Acquis.2.2/alignments/index.html>
(24.05.2007).
- Inglise-eesti-inglise sõnastik; <http://www.eki.ee/dict/inglise/> (21.08.2006).
- TÜPK = Inglise-eesti ja eesti-inglise paralleelkorpus; <http://test.cl.ut.ee/korpused/paralleel/> (24.05.2007).
- Kay, Martin* 1991. Text-translation alignment. ACH/ALLC '91: „Making Connections” Conference Handbook. Tempe, Arizona.
- Kitsnik, Mare* 2006. Keelekorpused ja võõrkeeleõpe. – Eesti Rakenduslingvistika Ühingu aastaraamat 2 (2006). Toim Helle Metslang; Margit Langemets, Tallinn: Eesti Keele Sihtasutus, 93–1007.
- Koehn, Philipp* 2002. Europarl: A multilingual corpus for evaluation of machine translation. Draft, unpublished; <http://people.csail.mit.edu/~koehn/publications/europarl.ps> (24.05.2007).
- Koehn, Philipp* 2003. Noun Phrase Translation. PhD thesis, University of Southern California; <http://www.iccs.inf.ed.ac.uk/~pkoehn/publications/thesis-readable.ps> (24.05.2007).
- Koit, Mare* 2003. Masintõlge ja kus temast kasu on? - Arvutimaailm, 2003, nr 4, 51–55;
<http://www.am.ee/6569> (21.08.2006).
- Kuhn, Jonas* 2004. Experiments in parallel-text based grammar induction. – Proceedings of the 42nd Annual Meeting of the Association for Computational Linguistics: ACL 2004, 470–477.
- Langlais, P., Simard, M., & Véronis, J.* 1998. Methods and practical issues in evaluating alignment techniques. – Joint 17th International Conference on Computational Linguistics (COLING'98) and 36th Annual Meeting of the Association for Computational Linguistics (ACL'98), Montréal, 10-14.
- Linear B. Word alignment tool; <http://demo.linearb.co.uk:8080/sandbox/start.jsp> (21.08.2006).
- Martin, J.; Johnson, H.; Farley, B.; Maclachlan A.* 2003. Aligning and using an english-inuktitut parallel corpus. – Proceedings of the HLT-NAACL Workshop on Building and Using Parallel Texts: Data Driven Machine Translation and Beyond. Edmonton, Canada, 115-118;
<http://acl.ldc.upenn.edu/W/W03/W03-0320.pdf> (26.05.2007).

- McEnery, A.M.; Piao, Scott; Xin, Xu* 2000, Parallel alignment in English and Chinese. – Multilingual Corpora in Teaching and Research. S.P. Botley, A.M. McEnery and A. Wilson (eds.). Rodopi, Amsterdam – Atlanta, 177–191.
- Melamed, Dan* 1997. A word-to-word model of translational equivalence. – Proceedings of the Thirty-Fifth Annual Meeting of the Association for Computational Linguistics and Eighth Conference of the European Chapter of the Association for Computational Linguistics. Madrid, Spain, 490–497.
- Melamed, Dan* 2001. Empirical methods for exploiting parallel texts. Cambridge, Massachusetts: MIT Press, 198 p.
- Melby, Alan K.* 2000. Sharing of translation memory databases derived from aligned parallel text. – Parallel Text Processing: Alignment and Use of Parallel Corpora. Véronis, J. (ed.). Dordrecht: Kluwer, 347–368.
- Moore Robert C.* 2002. Fast and accurate sentence alignment of bilingual corpora. – Machine Translation: From Research to Real Users (Proceedings, 5th Conference of the Association for Machine Translation in the Americas, Tiburon, California). Heidelberg, Germany: Springer-Verlag, 135-244.
- Muischnek, Kadri* 2006. Improving the quality of the statistical machine translation with linguistic preprocessing: the case of particle verbs in Estonian;
http://www.id.cbs.dk/~dh/ngslt/projects/muischnek_final_paper.doc (21.08.2006).
- Muischnek, Kadri; Orav, Heili; Kaalep, Heiki-Jaan; Õim, Haldur* 2003. Eesti keele tehnoloogilised ressursid ja vahendid. Arvutikorpused, arvutisõnastikud, keele tehnoloogiline tarkvara. Toim. U. Talvik. Tallinn: Eesti Keele Sihtasutus.
 Multext-East home page; <http://nl.ijs.si/ME/> (21.08.2006).
- Müller, Karin* 2005. Revealing phonological similarities between related languages from automatically generated parallel corpora. – Proceedings of the ACL Workshop on Building and Using Parallel Texts. ACL, 33–40 ;
<http://acl.eldoc.ub.rug.nl/mirror/W/W05/W05-0805.pdf> (26.05.2007).

- Nerbonne, John; Heeringa, Wilbert* 1997. Measuring dialect distance phonetically. – Proceedings of the third meeting of the SIGPHON at ACL, 11–18.
- Nerbonne, John* 2000. Parallel texts in computer-assisted language learning. – Parallel Text Processing: Alignment and Use of Parallel Corpora. Véronis, J. (ed.). Dordrecht: Kluwer, 354-369.
- Nivre, Joakim; De Smedt, Koenraad; Volk, Martin*. 2005. Treebanking in Northern Europe: A white paper. – Henrik Holmboe (ed.), Nordisk Sprogteknologi 2004. Arbog for Nordisk Sprogteknologisk Forskningsprogram 2000-2004. Copenhagen: Museum Tusulanums Forlag, 97–112;
<http://ling.uib.no/~desmedt/trepil/whitepaper/whitepaper-yearbook2004.pdf> (26.05.2007).
- Oakes, M.P.; McEnery, A.M.* 1998. Bilingual text alignment: an overview. – McEnery, A.M, Botley, S.P and Wilson, A. (eds), Multilingual Corpora in Teaching and Research, Amsterdam: Rodopi, 1-37.
- Och, Franz Josef; Ney, Hermann* 2000. Improved statistical alignment models. – Proc. of the 38th Annual Meeting of the Association for Computational Linguistics, Hongkong, China, 440–447.
- Resnik, Philip* 1999. Mining the web for bilingual text. – Proc. Of 37th Meeting of the ACL. Maryland, 527-534; <http://umiacs.umd.edu/~resnik/pubs/acl99.ps.gz> (26.05.2007).
- Resnik, P.; Smith, N. A.* 2003. The Web as a parallel corpus. – Computational Linguistics, 29(3), Cambridge: MIT Press, 349–380.
- Rosen, Alexandr* 2005. In search of the best method for sentence alignment in parallel texts. – Proceedings of SLOVKO 2005, the Third International Seminar on Computer Treatment of Slavic and East European Languages, VEDA. Bratislava;
<http://utkl.ff.cuni.cz/~rosen/public/slovko05.pdf> (26.05.2007).
- Salkie, Raphael* 2002. How can linguists profit from parallel corpora? – Parallel corpora, parallel worlds. Language and Computers: Studies in Practical Linguistics nr. 43. Borin, Lars (red.). Amsterdam: Rodopi, 93–109.

- Samy, D.; Sandoval, A.M.; Guirao, J.M.; Alfonseca, E.* 2006. Building a parallel multilingual corpus (Arabic-Spanish-English). – Proceedings of the 5th Intl. Conf. on Language Resources and Evaluations, LREC 2006. Genoa, Italy, 2176-2181.
- Scannell, K. P.* 2003. Automatic thesaurus generation for minority languages: an Irish example. – Actes de la 10e conférence TALN à Batz-sur-Mer, volume 2, 203–212.
- Singh, A. K.; Husain, S.* 2005. Comparison, selection and use of sentence alignment algorithms for new language pairs. – Proceedings of the ACL Workshop on Building and Using Parallel Texts, Ann Arbor, Michigan, Association for Computational Linguistics, 99–106.
- Smadja, F.; McKeown, K. R.; Hatzivassiloglou, V.* 1996. Translating collocations for bilingual lexicons: A statistical approach. – Computational Linguistics, 22(1), Cambridge: MIT Press, 1–38.
- Software for word alignment; <http://www.cse.unt.edu/~rada/wa/#softwareWA> (21.08.2006).
- Sperberg-McQueen, C. M.; Burnard, Lou* (eds.) 2004. Guidelines for Electronic Text Encoding and Interchange. XML-compatible edition; <http://www.tei-c.org/P4X/index.html> (28.05.2007).
- Stahl, Peter* 2002. Building and processing a multilingual corpus of parallel texts. – Parallel corpora, parallel worlds. Language and Computers: Studies in Practical Linguistics nr. 43. Borin, Lars (red.). Amsterdam: Rodopi, 47–59.
- Steinberger, Ralf; Pouliquen, Bruno; Widiger, Anna; Ignat, Camelia; Erjavec, Tomaž; Tufiş, Dan; Varga, Dániel* 2006. The JRC-Acquis: A multilingual aligned parallel corpus with 20+ languages. – Proceedings of the 5th International Conference on Language Resources and Evaluation (LREC'2006). Genoa, Italy; http://langtech.jrc.it/Documents/0605_LREC_JRC-Acquis_Steinberger-et-al.pdf (24.05.2007).
- Sumita, Eiichihiro* 2003. An Example-Based Machine Translation System Using DP-Matching Between Word Sequences. – Recent Advantages in Example-Based Machine Translation. Michael Carl and Andy Way (eds.). Dordrecht: Kluwer Academic Publishers, 189-209.

- Sågvall Hein, Anna* 2002. The PLUG-project. Parallel corpora in Linköping, Uppsala, Göteborg: aims and achievements. – Parallel corpora, parallel worlds. Language and Computers: Studies in Practical Linguistics nr. 43. Borin, Lars (red.). Amsterdam: Rodopi, 61–78.
- The JRC-Acquis multilingual parallel corpus; <http://langtech.jrc.it/JRC-Acquis.html> (21.08.2006).
- The PLUG Word Aligner – PWA; <http://stp.ling.uu.se/plug/pwa/index.html> (21.08.2006).
- Tiedemann, Jörg* 1997. Automatical lexicon extraction from aligned bilingual corpora. Diploma thesis, University ‘Otto-von-Guericke’, Magdeburg: Department of Computer Science; <http://stp.ling.uu.se/~joerg/diplom/diplom.html> (26.05.2007).
- Tiedemann, Jörg* 2002. Uplug – a modular corpus tool for parallel corpora. – Parallel corpora, parallel worlds. Language and Computers: Studies in Practical Linguistics nr. 43. Borin, Lars (red.). Amsterdam: Rodopi, 181-197.
- Tiedemann, Jörg* 2003. Recycling translations. Extraction of lexical data from parallel corpora and their application in natural language processing; http://stp.ling.uu.se/~joerg/phd/html/thesis_html.html (21.08.2006).
- Tiedemann, J.; Nygaard, L.* 2004. The OPUS corpus – parallel and free. – Proceedings of the Fourth International Conference on Language Resources and Evaluation (LREC-04), Lisbon, Portugal, 1183-1186; http://stp.ling.uu.se/~joerg/paper/opus_lrec04.pdf (26.05.2007).
- Trosterud, T.* 2002. Parallel corpora as tools for investigating and developing minority languages. – Parallel corpora, parallel worlds. Language and Computers: Studies in Practical Linguistics nr. 43. Borin, Lars (red.). Amsterdam: Rodopi, 111–122.
- Tufiş, Dan; Barbu, Ana-Maria* 2001. Automatic construction of translation lexicons. – Advances in Automation, Multimedia and Modern Computer Science, A Series of Reference Books and Textbooks in Electrical and Computer Engineering. V. V. Kluew; C. E. D'Attellis; N.E. Mastorakis (eds). WSEAS Press, 156–172.
- Uchimoto, Kiyotaka; Zhang, Yujie; Sudo, Kiyoshi; Murata, Masaki, Sekine, Satoshi; Isahara, Hitoshi* 2004. Multilingual aligned parallel treebank corpus reflecting contextual information and its applications. – Proceedings of the Multilingual Language Resources (MLR2004); 2004, Geneva,

- Switzerland;
<http://cs.nyu.edu/~sekine/papers/mlr04-uchimoto.pdf> (26.05.2007).
- Varga, D., L. Németh, P. Halácsy, A. Kornai, V. Trón* (2005): Parallel corpora for medium density languages. – Proceedings of Recent Advances in Natural Language Processing (RANLP-05), Borovets, Bulgaria; <http://www.metacarta.com/docs/ranlp05parallel.pdf> (26.05.2007).
- Volk, Martin; Samuelsson, Yvonne* 2004. Bootstrapping parallel treebanks. – 20th International Conference on Computational Linguistics, Geneva, 29 August 2004; <http://www.coli.uni-saarland.de/conf/linc-04/volk.pdf> (26.05.2007).
- Volk, Martin* 2005. The role of the Web in machine translation. – M. Langemets, P. Penjam (eds.), Proc. Of the 2nd Baltic Conference on Human Language Technologies, Tallinn, 79-85.
- Véronis, Jean; Langlais, Philippe* 2000. Evaluation of parallel text alignment systems: The ARCADE project. – Parallel Text Processing: Alignment and Use of Parallel Corpora. Véronis, J. (ed.). Dordrecht: Kluwer, 369–388.
- Veskis, Kaarel* 2007. Kakskeelsete leksikonide genereerimine paralleelkorpuse baasil. – Eesti Rakenduslingvistika Ühingu aastaraamat 3. Toim. Helle Metslang, Margit Langemets, Maria-Maren Sepper. Tallinn: Eesti Rakenduslingvistika Ühing, lk. 355–372.
- Widdowson, H. G.* 1990. Aspects of language teaching, Oxford, Oxford University Press, 213 p.

Lisa 1. Fragment UWA-ga TÜ paralleelkorpusest ekstraheeritud leksikonist

...	{
{viljapuuistandike}	2X:infertile
{	}
2X:fruit tree plantations	{viljatükk}
}	{
{viljastamine}	1X:blemished
{	}
1X:fertilisation	{viljelemine}
}	{
{viljastamise}	1X:eliminated
{	}
1X:violation	{viljelty}
}	{
{viljastanud}	1X:viljelty
{	}
2X:fertilizing	{viljelusest}
}	{
{viljastumiseks}	1X:permanent abandonment
{	}
1X:fertilization	{vilju}
}	{
{viljatu}	3X:units
{	}
3X:infertile	{vilken}
}	{
{viljatut}	6X:vilken
	}

{vilkingid}	{
{	1X:villes déterminées
5X:wilkingis	par
}	}
{vilkingite}	{villija}
{	{
1X:wilkingis	6X:bottler
}	}
{vilkudes}	{villimine}
{	{
1X:slow-flashing	8X:bottling
}	1X:racking
{vilkur}	}
{	{vin}
6X:flashing lamps	{
}	3X:vin
{vill}	}
{	{vingt-neuf}
9X:wool	{
}	2X:vingt-neuf
{villa}	}
{	{vinho}
1X:villa	{
}	2X:vinho
{villane}	}
{	{vinifera}
1X:woollen	{
}	1X:vinifera
{villarasvaalkoholid}	}
{	{vinifera L}
1X:cetostearyl	{
}	2X:vinifera L
{villasest}	}
{	{vinification}
1X:woollen	{
}	1X:vinification
{villes déterminées par}	}

{vinificazione}	{
{	IX:violet/purple
IX:vinificazione	}
}	{violetne/roheline}
{vinkamiin}	{
{	IX:violet/green
IX:vincamine	}
}	{violetne/roheline spargel}
{vino}	{
{	2X:violet/green asparagus
6X:vino	}
}	{violetse}
{vinorelbiini}	{
{	3X:violet
IX:vinorelbine	}
}	{violetse spargli}
{vinte}	{
{	2X:violet asparagus
5X:vinte	}
}	{violetse/rohelise}
{vinte e}	{
{	IX:violet/green
IX:ventidue	}
}	{violetse/rohelise spargli}
{vintraudset}	{
{	2X:violet/green asparagus
IX:a rifled barrel	}
}	{violetsed/purpursed}
{violetne}	{
{	IX:violet/purple tops
4X:violet	}
}	{violetset}
{violetne spargel}	{
{	IX:violet
IX:violet asparagus	}
}	{violetset sparglit}
{violetne/purpurne}	{

<i>IX:violet asparagus</i>	}
}	<i>{virginiamütsiiniresistentne}</i>
<i>{violetsete/purpursete}</i>	{
{	<i>IX:virginiamycin-resistant</i>
<i>IX:violet/purple</i>	}
}	<i>{virginiamütsiiniresistentsed}</i>
<i>{violetpunase}</i>	{
{	<i>IX:virginiamycin-resistant</i>
<i>IX:violet-red</i>	}
}	<i>{virginiana}</i>
<i>{virde}</i>	{
{	<i>2X:virginiana</i>
<i>2X:provided</i>	}
}	<i>{virksomhed}</i>
<i>{virens}</i>	{
{	<i>2X:virksomhed</i>
<i>2X:virens</i>	}
}	<i>{virksomheder}</i>
<i>{virgaurea}</i>	{
{	<i>IX:virksomheder</i>
<i>IX:virgaurea</i>	}
}	<i>{virstatavate}</i>
<i>{virgem}</i>	{
{	<i>IX:stackable</i>
<i>IX:virgem</i>	}
}	<i>{virstatavates}</i>
<i>{virgen}</i>	{
{	<i>5X:stackable</i>
<i>IX:virgen</i>	}
}	<i>{viroidid}</i>
<i>{virginiamütsiin}</i>	{
{	<i>IX:viroids</i>
<i>IX:virginiamycin</i>	}
}	<i>{viroididele}</i>
<i>{virginiamütsiini}</i>	{
{	<i>IX:viroids</i>
<i>4X:virginiamycin</i>	}

{virola}	{virtsa}
{	{
1X:virola	2X:liquid manure
}	}
{viroloogiainstituut}	{virtuaalne kollokatsioon}
{	{
1X:Virology	1X:virtual collocation
}	}
{virooloogiline}	{virtuaalse kollokatsiooni}
{	{
1X:virological	1X:virtual collocation
}	}
{virooloogiliseks}	{virtud}
{	{
2X:virological	2X:virtud
}	}
{virooloogiliste}	{virtud del}
{	{
2X:virological	1X:virtud del
}	}
{virsiku}	{visés}
{	{
1X:peel	1X:visés
}	}
{virsikud}	{visatakse}
{	{
5X:Peaches	7X:Discard
13X:peaches	2X:discard
}	}
{virsikute}	{viscum}
{	{
7X:peaches	1X:viscum
}	}
{virsikutest}	{visiidid}
{	{
1X:the fruit nectars	3X:visits
}	}

{visiiditasu}	}
{	{vismutsubgallaat}
IX:visit fee	{
IX:visit fees	IX:bismuth subgallate
}	}
{visiirid}	{vismutsubkarbonaat}
{	{
IX:visors	IX:bismuth subcarbonate
}	}
{visiite}	{vismutsubnitraat}
{	{
IX:visits	IX:bismuth subnitrate
}	}
{visiitkaardid}	{vismutsubnitraat ja}
{	{
IX:visiting cards	IX:bismuth subnitrate
}	}
{viskekeha}	{vismutsubsalitsülaat}
{	{
IX:kinetic	IX:bismuth subsalicylate
}	}
{viski}	{visuaal-}
{	{
7X:whiskey	2X:visual
6X:whisky	}
}	{visuaallennuks}
{viskoosist}	{
{	IX:visual flight
IX:viscose	}
}	{visuaallähenemise}
{viskoosne}	{
{	IX:visual approach exercise
3X:viscous	}
}	{visuaalne}
{viskoosus}	{
{	IX:Visual
2X:viscosity	IX:the visual

9X:visual	}
}	{vitamiin}
{visuaalse}	{
{	4X:Vitamin
11X:visual	4X:vitamin
}	}
{visuaalsed}	{vitamiin B1}
{	{
1X:visual	1X:vitamin B1
}	}
{visuaalseks}	{vitamiin B2}
{	{
1X:visual	1X:vitamin B2
}	}
{visuaalsel}	{vitamiin B3}
{	{
2X:visual	1X:vitamin B3
}	}
{visuaalselt}	{vitamiin B6}
{	{
6X:visually	1X:vitamin B6
}	}
{visuaalset}	{vitamiin C}
{	{
3X:visual	1X:vitamin C
}	}
{visuaalseteks}	{vitamiin H}
{	{
1X:visual checks	1X:vitamin H
}	}
{visuaalsüsteem}	{vitamiine}
{	{
1X:visual system	1X:vitamin
}	12X:vitamins
{visvaart}	}
{	{vitamiini}
2X:visvaart	{

<i>1X:vegetable juice</i>	<i>2X:showcases</i>
<i>2X:vegetable juices</i>	}
<i>1X:vitamin D</i>	{ <i>vitro</i> }
}	{
{ <i>vitamiini-</i> }	<i>6X:vitro</i>
{	}
<i>2X:vitamin</i>	{ <i>vivo</i> }
}	{
{ <i>vitamiinid</i> }	<i>8X:vivo</i>
{	}
<i>1X:Vitamins</i>	{ <i>vl39</i> }
<i>15X:vitamins</i>	{
}	<i>1X:ISO/IEC</i>
{ <i>vitamiinide</i> }	}
{	{ <i>vleesras</i> }
<i>1X:the vitamin</i>	{
<i>6X:vitamin</i>	<i>1X:vleesras</i>
<i>27X:vitamins</i>	}
}	{ <i>voce</i> }
{ <i>vitamiinidest</i> }	{
{	<i>1X:voce</i>
<i>1X:vitamins</i>	}
}	{ <i>voeder</i> }
{ <i>vitaminiseerimine</i> }	{
{	<i>2X:voeder</i>
<i>1X:vitaminization</i>	}
}	{ <i>voeding</i> }
{ <i>vitelli</i> }	{
{	<i>2X:voeding</i>
<i>3X:vitelli</i>	}
}	{ <i>voedselhulp</i> }
{ <i>vitreous</i> }	{
{	<i>1X:voedselhulp</i>
<i>2X:vitreous</i>	}
}	{ <i>voi</i> }
{ <i>vitriinides</i> }	{
{	<i>11X:voi</i>

}	{voldaan}
{voimassa}	{
{	1X:voldaan
8X:voimassa	}
}	{voldi}
{voimassa viisi}	{
{	3X:volts
1X:voimassa viisi	}
}	{voldik}
{voimassaolon}	{
{	2X:folder
2X:voimassaolon	}
}	{voldise}
{voiöljy}	{
{	1X:a CV product
11X:voioeljy	1X:voltage
}	}
{voiöljyksi ja}	{volditud}
{	{
1X:voioeljyksi ja	1X:concertinafolded
}	}
{vol}	{voldoen}
{	{
5X:vol	1X:voldoen
}	}
{volatiilsus-}	{volgende}
{	{
1X:volatility	1X:volgende
}	}
{volatiilsuse}	{volgnummer}
{	{
1X:volatility	1X:volgnummer
}	}
{volatiilsusega}	{volikiri}
{	{
1X:volatility	1X:Authorisation document
}	1X:Authorisation documents

<i>29X:authorisation</i>	<i>2X:authorise a joint</i>
<i>60X:authorisation document</i>	}
<i>1X:authorisation documents</i>	{ <i>volitatu</i> }
}	{
{ <i>volikogu</i> }	<i>2X:delegator</i>
{	}
<i>6X:Council</i>	{ <i>volitatud</i> }
<i>1X:Councils</i>	{
<i>259X:council</i>	<i>6X:Approved</i>
<i>23X:councils</i>	<i>18X:Authorised</i>
}	<i>1X:Obligation</i>
{ <i>volikogu esimehe valimiseni</i> }	<i>20X:agencies</i>
{	<i>29X:agency</i>
<i>1X:volikogu esimehe valimiseni</i>	<i>193X:approved</i>
}	<i>772X:authorised</i>
{ <i>volikogu liikmeks</i> }	<i>1X:invalidate</i>
{	<i>2X:invalidated</i>
<i>1X:unelected independent</i>	<i>1X:negotiated</i>
}	<i>1X:obligated</i>
{ <i>volinike</i> }	<i>6X:obligation</i>
{	<i>4X:obligations</i>
<i>4X:proxies</i>	<i>3X:political</i>
}	<i>1X:violated</i>
{ <i>volitada teisi</i> }	<i>1X:violates</i>
{	<i>9X:violation</i>
<i>2X:authorise other persons</i>	<i>1X:violations</i>
}	}
{ <i>volitajad</i> }	{ <i>volitatud ametiisiku juuresolekul</i> }
{	{
<i>1X:delegators</i>	<i>1X:data protection supervision</i>
}	}
{ <i>volitamise</i> }	{ <i>volitatud importijad</i> }
{	{
<i>1X:the limits</i>	<i>1X:designated importers</i>
}	}
{ <i>volitamiseks</i> }	{ <i>volitatud isikud koha</i> }
{	{

2X:projects receiving financial	62X:powers
}	}
{volitatus politseiasutused}	{volitusnormita}
{	{
1X:police authorities authorised	1X:a provision delegating
}	}
{volitus}	{volitust}
{	{
2X:Authorisation	1X:limits
36X:authorisation	}
18X:mandate	{volituste}
1X:mandates	{
}	3X:Committee
{volituse}	10X:committee
{	1X:committees
1X:clause	1X:situated
1X:course	1X:virtue
3X:limits	}
6X:validity	{volituste lõppemise}
}	{
{volitused}	7X:termination
{	}
3X:Powers	{volledig}
3X:appointed	{
6X:authorised	2X:volledig
5X:entrusted	}
95X:powers	{volti}
1X:unlimited	{
}	1X:volts
{volitused kestavad uue}	}
{	{voltmeetrite}
3X:a court-appointed member	{
}	1X:voltmeters
{volitusi}	}
{	{voluntary}
1X:Powers	{
1X:complaints	5X:voluntary

}	{vooluallikate}
{vom}	{
{	2X:power supplies
IX:vom	}
}	{vooluhulk}
{vom Antragssteller beantragter}	{
{	IX:unpermitted
IX:vom Antragssteller beantragter	}
}	{vooluimpulss}
{von}	{
{	IX:impulse
47X:von	}
}	{voolujaotusega}
{von Artikel}	{
{	IX:divider
IX:von Artikel	}
}	{voolukulu}
{von Lebensmittelzubereitungen}	{
{	10X:flowrate
IX:von Lebensmittelzubereitungen	2X:flowrates
}	}
{von Maschinenanlagen}	{voolukulu Q}
{	{
IX:von Maschinenanlagen	IX:the gas meter
}	}
{vooderdatud}	{vooluringide}
{	{
2X:lined	IX:design drawings
}	2X:sub-assemblies
{voodipesu}	}
{	{voolutatakse}
3X:bed linen	{
}	IX:flushed
{voolavad}	}
{	{vooluvee}
IX:free-flowing	{
}	IX:barring

}	{voorbehoud dat}
{voor}	{
{	IX:voorbehoud dat
31X:voor	}
}	{voorbehoud van bijzondere}
{voor Zuivel}	{
{	IX:voorbehoud van bijzondere
IX:voor Zuivel	}
}	{voorru}
{voor gehandicapten}	{
{	19X:Round
IX:voor gehandicapten	2X:round
}	IX:voorru
{voor het}	}
{	{voorru tulemusi ühendava}
IX:voor het	{
}	2X:Embodying
{voor ten}	IX:Embodying the Results
{	}
2X:voor ten	{voorru}
}	{
{voor uitvoer}	2X:course
{	IX:offers
IX:voor uitvoer	}
}	{voorwaarden}
{vooraf}	{
{	IX:voorwaarden
IX:vooraf	}
}	{voorzien}
{vooraf van}	{
{	IX:voorzien
IX:vooraf van	}
}	{vorbehaltlich eines}
{voorbehoud}	{
{	IX:vorbehaltlich eines
IX:voorbehoud	}
}	{vorgesehene ermaessigte Zollsätze}

{	166X:format
1X:vorgesehene ermaessigte Zollsaeetze	1X:formats
}	11X:forms
{vorm}	3X:origin
{	3X:permit
43X:Form	}
9X:Format	{vormi CO}
4X:Forms	{
12X:for	1X:Form CO
113X:form	}
128X:format	{vormi V}
1X:formats	{
6X:forms	1X:Form V
1X:status	}
3X:work	{vormi-}
}	{
{vorm(id)}	4X:Format
{	1X:form
2X:form(s)	10X:format
}	}
{vorme}	{vormid}
{	{
2X:Forms	1X:Board
3X:form	2X:Form
22X:forms	8X:Forms
1X:more	4X:form
}	14X:format
{vormi}	16X:formats
{	47X:forms
2X:FORM	1X:origin
24X:Form	2X:permit
1X:Formal	1X:permits
10X:Format	2X:voting
1X:Forms	}
7X:Origin	{vormide}
182X:form	{
15X:formal	2X:invoice

}	}
{ <i>vormides</i> }	{ <i>vormirietus</i> }
{	{
7 <i>X</i> : <i>forms</i>	2 <i>X</i> : <i>Uniform</i>
4 <i>X</i> : <i>place particularly</i>	3 <i>X</i> : <i>uniform</i>
}	1 <i>X</i> : <i>uniforms</i>
{ <i>vormiga</i> }	}
{	{ <i>vormis</i> }
2 <i>X</i> : <i>Origin</i>	{
}	4 <i>X</i> : <i>Form</i>
{ <i>vormil</i> }	3 <i>X</i> : <i>a form</i>
{	182 <i>X</i> : <i>form</i>
2 <i>X</i> : <i>a form</i>	2 <i>X</i> : <i>formal</i>
16 <i>X</i> : <i>form</i>	21 <i>X</i> : <i>format</i>
2 <i>X</i> : <i>forms</i>	21 <i>X</i> : <i>forms</i>
}	1 <i>X</i> : <i>visits</i>
{ <i>vormile</i> }	}
{	{ <i>vormis X</i> }
1 <i>X</i> : <i>Article</i>	{
1 <i>X</i> : <i>form set</i>	1 <i>X</i> : <i>Form X</i>
1 <i>X</i> : <i>variety</i>	}
}	{ <i>vormis teabesertifikaadi INF</i> }
{ <i>vormilise</i> }	{
{	1 <i>X</i> : <i>information certificate INF</i>
2 <i>X</i> : <i>a formal arrangement</i>	}
}	...
{ <i>vormimiseks</i> }	
{	
2 <i>X</i> : <i>moulding</i>	
}	
{ <i>vorming</i> }	
{	
11 <i>X</i> : <i>format</i>	
}	
{ <i>vorminõuete</i> }	
{	
5 <i>X</i> : <i>formalities</i>	

Lisa 2. Fragment LWA-ga TÜ paralleelkorpusest ekstraheeritud leksikonist

...		<i>kaitse</i>	<i>kaitse</i>
<i>kahe aasta</i>	<i>biennial</i>	<i>kaitse</i>	<i>defence</i>
<i>kahe aasta jooksul</i>	<i>two</i>	<i>kaitse kohta</i>	<i>protection</i>
<i>kahekolmandikulise häälteenamusega</i>		<i>kaitse kohta</i>	<i>community</i>
<i>two-thirds</i>		<i>kaitseaspektide</i>	<i>security</i>
<i>kaheksanda</i>	<i>eighth</i>	<i>kaitsega</i>	<i>environment</i>
<i>kahekümnendal</i>	<i>twentieth</i>	<i>kaitseidentiteedi</i>	<i>esdi</i>
<i>kahelda</i>	<i>doubt</i>	<i>kaitseks</i>	<i>by</i>
<i>kahepoolse</i>	<i>bilateral</i>	<i>kaitseks</i>	<i>protect</i>
<i>kahepoolseid</i>	<i>concluded</i>	<i>kaitsemeetmeid</i>	<i>stringent</i>
<i>kahepoolsetes</i>	<i>double-taxation</i>	<i>kaitsepoliitika</i>	<i>defence policy</i>
<i>kahes</i>	<i>two</i>	<i>kaitsepoliitikaga</i>	<i>framework</i>
<i>kahest</i>	<i>two</i>	<i>kaitsepoliitilise tähendusega</i>	
<i>kahjum</i>	<i>taxable</i>	<i>implications</i>	
<i>kahjustab</i>	<i>kahjustab</i>	<i>kaitsepoliitilise tähendusega</i>	
<i>kahjustada</i>	<i>effectiveness</i>	<i>implications.</i>	
<i>kahjustamata</i>	<i>clarifications</i>	<i>kaitsepoliitilise tähendusega otsuste</i>	
<i>kahjustavate</i>	<i>affecting</i>	<i>decisions</i>	
<i>kaht</i>	<i>two</i>	<i>kaitset</i>	<i>otherwise</i>
<i>kahte</i>	<i>two</i>	<i>kaitsevad</i>	<i>defend</i>
<i>kahtlemise</i>	<i>doubting</i>	<i>kaitsmiseks</i>	<i>internal</i>
<i>kahtlust</i>	<i>doubt</i>	<i>kaitsta</i>	<i>protect</i>
<i>kahvaturoheline</i>	<i>pale</i>	<i>kaitsta</i>	<i>protected</i>
<i>kaimanisaared</i>	<i>cayman islands</i>	<i>kaitsta</i>	<i>safeguard</i>
<i>kaitse</i>	<i>protection</i>	<i>kaitstuse</i>	<i>protection</i>

<i>kajastamiseks</i>	<i>reflect</i>	<i>kapitali</i>	<i>capital</i>
<i>kaks</i>	<i>two</i>	<i>kapitaliosa</i>	<i>share</i>
<i>kaks</i>	<i>kaks</i>	<i>kapitalis</i>	<i>capital</i>
<i>kalanduspoliitika</i>	<i>policy</i>	<i>kapitalist</i>	<i>capital</i>
<i>kalatoodete</i>	<i>fisheries</i>	<i>kapitaliturul</i>	<i>capital</i>
<i>kalatoodetelt</i>	<i>islands</i>	<i>karistusmaksed</i>	<i>periodic penalty payments</i>
<i>kalendriaasta</i>	<i>calendar year</i>		
<i>kalendriaastaks</i>	<i>calendar</i>	<i>karjamaapasture</i>	
<i>kaltsiumfluoriidi</i>	<i>fluorine</i>	<i>karusnahad</i>	<i>pelts</i>
<i>kammitud</i>	<i>carded</i>	<i>karusnahku</i>	<i>annex</i>
<i>kanaari saarte</i>	<i>canary</i>	<i>kas</i>	<i>governors</i>
<i>kanaari saartele</i>	<i>customs</i>	<i>kas kogu</i>	<i>auditors</i>
<i>kanada</i>	<i>canada</i>	<i>kasu</i>	<i>benefits</i>
<i>kanadast</i>	<i>canada</i>	<i>kasu</i>	<i>benefit</i>
<i>kanded</i>	<i>reflect</i>	<i>kasuks</i>	<i>favour</i>
<i>kandmiseregistratsioon</i>		<i>kasum</i>	<i>hands</i>
<i>kandmist activities</i>		<i>kasumi</i>	<i>profits</i>
<i>kantakse</i>	<i>registered</i>	<i>kasumit</i>	<i>profits distributed</i>
<i>kantud</i>	<i>registered</i>	<i>kasumit</i>	<i>profits</i>
<i>kantud</i>	<i>listed</i>	<i>kasumitaotluseta</i>	<i>intermediate</i>
<i>kantud</i>	<i>list</i>	<i>kasutab</i>	<i>exercise</i>
<i>kantud</i>	<i>included</i>	<i>kasutab</i>	<i>exercises</i>
<i>kaotamise</i>	<i>penalty</i>	<i>kasutada</i>	<i>used</i>
<i>kaotamise kohta</i>	<i>june</i>	<i>kasutada</i>	<i>use</i>
<i>kaotatakse</i>	<i>abolished</i>	<i>kasutades</i>	<i>using</i>
<i>kaotatud</i>	<i>lost</i>	<i>kasutajate</i>	<i>users</i>
<i>kapital</i>	<i>capital</i>	<i>kasutamine</i>	<i>use</i>

<i>kasutamine</i>	<i>accordance</i>	<i>katta</i>	<i>rise</i>
<i>kasutamise deklaratsiooni</i>	<i>similar</i>	<i>katteleht</i>	<i>binder</i>
<i>kasutamise</i>	<i>use</i>	<i>kauba</i>	<i>goods</i>
<i>kasutamise</i>	<i>regards</i>	<i>kaubad</i>	<i>goods</i>
<i>kasutamisel</i>	<i>formulation</i>	<i>kaubandus</i>	<i>sustainable</i>
<i>kasutamist</i>	<i>use</i>	<i>kaubandus</i>	<i>sub-committee</i>
<i>kasutatakse</i>	<i>used</i>	<i>kaubandus</i>	<i>aspektide</i> <i>aspects</i>
<i>kasutatava</i>	<i>fine-cut</i>	<i>kaubanduse</i>	<i>trade</i>
<i>kasutatavad</i>	<i>used</i>	<i>kaubanduseeskirjade</i>	<i>regulations</i>
<i>kasutatavate</i>	<i>used</i>	<i>kaubandusele</i>	<i>trade</i>
<i>kasutuseks</i>	<i>national</i>	<i>kaubandusettevõtjate</i>	<i>traders</i>
<i>kasutusele</i>	<i>resort</i>	<i>kaubanduskokkuleppe</i>	<i>tariffs</i>
<i>kasutusotstarbest</i>	<i>exemptions</i>	<i>kaubanduslepingu</i>	<i>multilateral</i>
<i>kasvatajale</i>	<i>breeder</i>	<i>trade</i>	<i>agreement</i>
<i>kasvatatud</i>	<i>animals</i>	<i>kaubanduslepingu</i>	<i>plurilateral</i>
<i>kasvatatud</i>	<i>animal</i>	<i>kaubanduslepingud</i>	<i>agreement</i>
<i>kasvatatud</i>	<i>birds</i>	<i>kaubanduslepingud</i>	<i>agreements</i>
<i>kasvatatud</i>	<i>captivity</i>	<i>kaubanduslepingus</i>	<i>multilateral</i>
<i>kasvava</i>	<i>growing</i>	<i>trade</i>	<i>agreement</i>
<i>kasvu</i>	<i>growth</i>	<i>kaubanduslepingus</i>	<i>agreement</i>
<i>kasvupiirkonnas</i>	<i>situated</i>	<i>kaubanduslepingute</i>	<i>trade</i>
<i>katavad</i>	<i>communities</i>	<i>kaubanduslepingute</i>	<i>multilateral</i>
<i>kategooria</i>	<i>category</i>	<i>kaubanduslepingutega</i>	<i>covered</i>
<i>kategooria ainetet</i>	<i>transactions</i>	<i>kaubanduslepingutes</i>	<i>multilateral</i>
<i>katkematu</i>	<i>years</i>	<i>kaubandusliku</i>	<i>commercial</i>
<i>katkestanud</i>	<i>started</i>	<i>kaubanduslääbirääkimiste</i>	<i>trade</i>
<i>katmiseks</i>	<i>cover</i>	<i>negotiations</i>	

<i>kaubandusl�bir��kimiste uruguay vooru round</i>	<i>kausside outside</i>
<i>kaubandusmeetmete trade measures</i>	<i>kava scheme</i>
<i>kaubandusorganisatsioon world</i>	<i>kava plan</i>
<i>kaubanduspoliitika trade policy review</i>	<i>kavade scheme</i>
<i>kaubanduspoliitika trading</i>	<i>kavade schemes</i>
<i>kaubandussektor diminished</i>	<i>kavandamisel framing</i>
<i>kaubandust other</i>	<i>kavandamisele framing</i>
<i>kaubandust environmental measures</i>	<i>kavandatav proposed</i>
<i>kaubandustegevus commercial presence</i>	<i>kavandatavad envisaged</i>
<i>kaubandustegevuse commercial presence</i>	<i>kavandatud designed</i>
<i>kaubandust�kete barriers</i>	<i>kavandatud enter</i>
<i>kaubandust�kkeid barriers</i>	<i>kavatseb important</i>
<i>kaubandusvoogude flows</i>	<i>kavatsevad intend</i>
<i>kaubasaajale consignee</i>	<i>kavatsusest intention</i>
<i>kaubasaatjale consignor</i>	<i>keda assisted</i>
<i>kaubavahetuse goods</i>	<i>keda whom</i>
<i>kaudselt indirectly</i>	<i>keelab prohibit</i>
<i>kaudu hendamise</i>	<i>keelab prohibits</i>
<i>kaup goods</i>	<i>keelamine prohibition</i>
<i>kaupa goods</i>	<i>keelamistprohibition</i>
<i>kaupade goods</i>	<i>keelata prohibit</i>
<i>kaupade arrangements</i>	<i>keelata prohibited</i>
<i>kauplemise trade</i>	<i>keelatakse prohibited</i>
<i>kauplemist trade</i>	<i>keelatakse prohibits</i>
	<i>keelatud prohibited</i>
	<i>keelavad prohibit</i>
	<i>keelduda refuse</i>

<i>keeldude</i>	<i>prohibitions</i>	<i>kehtestatud</i>	<i>monitoring</i>
<i>keeles</i>	<i>ametlikus</i>	<i>kehtestavad</i>	<i>member</i>
<i>keeltes</i>	<i>keeles</i>	<i>kehtetuks</i>	<i>repealed</i>
<i>keeltest</i>	<i>languages</i>	<i>kehtetuks tunnistatud</i>	<i>eec</i>
<i>keelust</i>	<i>prohibitions</i>	<i>kehtib</i>	<i>apply</i>
<i>keemilised</i>	<i>mechanical</i>	<i>kehtib</i>	<i>valid</i>
<i>keemiliselt</i>	<i>chemically</i>	<i>kehtib</i>	<i>paragraph</i>
<i>keemiliste</i>	<i>chemical</i>	<i>kehtinud</i>	<i>applicable</i>
<i>keeratav</i>	<i>fine-cut</i>	<i>kehtiv</i>	<i>valid</i>
<i>keeratava</i>	<i>fine-cut</i>	<i>kehtiva</i>	<i>operating</i>
<i>keeratud</i>	<i>particles</i>	<i>kehtiva</i>	<i>valid</i>
<i>keersmaeker</i>	<i>keersmaeker</i>	<i>kehtivad</i>	<i>apply</i>
<i>kehtestabfalls</i>		<i>kehtivad</i>	<i>sitting</i>
<i>kehtestada</i>	<i>establish</i>	<i>kehtivate</i>	<i>force</i>
<i>kehtestada</i>	<i>made</i>	<i>kehtivate</i>	<i>implemented</i>
<i>kehtestama</i>	<i>make</i>	<i>kehtivate</i>	<i>implement</i>
<i>kehtestama</i>	<i>established</i>	<i>kehtivus</i>	<i>validity</i>
<i>kehtestamast</i>	<i>introducing</i>	<i>kehtivusaeg</i>	<i>validity</i>
<i>kehtestamise</i>	<i>last</i>	<i>kehtivusaega</i>	<i>validity</i>
<i>kehtestamist</i>	<i>restriction</i>	<i>kehtivuse</i>	<i>validity</i>
<i>kehtestamist</i>	<i>restrictions</i>	<i>kehtivuse</i>	<i>expired</i>
<i>kehtestatakse</i>	<i>quotas</i>	<i>kehtivuse</i>	<i>expire</i>
<i>kehtestatakse</i>	<i>reduction</i>	<i>kehtivuspäevani</i>	<i>validity</i>
<i>kehtestatakse</i>	<i>reliefs</i>	<i>kehtivust</i>	<i>validity</i>
<i>kehtestatakse</i>	<i>forecast</i>	<i>kel</i>	<i>believe</i>
<i>kehtestatiproducers</i>		<i>kellel</i>	<i>legal</i>
<i>kehtestatud</i>	<i>established</i>	<i>kellel</i>	<i>holding</i>

<i>kellel</i>	<i>lawyers</i>	<i>keskpanga</i>	<i>bank.</i>
<i>kellel</i>	<i>lawyer</i>	<i>keskpangad</i>	<i>central banks</i>
<i>kellel</i>	<i>issues</i>	<i>keskpankade</i>	<i>central banks</i>
<i>kellele</i>	<i>whom</i>	<i>kestab</i>	<i>years</i>
<i>keraamilisel</i>	<i>ge.</i>	<i>kestab</i>	<i>year</i>
<i>kergesti</i>	<i>kergesti</i>	<i>kestus</i>	<i>duration</i>
<i>kerma</i>	<i>tarkoitettuihin</i>	<i>kestust</i>	<i>duration</i>
<i>kes</i>	<i>contracting</i>	<i>kiechle</i>	<i>kiechle</i>
<i>kes</i>	<i>see</i>	<i>kiidab</i>	<i>acting</i>
<i>kes</i>	<i>legal person</i>	<i>kiiremini</i>	<i>rapidly</i>
<i>kes</i>	<i>resident</i>	<i>kiiresti</i>	<i>soon</i>
<i>kes</i>	<i>running</i>	<i>kiirgusavarii</i>	<i>emergency</i>
<i>kes</i>	<i>exercise</i>	<i>kiitnud</i>	<i>by</i>
<i>kes</i>	<i>return</i>	<i>kilogrammi</i>	<i>carcasses</i>
<i>keskkonda</i>	<i>environment</i>	<i>kilogrammi</i>	<i>carcase</i>
<i>keskkonna</i>	<i>environment</i>	<i>kilogrammi</i>	<i>kilograms</i>
<i>keskkonnakaitse</i>	<i>environmental</i>	<i>kindla</i>	<i>fixed</i>
<i>keskkonnakaitseks</i>	<i>protect</i>	<i>kindlaks</i>	<i>analysis</i>
<i>keskkonnakomitee</i>	<i>environment</i>	<i>kindlaks</i>	<i>fixed</i>
<i>keskkonnameetmete</i>	<i>environmental</i>	<i>kindlaks</i>	<i>determine</i>
<i>keskkonnaotstarbeliste</i>	<i>system</i>	<i>kindlaks</i>	<i>stipulated</i>
<i>keskkonnapoliitika</i>	<i>policies</i>	<i>kindlaks</i>	<i>unanimous</i>
<i>keskkonnas</i>	<i>producing</i>	<i>kindlaks</i>	<i>visa</i>
<i>keskmise</i>	<i>average</i>	<i>kindlaks määrata</i>	<i>determine</i>
<i>keskmise</i>	<i>average</i>	<i>kindlaks määrata</i>	<i>determined</i>
<i>keskmise suurusega</i>	<i>small</i>	<i>kindlaks määratud</i>	<i>by</i>
<i>keskpanga</i>	<i>bank</i>	<i>kindlaksmääramiseni</i>	<i>date</i>

<i>kindlaksmääratud</i>	<i>fixed</i>	<i>kirjeldus</i>	<i>description</i>
<i>kindlatesse</i>	<i>limited</i>	<i>kirjelduse</i>	<i>description</i>
<i>kindlustada</i>	<i>ensure</i>	<i>kirjeldused</i>	<i>content</i>
<i>kindlustades</i>	<i>access</i>	<i>kirjelduste</i>	<i>particular</i>
<i>kindlustamisele</i>	<i>rural areas</i>	<i>kirjelduste</i>	<i>particulars</i>
<i>kinni</i>	<i>recovered</i>	<i>kirjutatud</i>	<i>signed</i>
<i>kinnipeetava</i>	<i>withholding tax</i>	<i>kirjutatud</i>	<i>forms</i>
<i>kinnisvara</i>	<i>property</i>	<i>kirjutatud</i>	<i>formed</i>
<i>kinnitab</i>	<i>approved</i>	<i>kirjutavad</i>	<i>signed</i>
<i>kinnitab</i>	<i>approve</i>	<i>kitsed</i>	<i>goats0</i>
<i>kinnitab</i>	<i>assurance</i>	<i>kitsede</i>	<i>goats</i>
<i>kinnitab</i>	<i>lay down</i>	<i>kiu</i>	<i>fibre</i>
<i>kinnitades</i>	<i>confirming</i>	<i>kiu</i>	<i>stated</i>
<i>kinnitades</i>	<i>reaffirming</i>	<i>kiu</i>	<i>state</i>
<i>kinnitatud</i>	<i>affixed</i>	<i>kiu kirjeldus</i>	<i>column'fibre</i>
<i>kinnitava</i>	<i>assurance</i>	<i>kiud</i>	<i>fibre</i>
<i>kinnitavad</i>	<i>reaffirm</i>	<i>kiud</i>	<i>fibres</i>
<i>kirikute</i>	<i>churches</i>	<i>kiudude</i>	<i>fibres</i>
<i>kirjalik</i>	<i>written</i>	<i>kiududest</i>	<i>backing</i>
<i>kirjalike</i>	<i>written</i>	<i>kiudusid</i>	<i>fibres</i>
<i>kirjalikke</i>	<i>written</i>	<i>kiukoostis</i>	<i>fibre composition</i>
<i>kirjaliku</i>	<i>written</i>	<i>kiukoostis</i>	<i>base</i>
<i>kirjalikult</i>	<i>writing</i>	<i>kiukoostise</i>	<i>fibre composition</i>
<i>kirjalikust</i>	<i>written</i>	<i>kiuprotsendid</i>	<i>fibre</i>
<i>kirjas</i>	<i>tr&uuml</i>	<i>kiusisaldus</i>	<i>stated</i>
<i>kirjavahetuse</i>	<i>letters</i>	<i>kiusisalduse</i>	<i>stated</i>
<i>kirjavahetuses</i>	<i>letters</i>	<i>kiusisaldusega</i>	<i>contents</i>

<i>kiusisaldusega</i>	<i>content</i>	<i>kodukorras</i>	<i>rules</i>
<i>klauslile</i>	<i>referred</i>	<i>kogeda</i>	<i>financing</i>
<i>klauslite</i>	<i>clauses</i>	<i>kogemused</i>	<i>experience</i>
<i>kleinwalsertali</i>	<i>austria</i>	<i>kogemuste</i>	<i>experiences</i>
<i>kleinwalsertalis</i>	<i>austria</i>	<i>kogemuste</i>	<i>experience</i>
<i>kliendideklaratsioon</i>	<i>showing</i>	<i>kogemustele</i>	<i>experiences</i>
<i>kliendile</i>	<i>declaration</i>	<i>kogu</i>	<i>entire</i>
<i>klient</i>	<i>furnish</i>	<i>kogu toote</i>	<i>product</i>
<i>kliima</i>	<i>climate</i>	<i>kogumassist</i>	<i>total weight</i>
<i>kodade</i>	<i>consisting</i>	<i>kogus</i>	<i>quantity</i>
<i>kodakondsus</i>	<i>citizenship</i>	<i>kogused</i>	<i>quantities</i>
<i>kodakondsus</i>	<i>nationality</i>	<i>kogusega</i>	<i>slips</i>
<i>kodakondsuse</i>	<i>nationality</i>	<i>koguseid</i>	<i>quantities</i>
<i>kodakondsuse</i>	<i>citizenship</i>	<i>kogusumma</i>	<i>burden</i>
<i>kodakondsuseta</i>	<i>stateless</i>	<i>kogutoodangu</i>	<i>accordance</i>
<i>kodakondsust.</i>	<i>citizenship.</i>	<i>koha</i>	<i>office</i>
<i>kodanik</i>	<i>nationals</i>	<i>kohal</i>	<i>chambers</i>
<i>kodanik</i>	<i>national</i>	<i>kohal</i>	<i>chamber</i>
<i>kodanike</i>	<i>nationals</i>	<i>kohaldada</i>	<i>apply</i>
<i>kodanike</i>	<i>citizens</i>	<i>kohaldada</i>	<i>applicable</i>
<i>kodanikel</i>	<i>nationals</i>	<i>kohaldada</i>	<i>applying</i>
<i>kodanikele</i>	<i>nationals</i>	<i>kohaldades</i>	<i>applying</i>
<i>kodanikud</i>	<i>nationals</i>	<i>kohaldamise</i>	<i>application</i>
<i>kodukord</i>	<i>procedure</i>	<i>kohaldamise</i>	<i>applications</i>
<i>kodukorra</i>	<i>procedure</i>	<i>kohaldamise</i>	<i>id</i>
<i>kodukorra.</i>	<i>procedure.</i>	<i>kohaldamise</i>	<i>applying</i>
<i>kodukorraga</i>	<i>rules</i>	<i>kohaldamisega</i>	<i>application</i>

<i>kohaldamiseks</i>	<i>purpose</i>	<i>kohandatud</i>	<i>adapted</i>
<i>kohaldamiseks</i>	<i>purposes</i>	<i>kohandusi</i>	<i>adjustments</i>
<i>kohaldamiseks</i>	<i>application</i>	<i>kohapeal</i>	<i>institutions</i>
<i>kohaldamisel</i>	<i>purposes</i>	<i>kohasel</i>	<i>post</i>
<i>kohaldamisel</i>	<i>purpose</i>	<i>kohaselt</i>	<i>cohesion</i>
<i>kohaldamisel</i>	<i>application</i>	<i>kohaselt</i>	<i>according</i>
<i>kohaldamisel</i>	<i>throughout</i>	<i>kohaselt.</i>	<i>management.</i>
<i>kohaldamist</i>	<i>preceding</i>	<i>kohasust</i>	<i>transactions</i>
<i>kohaldata</i>	<i>apply</i>	<i>kohdan</i>	<i>kohdan</i>
<i>kohaldataks</i>	<i>applied</i>	<i>kohe</i>	<i>immediately</i>
<i>kohaldatakse</i>	<i>apply</i>	<i>kohestatud</i>	<i>other</i>
<i>kohaldatakse</i>	<i>apply.</i>	<i>kohestatud</i>	<i>animal hair</i>
<i>kohaldatakse</i>	<i>applied</i>	<i>kohta</i>	<i>certain arable crops</i>
<i>kohaldatakse</i>	<i>applies</i>	<i>kohta</i>	<i>per</i>
<i>kohaldatav</i>	<i>applicable</i>	<i>kohta</i>	<i>organization</i>
<i>kohaldatava</i>	<i>applicable</i>	<i>kohta</i>	<i>investment</i>
<i>kohaldatavad</i>	<i>applicable</i>	<i>kohta</i>	<i>clothing</i>
<i>kohaldatavaks</i>	<i>become</i>	<i>kohta</i>	<i>c</i>
<i>kohaldatavate</i>	<i>applicable</i>	<i>kohta</i>	<i>tomatoes</i>
<i>kohalikud</i>	<i>regional</i>	<i>kohta</i>	<i>western</i>
<i>kohandada</i>	<i>adapted</i>	<i>kohta kehtivad</i>	<i>subject</i>
<i>kohandada</i>	<i>adapt</i>	<i>kohta vastavalt</i>	<i>control</i>
<i>kohandada</i>	<i>adjusted</i>	<i>kohta vastavalt</i>	<i>controls</i>
<i>kohandamiseks</i>	<i>adapting</i>	<i>kohta.</i>	<i>per</i>
<i>kohandamist</i>	<i>adaptation</i>	<i>kohta.</i>	<i>investment</i>
<i>kohandatakse</i>	<i>adapting</i>	<i>kohtades</i>	<i>places</i>
<i>kohandatakse</i>	<i>adjusted</i>	<i>kohtadesse</i>	<i>points</i>

<i>kohtlemine</i>	<i>treatment</i>	<i>kohtunikku</i>	<i>judge</i>
<i>kohtu</i>	<i>justice</i>	<i>kohtuniku</i>	<i>judge</i>
<i>kohtu</i>	<i>secrecy</i>	<i>kohtunikud</i>	<i>judges</i>
<i>kohtu</i>	<i>court</i>	<i>kohtunud ukogu</i>	
<i>kohtu</i>	<i>pending before</i>	<i>kohtuotsuse</i>	<i>judgment</i>
<i>kohtu</i>	<i>judicial authority</i>	<i>kohtuotsusega</i>	<i>decision</i>
<i>kohtu otsusega</i>	<i>settled</i>	<i>kohtus</i>	<i>justice</i>
<i>kohtuasja</i>	<i>ruling</i>	<i>kohtus</i>	<i>practise</i>
<i>kohtuasjasse</i>	<i>instance</i>	<i>kohtus</i>	<i>pending before</i>
<i>kohtuistung</i>	<i>hearing</i>	<i>kohtus</i>	<i>judicial authority</i>
<i>kohtujurist</i>	<i>reasoned</i>	<i>kohtusekretär</i>	<i>registrar</i>
<i>kohtujurist</i>	<i>reason</i>	<i>kohtusekretäri</i>	<i>assistant</i>
<i>kohtujuristi</i>	<i>advocate</i>	<i>kohtusekretäri</i>	<i>registrar</i>
<i>kohtujuristid</i>	<i>advocates-general</i>	<i>kohtusekretärile</i>	<i>instance</i>
<i>kohtujuristide</i>	<i>advocates-general</i>	<i>kohtusekretärile</i>	<i>registrar</i>
<i>kohtujuristide</i>	<i>advocates</i>	<i>kohtusse</i>	<i>justice</i>
<i>kohtukulude</i>	<i>costs</i>	<i>kohtusse</i>	<i>court</i>
<i>kohtule</i>	<i>case</i>	<i>kohtusse</i>	<i>courts</i>
<i>kohtule</i>	<i>appealed</i>	<i>kohtusse</i>	<i>appear</i>
<i>kohtule</i>	<i>appeal</i>	<i>kohtute</i>	<i>judicial</i>
<i>kohtule</i>	<i>document</i>	<i>kohus</i>	<i>court</i>
<i>kohtult</i>	<i>decision</i>	<i>kohus</i>	<i>instance</i>
<i>kohtumiste</i>	<i>intervals</i>	<i>kohusetundlikult</i>	<i>preserve</i>
<i>kohtunik</i>	<i>judge</i>	<i>kohustatud</i>	<i>obliged</i>
<i>kohtunike</i>	<i>judges</i>	<i>kohustub</i>	<i>undertaken</i>
<i>kohtunikel</i>	<i>immune</i>	<i>kohustus</i>	<i>obligation</i>
<i>kohtunikku</i>	<i>judges</i>	<i>kohustuse</i>	<i>obligation</i>

<i>kohustused</i>	<i>mfn</i>	<i>kohustuvad</i>	<i>undertake</i>
<i>kohustused</i>	<i>obligations</i>	<i>kohvikannusoojendajad</i>	<i>coffee cosy</i>
<i>kohustused</i>	<i>responsibilities</i>	<i>covers</i>	
<i>kohustusest</i>	<i>obligation</i>	<i>koja</i>	<i>chamber</i>
<i>kohustusiobligations</i>		<i>kokku</i>	<i>agreed</i>
<i>kohustusiimpartially</i>		<i>kokku</i>	<i>agree</i>
<i>kohustusiduties</i>		<i>kokku</i>	<i>convene</i>
<i>kohustusiresponsibilities</i>		<i>kokku</i>	<i>member</i>
<i>kohustusi offering</i>		<i>kokku</i>	<i>members</i>
<i>kohustusi addition</i>		<i>kokku</i>	<i>set</i>
<i>kohustuslikku</i>	<i>partly</i>	<i>kokku järgmises</i>	<i>agrees</i>
<i>kohustuslikuks</i>	<i>destination</i>	<i>kokku järgmises</i>	<i>agree</i>
<i>kohustustobligations</i>		<i>kokku lepitud</i>	<i>agreed</i>
<i>kohustustprovide</i>		<i>kokku leppinud</i>	<i>agreed</i>
<i>kohustuste</i>	<i>specific commitments</i>	<i>kokkulepe</i>	<i>agreement</i>
<i>kohustuste</i>	<i>obligations</i>	<i>kokkulepete</i>	<i>agreements</i>
<i>kohustuste</i>	<i>responsibilities</i>	<i>kokkulepitud</i>	<i>agreed</i>
<i>kohustuste täitmiseks</i>	<i>obligations</i>	<i>kokkuleppe</i>	<i>agreement</i>
<i>kohustuste täitmisel</i>	<i>performance</i>	<i>kokkuleppe</i>	<i>agreements</i>
<i>kohustuste täitmisel</i>	<i>discharge</i>	<i>kokkuleppega</i>	<i>bank</i>
<i>kohustuste täitmist</i>	<i>responsibilities</i>	<i>kokkuleppeid</i>	<i>amendment</i>
<i>kohustuste täitmist</i>	<i>performance</i>	<i>kokkuleppel</i>	<i>agreement</i>
<i>kohustustega</i>	<i>debts</i>	<i>kokkuleppelised</i>	<i>agreed allowances</i>
<i>kohustustega</i>	<i>obligation</i>	<i>kokkuleppelist</i>	<i>allowances</i>
<i>kohustustega</i>	<i>obligations</i>	<i>kokkuvõttetabelid</i>	<i>tables</i>
<i>kohustustele</i>	<i>legal</i>	<i>kokkuvõttetabelit</i>	<i>summary</i>
<i>kohustustest</i>	<i>admission</i>	<i>kokkuvõttetabelite</i>	<i>summary</i>

<i>kollane</i>	<i>yellow</i>	<i>kolmeneljandikulise</i>	<i>three-fourths</i>
<i>kollektiivse</i>	<i>collective</i>	<i>kolmeneljandikulise</i>	<i>fourths</i>
<i>kolm</i>	<i>three</i>	<i>kolmest</i>	<i>three</i>
<i>kolmanda</i>	<i>third</i>	<i>kolmeteistkümne</i>	<i>thirteen</i>
<i>kolmanda isiku vastu tugineda</i>	<i>third parties</i>	<i>koma</i>	<i>citation</i>
<i>kolmanda riigi</i>	<i>third country</i>	<i>komisjon</i>	<i>commission</i>
<i>kolmandal</i>	<i>third</i>	<i>komisjon</i>	<i>komisjon</i>
<i>kolmandal etapil</i>	<i>third stage</i>	<i>komisjoni</i>	<i>commission</i>
<i>kolmandas lõigus</i>	<i>third subparagraph</i>	<i>komisjoni</i>	<i>komisjoni</i>
<i>kolmandast</i>	<i>third</i>	<i>komisjoni</i>	<i>amended by regulation</i>
<i>kolmandate</i>	<i>third parties</i>	<i>komisjoni</i>	<i>board</i>
<i>kolmandate isikute</i>	<i>third</i>	<i>komisjoni</i>	<i>n</i>
<i>kolmandate riikide</i>	<i>countries</i>	<i>komisjoni</i>	<i>schoolchildren</i>
<i>kolmandate riikide kodanike</i>	<i>nationals</i>	<i>komisjoni ettepaneku põhjal</i>	<i>proposal</i>
<i>kolmandatesse</i>	<i>export</i>	<i>komisjoni ettepanekut</i>	<i>whereas article</i>
<i>kolmandatest riikidest</i>	<i>countries</i>	<i>komisjoni nimel komisjoni liige</i>	<i>oj</i>
<i>kolmandikku</i>	<i>upon</i>	<i>komisjoni nimel komisjoni liige</i>	<i>franz</i>
<i>kolmandiku</i>	<i>thirds</i>	<i>fischler</i>	<i>commission</i>
<i>kolmandiku</i>	<i>third</i>	<i>komisjoni nimel komisjoni liige</i>	<i>monti</i>
<i>kolmas</i>	<i>third</i>	<i>ren&eacute;ren&eacute;</i>	<i>ren&eacute;ren&eacute;</i>
<i>kolmas</i>	<i>by</i>	<i>komisjoniga</i>	<i>commission</i>
<i>kolme</i>	<i>three</i>	<i>komisjonil</i>	<i>commission</i>
<i>kolme</i>	<i>and'four</i>	<i>komisjonile</i>	<i>commission</i>
<i>kolme kuu jooksul pärast</i>	<i>months</i>	<i>komisjonile</i>	<i>komisjonile</i>
<i>kolmekomponentsete</i>	<i>ternary</i>	<i>komisjonile.</i>	<i>them.</i>
<i>kolmekümnendal</i>	<i>30th</i>	<i>komitee</i>	<i>committee</i>

<i>komitee</i>	<i>komitee</i>	<i>konsultatsioone</i>	<i>consultation</i>
<i>komiteed</i>	<i>committee</i>	<i>konsultatsiooni</i>	<i>consultations</i>
<i>komiteede</i>	<i>committees</i>	<i>konsultatsiooni</i>	<i>consultation</i>
<i>komiteega</i>	<i>committee</i>	<i>konsultatsioonide</i>	<i>consultations</i>
<i>komiteele</i>	<i>committee</i>	<i>konsultatsioonide</i>	<i>consultation</i>
<i>komiteele</i>	<i>komiteele</i>	<i>konsultatsioonikava</i>	<i>plan</i>
<i>komitees</i>	<i>committee</i>	<i>konsultatsioonimenetlust</i>	<i>procedures</i>
<i>kompensatsiooni</i>	<i>compensation</i>	<i>konsulteerib</i>	<i>consult</i>
<i>kompensatsiooni saamise tingimustele</i>	<i>producers</i>	<i>konsulteerib</i>	<i>consults</i>
<i>kompensatsioonide</i>	<i>producers</i>	<i>konsulteerib</i>	<i>consulted</i>
<i>kompenseerimise</i>	<i>compensation</i>	<i>konsulteerida</i>	<i>consulted</i>
<i>kompenseerimiseks</i>	<i>enforced</i>	<i>konsulteerida</i>	<i>relating</i>
<i>kompenseeriv</i>	<i>adjustment</i>	<i>konsulteerides</i>	<i>authority</i>
<i>kompenseeriva</i>	<i>compensatory</i>	<i>konsulteerimis</i>	<i>consultation</i>
<i>koncentreras</i>	<i>koncentreras</i>	<i>konsulteerimise</i>	<i>consultation</i>
<i>koncentrerat</i>	<i>koncentrerat</i>	<i>konsulteerimist</i>	<i>consulting</i>
<i>koncentrerat smör för</i>	<i>koncentrerat</i>	<i>konsulteerimist</i>	<i>consultation</i>
<i>smoer foer</i>		<i>konsulteerimist</i>	<i>consultations</i>
<i>konföderatsiooni</i>	<i>annexed</i>	<i>konsulteerimist</i>	<i>euroopa parlamendiga</i>
<i>konkreetse</i>	<i>particular</i>		<i>consulting</i>
<i>konkurentsi</i>	<i>competition</i>	<i>konsulteeris</i>	<i>cross-industry</i>
<i>konkurentsivõime</i>	<i>competitiveness</i>	<i>konsulteeriv</i>	<i>consulting</i>
<i>konkurentsivõime</i>	<i>competitiveness</i>	<i>konsulteeriva</i>	<i>consulting</i>
<i>konkurentsivõime</i>	<i>competition law</i>	<i>konsulteerivad</i>	<i>consulted</i>
<i>konsensus</i>	<i>consensus</i>	<i>konsulteerivad</i>	<i>consult</i>
<i>konsensus</i>	<i>consensus</i>	<i>kontaktkomitee</i>	<i>contact</i>
<i>konsensus</i>	<i>consensus</i>	<i>kontekstis</i>	<i>context</i>
<i>konsultatsioone</i>	<i>consultations</i>		

<i>kontrastsele</i>	<i>taustale</i>	<i>kontsessioonide</i>	<i>concessions</i>
<i>kontroll</i>	<i>control</i>	<i>kontsessioonide</i>	<i>cessionaires</i>
<i>kontrollakti</i>	<i>participates</i>	<i>kontsessiooniga</i>	<i>affected</i>
<i>kontrolli</i>	<i>checks</i>	<i>konventsioone</i>	<i>conventions</i>
<i>kontrolli</i>	<i>control</i>	<i>konventsiooni</i>	<i>convention</i>
<i>kontrolli</i>	<i>controls</i>	<i>konventsiooni osapoolte</i>	<i>quota</i>
<i>kontrollib</i>	<i>examine</i>	<i>konventsiooni osapoolte</i>	<i>quotas</i>
<i>kontrollib</i>	<i>verify</i>	<i>konventsioonide</i>	<i>conventions</i>
<i>kontrollida</i>	<i>inspect</i>	<i>konventsiooniga</i>	<i>whereas</i>
<i>kontrollikoda</i>	<i>auditors</i>	<i>konverents</i>	<i>conference</i>
<i>kontrollikoda</i>	<i>court</i>	<i>konverents lepib kokku</i>	<i>agrees</i>
<i>kontrollikoda</i>	<i>doing</i>	<i>konverentsi</i>	<i>conference</i>
<i>kontrollikoja</i>	<i>reports</i>	<i>konverentsidel</i>	<i>conferences</i>
<i>kontrollikoja</i>	<i>auditors</i>	<i>konverentsil</i>	<i>conference</i>
<i>kontrollikoja</i>	<i>agreement between</i>	<i>konverentsil</i>	<i>conferences</i>
<i>kontrollikoja</i>	<i>managing</i>	<i>konverentsile</i>	<i>conference</i>
<i>kontrollikojale</i>	<i>court</i>	<i>koodi</i>	<i>code</i>
<i>kontrollimenetlusele</i>	<i>non-economic</i>	<i>koodi</i>	<i>codes</i>
<i>kontrollimiseks</i>	<i>bank.</i>	<i>koodidega</i>	<i>by</i>
<i>kontrollimiseks.</i>	<i>bank.</i>	<i>koodidega</i>	<i>i.a</i>
<i>kontrollitakse</i>	<i>performed</i>	<i>koodiga</i>	<i>code</i>
<i>kontrollitakse</i>	<i>carried</i>	<i>koodiga</i>	<i>codes</i>
<i>kontsentreeritud</i>	<i>concentrated</i>	<i>koodiga</i>	<i>replaced</i>
<i>kontsessioone</i>	<i>free</i>	<i>koondnomenklatuuri</i>	<i>nomenclature</i>
<i>kontsessiooni</i>	<i>concession</i>	<i>koopia</i>	<i>copy</i>
<i>kontsessioonid</i>	<i>concessions</i>	<i>koopia</i>	<i>proof</i>
<i>kontsessioonide</i>	<i>commitments</i>	<i>koopiad</i>	<i>copies</i>

<i>koopiat</i>	<i>issued</i>	<i>koostab</i>	<i>up</i>
<i>koopiate</i>	<i>application</i>	<i>koostab</i>	<i>prepare</i>
<i>koordineerimist</i>	<i>policies</i>	<i>koostada</i>	<i>up</i>
<i>kooritud</i>	<i>skimmed-milk</i>	<i>koostamise</i>	<i>preparation</i>
<i>koos</i>	<i>together</i>	<i>koostamisel</i>	<i>up</i>
<i>koos</i>	<i>tariff quota</i>	<i>koostas</i>	<i>drew</i>
<i>kooskõla</i>	<i>coherence</i>	<i>koostatakse</i>	<i>drawn up</i>
<i>kooskõlas</i>	<i>accordance</i>	<i>koostatud</i>	<i>up</i>
<i>kooskõlas</i>	<i>conformity</i>	<i>koostatud</i>	<i>established</i>
<i>kooskõlas artiklis</i>	<i>accordance</i>	<i>koostatud</i>	<i>esimesele</i>
<i>kooskõlas käesoleva</i>	<i>compatible</i>	<i>koostis</i>	<i>composition</i>
<i>kooskõlas käesoleva lepingu</i>		<i>koostise</i>	<i>composition</i>
<i>accordance</i>		<i>koostisega</i>	<i>composition</i>
<i>kooskõlastamiseks</i>	<i>coordination</i>	<i>koostises</i>	<i>composition</i>
<i>kooskõlastamiseks.</i>	<i>coordination</i>	<i>koostisosad</i>	<i>components</i>
<i>kooskõlastamist</i>	<i>coordination</i>	<i>koostisosi</i>	<i>representing</i>
<i>kooskõlastavad</i>	<i>coordinate</i>	<i>koostöö</i>	<i>cooperation</i>
<i>koosneb</i>	<i>consist</i>	<i>koostööd</i>	<i>cooperation</i>
<i>koosneb</i>	<i>fibre</i>	<i>koostööd</i>	<i>cooperate</i>
<i>koosneb</i>	<i>fibres</i>	<i>koostööd</i>	<i>member</i>
<i>koosnev</i>	<i>preparing</i>	<i>koostööga</i>	<i>counter</i>
<i>koosnev</i>	<i>composed</i>	<i>koostööga</i>	<i>cooperation</i>
<i>koosneva</i>	<i>composed</i>	<i>koostöös</i>	<i>liaison</i>
<i>koosnevate</i>	<i>judges</i>	<i>koostöös</i>	<i>cooperation</i>
<i>koosseisus</i>	<i>sitting</i>	<i>kopenhaageni</i>	<i>copenhagen</i>
<i>koosseisu</i>	<i>composition</i>	<i>kord</i>	<i>arrangements</i>
<i>koosseisus</i>	<i>composition</i>	<i>kord aastas</i>	<i>once</i>

<i>korda</i>	<i>arrangements</i>	<i>korralduse</i>	<i>common</i>
<i>kordamööda</i>	<i>examine</i>	<i>korralduskomitee</i>	<i>committee</i>
<i>kordumatu</i>	<i>means</i>	<i>korralduskomitee</i>	<i>whereas</i>
<i>korduvaid</i>	<i>chain</i>	<i>korralduskomiteede</i>	<i>233/94</i>
<i>korea</i>	<i>kong</i>	<i>korrapäraselt</i>	<i>periodically</i>
<i>korra</i>	<i>arrangements</i>	<i>korrapäraselt</i>	<i>regularly</i>
<i>korra kohaselt</i>	<i>uniform</i>	<i>korrektse</i>	<i>finds</i>
<i>korra kohta</i>	<i>applying</i>	<i>korrelatsioonitabel</i>	<i>correlation table</i>
<i>korra üksikasjalikud rakenduseeskirjad</i>		<i>korrigeerib</i>	<i>adjust</i>
<i>piima</i>	<i>between</i>	<i>korsetitoodete</i>	<i>corsetry</i>
<i>korral</i>	<i>consultation</i>	<i>korvausta</i>	<i>'ilman</i>
<i>korral</i>	<i>radiological</i>	<i>kosmeetikakotid</i>	<i>make-up cases</i>
<i>korral</i>	<i>necessary</i>	<i>kostjaks</i>	<i>proceedings</i>
<i>korral</i>	<i>event</i>	<i>kotid</i>	<i>cases</i>
<i>korraldada</i>	<i>conducted</i>	<i>kraasitud</i>	<i>combed</i>
<i>korraldada</i>	<i>conduct</i>	<i>kraasmed</i>	<i>ex</i>
<i>korraldatakse</i>	<i>rights</i>	<i>kraasmenetletud</i>	<i>case</i>
<i>korraldav</i>	<i>authority within</i>	<i>kreeka</i>	<i>republic</i>
<i>korraldava</i>	<i>authority</i>	<i>kreeka</i>	<i>greece</i>
<i>korraldus</i>	<i>organization</i>	<i>kreekas</i>	<i>greece</i>
<i>korraldus</i>	<i>gatt</i>	<i>kriiSIDE</i>	<i>forces</i>
<i>korraldusasutus</i>	<i>management</i>	<i>kriiSiolukordade</i>	<i>crisis</i>
<i>korraldusasutuse</i>	<i>management authority</i>	<i>kriiSiolukordades</i>	<i>crisis</i>
<i>korraldusasutused</i>	<i>management authorities</i>	<i>kriiSiivaatluskeskus</i>	<i>cell</i>
<i>korraldusasutusele</i>	<i>management authority</i>	<i>kriiSiivaatluskeskuse</i>	<i>cell</i>
<i>korraldusasutusele</i>	<i>satisfied</i>	<i>kriminaalasjades</i>	<i>criminal matters</i>
		<i>kriminaalõiguse</i>	<i>criminal</i>

<i>kriteeriume</i>	<i>criteria</i>	<i>kuivõrd</i>	<i>far</i>
<i>kriteeriumid</i>	<i>criteria</i>	<i>kuivõrd</i>	<i>insofar</i>
<i>krokodilliliste</i>	<i>tanned</i>	<i>kujul</i>	<i>form</i>
<i>kude</i>	<i>cotton</i>	<i>kujul</i>	<i>forms</i>
<i>kui</i>	<i>wound up</i>	<i>kujundamine</i>	<i>framing</i>
<i>kui</i>	<i>price</i>	<i>kujundamist</i>	<i>defence</i>
<i>kui</i>	<i>unless</i>	<i>kuld</i>	<i>environmentally</i>
<i>kui</i>	<i>authority</i>	<i>kulla</i>	<i>gold</i>
<i>kui</i>	<i>save</i>	<i>kullaga</i>	<i>gold</i>
<i>kui</i>	<i>voting</i>	<i>kullas</i>	<i>gold</i>
<i>kui</i>	<i>member</i>	<i>kultuuride</i>	<i>crops</i>
<i>kui</i>	<i>majority</i>	<i>kultuurliikide</i>	<i>parental</i>
<i>kui</i>	<i>audit</i>	<i>kulud</i>	<i>charged</i>
<i>kui</i>	<i>persistent</i>	<i>kulud</i>	<i>charge</i>
<i>kui</i>	<i>used</i>	<i>kulud</i>	<i>expenditure</i>
<i>kui</i>	<i>longer</i>	<i>kulude</i>	<i>expenditure</i>
<i>kui</i>	<i>seven</i>	<i>kulusid</i>	<i>expenditure</i>
<i>kui</i>	<i>derogation</i>	<i>kulusid</i>	<i>expenses</i>
<i>kui</i>	<i>response</i>	<i>kulutukseen</i>	<i>kulutukseen</i>
<i>kui</i>	<i>health</i>	<i>kulutused</i>	<i>expenditure</i>
<i>kui</i>	<i>whenever</i>	<i>kumb</i>	<i>whichever</i>
<i>kui</i>	<i>kui</i>	<i>kummipaelad</i>	<i>elastic</i>
<i>kui</i>	<i>appeal</i>	<i>kuna</i>	<i>since</i>
<i>kui</i>	<i>annexed</i>	<i>kuni</i>	<i>until such time</i>
<i>kui</i>	<i>distributed</i>	<i>kuni</i>	<i>headings</i>
<i>kui</i>	<i>w</i>	<i>kuni</i>	<i>pending</i>
<i>kuivsöödaturu</i>	<i>oils</i>	<i>kuni</i>	<i>up</i>

<i>kuni</i>	<i>kuni</i>	<i>kuu</i>	<i>month</i>
<i>kuningriigi</i>	<i>kingdom</i>	<i>kuud</i>	<i>six months</i>
<i>kuningriik</i>	<i>kingdom</i>	<i>kuud</i>	<i>months</i>
<i>kunstlikult</i>	<i>artificially propagated</i>	<i>kuud</i>	<i>month</i>
<i>kunstlikult paljundatud</i>	<i>propagated</i>	<i>kuue</i>	<i>six</i>
<i>kunstlilled</i>	<i>artificial flowers</i>	<i>kuue kuu jooksul pärast</i>	<i>months</i>
<i>kuritegevuse</i>	<i>combating</i>	<i>kuueks</i>	<i>six</i>
<i>kuritegude</i>	<i>offences</i>	<i>kuulama</i>	<i>hear</i>
<i>kuriteo</i>	<i>criminal</i>	<i>kuulata</i>	<i>procedure</i>
<i>kursiga</i>	<i>monetary</i>	<i>kuuluksid</i>	<i>include</i>
<i>kus</i>	<i>registry</i>	<i>kuuluv</i>	<i>falling</i>
<i>kus</i>	<i>kus</i>	<i>kuuluva</i>	<i>falling</i>
<i>kus</i>	<i>leghold trap</i>	<i>kuuluvad</i>	<i>fall</i>
<i>kus</i>	<i>situated</i>	<i>kuuluvad</i>	<i>falling within</i>
<i>kus</i>	<i>member</i>	<i>kuuluvaid</i>	<i>export</i>
<i>kust</i>	<i>state</i>	<i>kuuluvaid</i>	<i>contaminated</i>
<i>kustutatakse</i>	<i>paragraph</i>	<i>kuuluvate</i>	<i>heading</i>
<i>kustutatakse</i>	<i>paragraphs</i>	<i>kuuluvate</i>	<i>manufactured tobacco</i>
<i>kustutatakse</i>	<i>deleted</i>	<i>kuuluvatele</i>	<i>cheeses</i>
<i>kustutatakse</i>	<i>article</i>	<i>kuuluvatele</i>	<i>cheese</i>
<i>kutseorgani</i>	<i>member</i>	<i>kuupäev</i>	<i>date</i>
<i>kutsub</i>	<i>invites</i>	<i>kuupäevadate</i>	
<i>kutsub</i>	<i>invite</i>	<i>kuupäeval</i>	<i>date</i>
<i>kutsuda</i>	<i>invited</i>	<i>kuupäevani</i>	<i>into</i>
<i>kutsuda</i>	<i>invite</i>	<i>kuupäevast</i>	<i>date</i>
<i>kutsuvad</i>	<i>invite</i>	<i>kuus kuud</i>	<i>six months</i>
<i>kuu</i>	<i>months</i>	<i>kvalifikatsiooni</i>	<i>qualifications</i>

<i>kvalifitseeritud</i>	<i>acting by a qualified majority</i>	<i>käesolev</i>	<i>force</i>
<i>kvalifitseeritud</i>	<i>qualified</i>	<i>käesolev</i>	<i>&uml</i>
<i>kvalifitseeritud häälteenamusega teha</i>	<i>qualified majority</i>	<i>käesolev</i>	<i>provisions</i>
<i>kvalifitseeritud häälteenamuseks</i>	<i>qualified majority</i>	<i>käesolev</i>	<i>keeping</i>
<i>kvaliteedi</i>	<i>quality</i>	<i>käesolev</i>	<i>closely</i>
<i>kvaliteeti</i>	<i>changing</i>	<i>käesolev määrus</i>	<i>communities.</i>
<i>kvantitatiivse</i>	<i>quantitative</i>	<i>käesolev määrus jõustub</i>	<i>enter</i>
<i>kvantitatiivsete</i>	<i>quantitative</i>	<i>käesolev määrus jõustub euroopa ühenduste teatajas avaldamise päeval</i>	<i>day</i>
<i>kvartali</i>	<i>second</i>	<i>käesolev määrus jõustub järgmisel päeval pärast</i>	<i>into</i>
<i>kvoodi</i>	<i>volumes</i>	<i>käesolev määrus jõustub kolmandal päeval pärast</i>	<i>force</i>
<i>kvoodi</i>	<i>volume</i>	<i>käesolev määrus jõustub seitsmendal päeval pärast</i>	<i>force</i>
<i>kvoodiga</i>	<i>replaced</i>	<i>käesoleva</i>	<i>conferred upon</i>
<i>kvoodile</i>	<i>quota</i>	<i>käesoleva</i>	<i>provided</i>
<i>kvoodimahu</i>	<i>volume</i>	<i>käesoleva</i>	<i>&uml</i>
<i>kvoodimahu</i>	<i>volumes</i>	<i>käesoleva</i>	<i>hereto</i>
<i>kvoot</i>	<i>respect</i>	<i>käesoleva</i>	<i>falling under</i>
<i>kvootide</i>	<i>quotas</i>	<i>käesoleva</i>	<i>act</i>
<i>kõrged lepinguosaliselised</i>	<i>high</i>	<i>käesoleva</i>	<i>force</i>
<i>käes</i>	<i>contracting</i>	<i>käesoleva</i>	<i>recommendations</i>
<i>käesolev</i>	<i>member states</i>	<i>käesoleva</i>	<i>&uml</i>
<i>käesolev</i>	<i>member states.</i>	<i>käesoleva</i>	<i>set</i>
<i>käesolev</i>	<i>hereto</i>	<i>käesoleva</i>	<i>provisions</i>
<i>käesolev</i>	<i>day following</i>	<i>käesoleva</i>	<i>annexes</i>
<i>käesolev</i>	<i>act</i>		

<i>käesoleva</i>	<i>lays down</i>	<i>käesolevaga</i>	<i>regulations</i>
<i>käesoleva</i>	<i>closely</i>	<i>käesolevaga</i>	<i>hereby</i>
<i>käesoleva</i>	<i>laid down</i>	<i>käesoleval</i>	<i>question</i>
<i>käesoleva</i>	<i>representatives</i>	<i>käesolevale</i>	<i>annexed</i>
<i>käesoleva</i>	<i>rulings</i>	<i>käesolevale</i>	<i>party</i>
<i>käesoleva akti</i>	<i>subject</i>	<i>käesolevale direktiivile</i>	<i>reference</i>
<i>käesoleva artikli kohaldamisel</i>	<i>purposes</i>	<i>käesolevale direktiivile</i>	<i>references</i>
<i>käesoleva artikli kohaldamisel</i>	<i>purpose</i>	<i>käesolevale protokollile lisatud</i>	<i>schedule</i>
<i>käesoleva direktiivi mõistes</i>		<i>käesolevale protokollile lisatud</i>	
<i>following</i>		<i>schedules</i>	
<i>käesoleva direktiiviga reguleeritavas</i>		<i>käesolevas</i>	<i>by</i>
<i>valdkonnas</i>	<i>text</i>	<i>käesolevas</i>	<i>provided</i>
<i>käesoleva direktiiviga reguleeritavas</i>		<i>käesolevas</i>	<i>s&auml;l</i>
<i>valdkonnas</i>	<i>texts</i>	<i>käesolevas</i>	<i>referred</i>
<i>käesoleva jaotise</i>	<i>title</i>	<i>käesolevas</i>	<i>falling under</i>
<i>käesoleva käsituslepe</i>	<i>understanding</i>	<i>käesolevas</i>	<i>set</i>
<i>käesoleva lepingu</i>	<i>treaty</i>	<i>käesolevas</i>	<i>lays down</i>
<i>käesoleva lepingu</i>	<i>treaty.</i>	<i>käesolevas</i>	<i>questions</i>
<i>käesoleva lepingu teiste sätete</i>		<i>käesolevas</i>	<i>laid down</i>
<i>provisions</i>		<i>käesolevas lepingus</i>	<i>treaty</i>
<i>käesoleva lõike alusel</i>	<i>paragraph</i>	<i>käesolevas lõikes</i>	<i>paragraph</i>
<i>käesoleva määruse</i>	<i>confidential</i>	<i>käesolevas määruses</i>	<i>purposes</i>
<i>information</i>		<i>käesolevas määruses ettenähtud</i>	<i>meetmed</i>
<i>käesoleva määruse lisas loetletud</i>		<i>measures</i>	
<i>regulation</i>		<i>käesolevas määruses sätestatud</i>	<i>meetmed</i>
<i>käesoleva protokolli</i>	<i>article 1</i>	<i>measures</i>	
<i>käesoleva põhikirja</i>	<i>statute</i>	<i>käesolevast lepingust</i>	<i>deriving</i>
<i>käesolevaga</i>	<i>regulation</i>	<i>käesolevast lepingust</i>	<i>member state</i>

<i>käesolevat lepingut</i>	<i>agreement</i>	<i>käsitlevates</i>	<i>approximation</i>
<i>käesolevat lõiget</i>	<i>paragraph</i>	<i>käsitsi</i>	<i>handmade</i>
<i>käibemaksu</i>	<i>tobacco</i>	<i>käsituslepe</i>	<i>understanding</i>
<i>käibemaksu</i>	<i>turnover</i>	<i>käsitusleppe</i>	<i>article</i>
<i>käibemaksu-omavahendite</i>	<i>vat</i>	<i>käsitusleppes</i>	<i>understanding</i>
<i>käigus</i>	<i>pending before</i>	<i>kättesaadav</i>	<i>nevertheless</i>
<i>käsitle</i>	<i>national</i>	<i>kättesaadavad</i>	<i>available</i>
<i>käsitle</i>	<i>nationals</i>	<i>kättesaadavaid</i>	<i>available</i>
<i>käsitleb</i>	<i>interim agreement</i>	<i>kättesaadavaks</i>	<i>available</i>
<i>käsitleb</i>	<i>scheduled</i>	<i>kõige</i>	<i>most</i>
<i>käsitleb</i>	<i>specific measures</i>	<i>kõigi</i>	<i>every</i>
<i>käsitletud</i>	<i>dealt</i>	<i>kõigi liikmete</i>	<i>membership</i>
<i>käsitlev</i>	<i>act concerning</i>	<i>kõigil</i>	<i>union</i>
<i>käsitlev</i>	<i>whereas</i>	<i>kõigil</i>	<i>practical arrangements</i>
<i>käsitleva</i>	<i>act concerning</i>	<i>kõigis</i>	<i>president j.</i>
<i>käsitleva</i>	<i>exemptions</i>	<i>kõigist</i>	<i>&ldquo</i>
<i>käsitleva protokoll</i>	<i>kingdom</i>	<i>kõik</i>	<i>membership</i>
<i>käsitlevad</i>	<i>relating</i>	<i>kõikide</i>	<i>copies</i>
<i>käsitlevad</i>	<i>sitlevad</i>	<i>kõikide</i>	<i>expansion</i>
<i>käsitlevaid</i>	<i>concerning</i>	<i>kõikidele liikmesriikidele</i>	<i>member states</i>
<i>käsitlevas</i>	<i>exemptions</i>	<i>kõikidele liikmesriikidele</i>	<i>member states.</i>
<i>käsitlevas</i>	<i>services</i>	<i>kõikidest</i>	<i>inform</i>
<i>käsitlevas</i>	<i>section</i>	<i>kõnealune</i>	<i>tribunal</i>
<i>käsitlevatwaste</i>		<i>kõnealune</i>	<i>authority</i>
<i>käsitlevatconcerning</i>		<i>kõnealuse</i>	<i>question</i>
<i>käsitlevate</i>	<i>mandate</i>	<i>kõnealuse</i>	<i>representative</i>
<i>käsitlevate</i>	<i>approximation</i>	<i>...</i>	