

Ilze AUZIŅA, Kristīne LEVĀNE-PETROVA, Roberts DARĢIS
(LU Matemātikas un informātikas institūts)

LATVIEŠU VALODAS APGUVĒJU KĻŪDU ANALĪZE: PAREIZRAKSTĪBAS KĻŪDAS

ANALYSIS OF THE LATVIAN LANGUAGE LEARNER ERRORS: SPELLING ERRORS

Atslēgvārdi: valodas apguvēju korpuss, kļūdu marķēšana, pareizrakstības kļūdas

Keywords: learner corpus, error annotation, spelling errors

Summary

This article presents the analysis of Latvian Language Learner Corpus errors, paying particular attention to spelling errors. The Corpus contains 1496 texts of different lengths (146 706 tokens in general). The Corpus was created in several stages: 1) data digitization; 2) text correction; 3) automated NLP analysis; 4) original and corrected text alignment; 5) automatic error annotation and manual review.

The methodology for the error annotation was also created. The automated error annotation includes following error types: spelling errors, punctuation errors, grammatical, syntactic, lexical errors and unclear text. The distribution of the error types is following: spelling errors (37%), punctuation errors (18%), inflection errors (17%), errors of combined type (for instance, lexical and spelling mistakes) (12%), syntactical error (8%), lexical error (4%), and unclear text (4%).

According to the analyzed corpus data spelling errors are the most common error type; therefore, the article analyses the spelling error by the subtypes, for instance, the spelling of diphthongs, omitted letters, capitalization, etc.

The spelling mistakes also reflect the ability of the learner to pronounce and listen to Latvian sounds and consonant clusters. There is also the influence of the mother tongue. The error analysis of the Corpus also shows the issues of the acquisition of Latvian and is the basis for the development of the teaching and methodological materials.

Ievads

Pēc atbilstošiem kritērijiem veidots valodas apguvēju korpuss ar valodas apguvēju pieļauto kļūdu marķējumu ļauj pētniekiem izveidot efektīvākus mācību materiālus un metodiku.

2017. gadā ar Latviešu valodas aģentūras (LVA) atbalstu pētījuma „Latviešu valodas prasmes kvalitāte: valsts valodas prasmes pārbaudes kārtotāju rezultāti” laikā LU Matemātikas un informātikas institūtā (LU MII) ir izveidota valsts valodas prasmes pārbaudes darbu datubāze – t.s. Latviešu valodas apguvēju korpuss, kurā apkopoti 900 valsts valodas pārbaudes darbu rakstītprasmes testi.

Kopumā korpusā ir iekļauti 1496 dažāda apjoma teksti, kuros ir 146706 tekstvienības (vārdformas un pieturzīmes). Korpusa izveides posmi ir šādi: 1) datu atlase un digitalizācija; 2) tekstu normalizēšana (mērķa hipotēzes izvirzīšana), 3) automātiska morfoloģiskā marķēšana, tostarp vārdšķiru

noteikšana, lemmatizēšana, 4) oriģinālā un normalizētā teksta sastatīšana, 5) automatizēta kļūdu anotēšana un manuāla pārskatīšana. (Dargis et al. 2018)

Visi teksti ir automātiski morfoloģiski marķēti. Tajos ir marķētas valodas apguvēju pieļautās kļūdas – pareizrakstības, formveidošanas un vārddarināšanas, leksikas, interpunkcijas un sintakses kļūdas – atbilstoši iepriekš izstrādātai kļūdu marķēšanas metodoloģijai. (Dargis et al. 2018) Kā atsevišķs sintakses kļūdu tips marķētas arī vārdu secības kļūdas. Reizēm, labojot autora darbu, nav skaidrs, kādu mērķa hipotēzi izvirzīt, t.i., kā valodas apguvēju izteikumu rekonstruēt mērķvalodā (Ellis 1994, 54; Znotiņa 2018), vai arī nav saprotama autora izteiktā doma, tādēļ šādiem gadījumiem paredzēts kļūdu tips “Nesaprotams teksts”, piemēram, *Līdz 15.00, no cenam 12 eiro; Es gribet tev pasākumi uz teatriem*. Savukārt reizēm, marķējot kļūdas, nav iespējams pateikt, vai valodas apguvējs nezina atbilstošo vārduformu, vai arī aizmirsis virs burta uzlikt garumzīmi, piemēram, *Es dzīvoju Rīga*. Šajā gadījumā tiek norādīti abi kļūdu tipi – formveidošanas un pareizrakstības kļūda, tādējādi veidojot atsevišķu kļūdu tipu – kombinētās kļūdas. Skat. 1. att.

Darbu analīze tiek veikta Latvijas Zinātnes padomes finansētā pētījuma projekta “Latviešu valodas apguvēju korpusa izveide: metodes, rīki un izmantojums” (projekta Nr. lzp-2018/1-0527) laikā.

Original	Draudzenīt		ka	tev	klājas	.
Edited	Draudzenīt	,	kā	Tev	klājas	?
Tag	ncfsv5	zc	r0m	pp20sdn	vmyipi130an	zs
Distance			1/2 (50%)	1/3 (33%)		1/1 (100%)
Syntax						
Order						
Spelling						
Word form						
Lexical						
Punctuation						
Unclear						

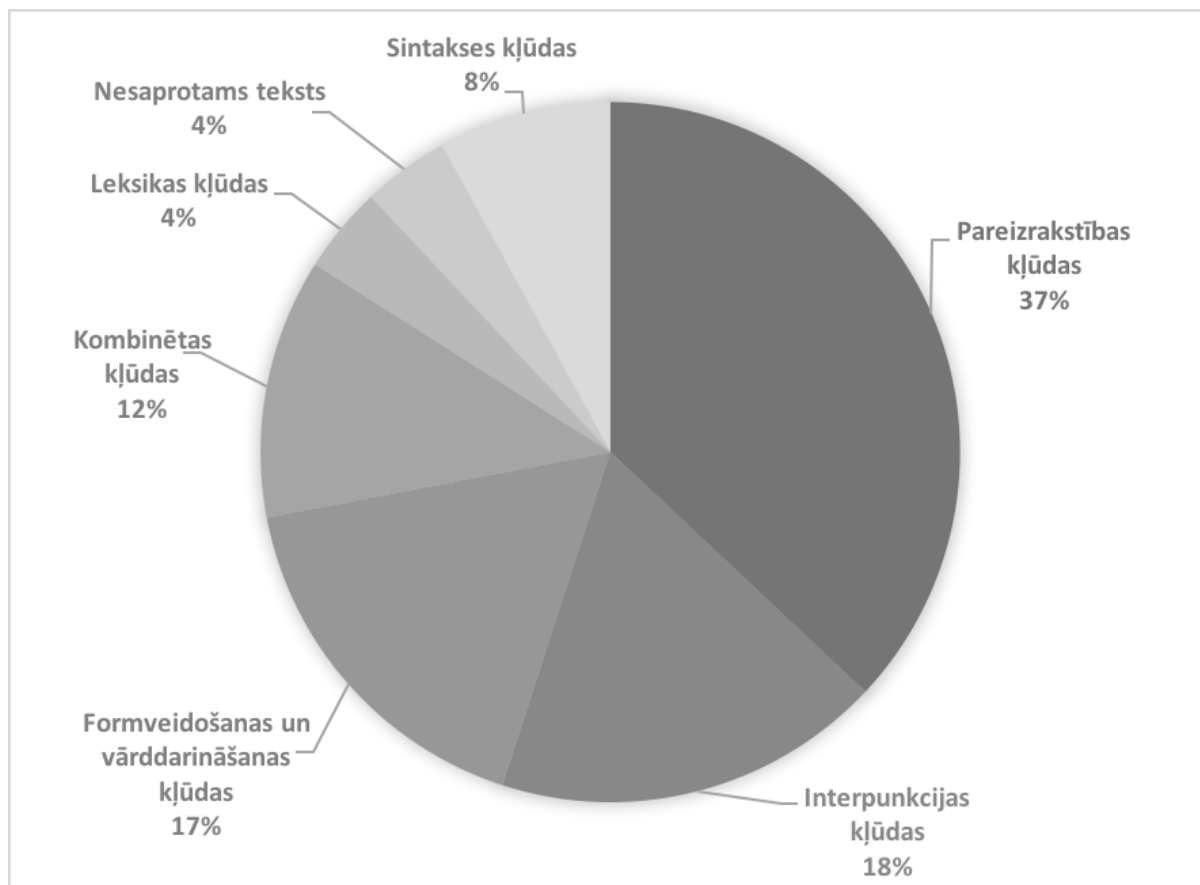
1. att. Kļūdu marķēšanas saskarnes fragmenti.

Kļūdu statistiskā analīze

Veicot valodas apguvēju kļūdu statistisko analīzi, noskaidrots, ka visizplatītākās ir pareizrakstības kļūdas (37%), tām seko interpunkcijas kļūdas (18%), formveidošanas un vārddarināšanas kļūdas (17%), kombinētās kļūdas (pareizrakstības un formveidošanas, leksikas un pareizrakstības u. tml.) (12%), sintakses kļūdas (ieskaitot vārdu secības kļūdas – 8%) leksikas kļūdas (4%) un nesaprotams teksts (4%) (skat. 2. att.). Kļūdu analīze rāda, ka 22% tekstvienību ir kļūdainas vai neatbilstošas. Reizēm nav skaidrs, kādam kļūdu tipam valodas

apguvēja pieļautā kļūda pieder, tādēļ ir izmantots dubultmarķējums – viena un tā pati kļūda marķēta, piemēram, kā pareizrakstības un formveidošanas kļūda (*Viņa dzīvo Rīga.* → *Viņa dzīvo Rīgā.*) vai kā pareizrakstības un leksikas kļūda (*Ka iet?* → *Kā iet?*).

Rakstā tiks analizēts tikai viens kļūdu tips – pareizrakstības kļūdas – un noteikts, kurš no pareizrakstības kļūdu apakštipiem (kopā vai šķirti rakstāmi vārdi, diakritiskās zīmes u. tml.) ir izplatīts konkrētā valodas apguves līmenī.



2. att. Kļūdu izplatība korpusā.

Pareizrakstības kļūdu apakšgrupas

Kā vārdi un vārdformas rakstāmas, kā un kādi grafiskie simboli izmantojami, atspoguļojot fonēmas, nosaka ortogrāfijas likumi jeb normas. Lai gan valodas attīstības gaitā normas var mainīties, tomēr katrā noteiktā laikposmā pieņemtās ortogrāfijas normas ir jāievēro. Tās ir jāapgūst arī ne tikai dzimtās valodas lietotājiem, bet arī tiem, kas apgūst valodu kā svešvalodu vai otro valodu.

Latviešu valodā pareizrakstības normas nosaka: 1) morfēmu rakstību vārdos un vārdformās (t. s. īpašvārdu rakstību), 2) vārdu rakstību kopā vai šķirti, 3) defisrakstību, 4) lielā sākumburta lietojumu nosaukumos, 5) vārdu pārnesumpārdali, 6) vārdu saīsināšanu. (Strautniece, Šulce 2009, 50)

Visu latviešu valodas prasmes līmeņu darbos izplatītākās ir pareizrakstības kļūdas: A līmenī – 31%, B līmenī – 37%, C līmenī – 38%. Aplūkojot valodas apguvēju pieļautās pareizrakstības kļūdas, tās iespējams grupēt šādās apakšgrupās:

- 1) patskaņa garumam neatbilstoša burta izmantojums:
 - a. garais patskanis tiek rakstīts ar īsa patskaņa burtu,
 - b. īsais patskanis tiek rakstīts ar gara patskaņa burtu;
- 2) līdzskaņu rakstība bez diakritiskajām zīmēm vai neatbilstošas diakritiskās zīmes lietojums;
- 3) divskaņu rakstība;
- 4) izlaisti burti;
- 5) lieki burti;
- 6) kopā vai šķirti rakstāmi vārdi;
- 7) lielo sākumburtu lietojums;
- 8) pārstatīti burti;
- 9) citvalodu īpašvārdu atveide.

No pareizrakstības kļūdām izplatītākās ir diakritisko zīmju trūkums vai pārdaudzums (1.–3. apakšgrupa), patskaņu un divskaņu šķīrums, īpašvārdu atveide (7.–8. apakšgrupa).

Kļūdu analīze

Latviešu valodas rakstības pamatā ir latīņu alfabēts, kas papildināts ar diakritiskām zīmēm patskaņu garuma (˘), palatālo līdzskaņu (.), šņāceņu (̃) apzīmēšanai. (Vanags 2018) Tieši šo papildu diakritisko zīmju kļūdaina izmantošana vai neesamība ir biežākās pareizrakstības kļūdas, kuru cēlonis ir neatbilstoša vārdu izruna. Gan patskaņa garumam neatbilstoša burta izmantojums (gara patskaņa vietā tiek rakstīts īsa patskaņa burts, savukārt īsa patskaņa vietā tiek rakstīts gara patskaņa burts), gan līdzskaņu rakstība bez diakritiskajām zīmēm vai neatbilstošas diakritiskās izmantojums, bieži vien ir saistīts ar valodas apguvēja dzimtās valodas fonētiski fonoloģiskās sistēmas un / vai grafētikas atšķirību no latviešu valodas.

Līdzskaņu rakstība bez diakritiskajām zīmēm

Latviešu grafētikā vienkārša grafēma un diakritiskā zīme veido līdzskaņu burtus *k, ģ, l, ņ, š, ž, č*, savukārt ar salikto grafēmu un diakritisko zīmi apzīmē afrikātu *dž*. Valodas apguvēji iespējams tīri mehāniski aizmirst rakstot pievienot diakritisko zīmi, bet varbūt tomēr tas ir saistīts ar viņu dzimtās valodas īpatnību, piemēram, palatālo līdzskaņu *k, ģ, ņ, l* neesamību viņa dzimtās valodas fonētiski fonoloģiskajā sistēmā un nespēju šo fonēmu saklausīt un izrunāt, kas savukārt traucē šo skaņu pareizi atspoguļot rakstos.

Visbiežāk valodas apguvēju darbos kļūdaini atveidots līdzskanis *ņ*, tad līdzskanis *l*, līdzskanis *s*, līdzskanis *ž*, tiem seko *ģ* un *č*, piemēram,

- līdzskanis *n* vietā rakstīts līdzskaņa *ņ* vietā: *nemtu* (*ņemtu*)¹, *atmina* (*atmiņa*), *devini* (*deviņi*), *septini* (*septiņi*), *viniem* (*viņiem*), *druscin* (*drusciņ*), *varoniem* (*varoņiem*);
- līdzskanis *l* vietā rakstīts līdzskaņa *ļ* vietā: *daleji* (*daļēji*), *bileti* (*biļeti*), *mīlais* (*mīļais*), *nedēla* (*nedēļa*), *atpakaļ* (*atpakaļ*), *mirklus* (*mirkļus*), *klūdu* (*kļūdu*);
- līdzskanis *s* vietā rakstīts līdzskaņa *š* vietā: *so* (*šo*), *seit* (*šeit*); *mēnesa* (*mēneša*), *vīriesi* (*vīrieši*); *rakstīsana* (*rakstīšana*), *domāsanu* (*domāšanu*); *trēskārt* (*treškārt*); *kurs* (*kurš*);
- līdzskanis *z* vietā rakstīts līdzskaņa *ž* vietā: *biezi* (*bieži*), *visbiežāk* (*visbiežāk*), *dazreiz* (*dažreiz*), *daziem* (*dažiem*), *izstāzu* (*izstāžu*), *sēz* (*sēž*), *mēzu* (*mežu*), *mezā* (*mežā*), *dazados* (*dažādos*).

Problēmas valodas apguvējiem sagādā arī līdzskaņu savienojumu *šņ*, *šķ* rakstība vārdu sākumā vai vidū, piemēram, *krašni* (*krāšņi*), *atsķiribu* (*atšķirību*), *skirošanu* (*šķirošanu*). Samērā bieži kļūdaini tiek rakstīti arī divlīdzskaņu savienojumi vārdu beigās, kur viens vai abi līdzskaņi rakstāmi ar diakritisko zīmi un otrais ir līdzskanis *š* (galotne), piemēram, *vinš* (*viņš*), *ceļš* (*ceļš*).

Patskaņa kvantitātei neatbilstoša burta izmantojums

Latviešu valodā pēc garuma izšķir īsus un garus patskaņus. Īso patskaņu izrunas laiks ir vidēji 2,3 līdz 2,5 reizes īsāki nekā garie patskaņi. Precīza skaitliskā attiecība ir atkarīga gan no patskaņu artikulārās kvalitātes, gan uzsvara, gan zilbes intonācijas. (Nītiņa, Grigorjevs 2013, 41) Turklāt garā un īsā patskaņa šķīrums tiek saglabāts gan uzsvērtās, gan neuzsvērtās zilbēs. lielākajai daļai valodas apguvēju, kas kārtoja valsts valodas pārbaudes eksāmenu un kuru darbi tika analizēti, dzimtā valoda ir kāda no austrumslāvu valodām, kurās netiek šķirti īsi un gari patskaņi, kā arī nav fiksēta pirmās zilbes uzsvara. Tādēļ pirms datu analīzes tika pieņemts, ka kļūdaini būs rakstīts īsais patskanis uzsvērtā zilbē, cenšoties ar garumzīmi parādīt tā intensīvāku un, iespējams, arī garāku izrunu, bet garš patskanis lielākajā daļā gadījums būs atspoguļots pareizi. Tomēr datu analīze liecina, ka tā nebūt nav – bieži vien gan uzsvērtā, gan neuzsvērtā zilbē garais patskanis tiek atspoguļots nepareizi, t. i., ar atbilstošas kvalitātes īsa patskaņa burtu.

Garš patskanis tiek rakstīts ar atbilstošu īsa patskaņa burtu:

- saknē, uzsvērtā zilbē: *bus* (*būs*), *atri* (*ātri*), *velos* (*vēlos*);
- saknē, neuzsvērtā zilbē: *izveleties* (*izvēlēties*), *aplukojot* (*aplūkojot*);
- piedēklī: *domaju* (*domāju*), *apmeklet* (*apmeklēt*), *veselība* (*veselība*), *mazak* (*mazāk*);
- galotnē: *Es dzīvoju Rīga*. → *Es dzīvoju Rīgā*.

¹ Šeit un turpmāk aiz kļūdainā vārda vai vārdformas iekavās dots pareizais variants.

Protams, valodas apguvēju darbos vērojama arī pretēja parādība – īss patskanis gan uzsvērtā, gan neuzsvērtā zilbē tiek rakstīts ar gara patskaņa burtu:

- uzsvērtā zilbē vienzilbes vārdos: *tū* (*tu*), *sēn* (*sen*), *mūms* (*nums*), *bēt* (*bet*).
- neuzsvērtā zilbē, divu un vairāk zilbju vārdos: *šovakār* (*šovakar*, *protāms* (*protams*));
- neuzsvērtā gala zilbē, galotnē: *viņī* (*viņi*), *pasākumū* (*pasākumu*).

Vairākkārt īsa vai gara latviešu valodas patskaņa atveidē tiek izmantots divskanis:

- garu patskani atveido ar divskani:
 - *ē* → *ei*: *bleīndari* (*blēndari*), *tapeic* (*tāpēc*);
 - *ē* → *ie*: *lietāki* (*lētāki*), *saņiemām* (*saņēmām*);
 - *ī* → *ie*: *piekrietu* (*piekrītu*), *vierieši* (*vīrieši*), *patiek* (*patīk*);
- īsu patskani atveido ar divskani:
 - *e* → *ie*: *jacienšas* (*jācenšas*), *aizņiemts* (*aizņemts*);
 - *i* → *ie*: *nopierkšu* (*nopīrkšu*).

Kļūdaina divskaņu atveide rakstos

Divskaņos *ai*, *ei* otrais komponents tiek rakstīts ar līdzskani *j*, iespējams, atspoguļojot izrunā vērojamo konsonantizāciju gan vaļējā, gan slēgtā zilbē, piemēram, *tikaj* (*tikai*), *šaj* (*šai*), *laj* (*lai*), *vaj* (*vai*); *sabiedriskajs* (*sabiedriskais*), *labākajš* (*labākais*). Šāda divskaņu *ai*, *ei* rakstība vērojama arī atsevišķu vārdu saknē, piemēram, *atvajnoet* (*atvainojiet*), *šejt* (*šeit*), *skajsta* (*skaista*).

Latviešu valodā, īsam patskanim *a*, *e* nonākot vienā zilbē ar līdzskani *j* vai *v*, līdzskanis vokalizējas, un veidojas pozicionālais divskanis, kas rakstos netiek atspoguļots. Tomēr valodas apguvēju tekstos reizēm pozicionālais divskanis tiek rakstīts atbilstoši tā izrunai:

- *aj* vietā kļūdaini tiek rakstīts *ai*, bet *ej* – *ei*: *tramvais* (*tramvajs*); *kafeinicā* (*kafejnīcā*), *voleibols* (*volejbols*);
- *av* vietā kļūdaini tiek rakstīts *au*, bet *ev* – *eu*: *nau* (*nav*), *teu* (*tev*).

Kā jau iepriekš minēts, lielākajai daļai valodas apguvēju, kuru darbi tiek analizēti, dzimtā valoda ir krievu, ukraiņu vai baltkrievu valoda. Šajās valodās divskaņu *nav*. Iespējams, tādēļ valodas apguvēji mēģina atspoguļo divskaņus ar kādu patskaņa burtu:

- *ie* → *e*: *jaleto* (*jālieto*), *ievešot* (*ieviešot*), *apcemoja* (*apciemojā*), *šoden* (*šodien*), *labden* (*labdien*), *ven* (*vien*), *draugem* (*draugiem*), *viņem* (*viņiem*);
- *ie* → *i*: *vins* (*viens*), *jaunišu* (*jauniešu*), *labdin* (*labdien*);
- *ie* → *ē*: *vētējā* (*vietējā*);
- *au* → *a*: *aizratu* (*aizrautu*);
- *au* → *ā*: *aptājā* (*aptaujā*);

- *ei* → *e*: *patecību* (*pateicību*), *teks* (*teiks*).

Reizēm viens no divskaņa komponentiem kļūdaini tiek rakstīts ar garā patskaņa burtu, piemēram, *drāudzību* (*draudzību*), *jāunas* (*jaunas*), *vāirāki* (*vairāki*), *parēizi* (*pareizi*), *smāida* (*smaida*), *ārzemnieki* (*ārzemnieki*); *taūtu* (*tautu*), *staīgāt* (*staigāt*), *aīcīnu* (*aicīnu*), *aiziēt* (*aiziet*).

Tikai dažas reizes, vēloties atspoguļot divskaņa [uo] izrunu, valodas apguvēju tekstos tas rakstīts kā *uo*, arī *ua*: *pirmaja un uotraja augusta* (*pirmajā un otrajā augustā*), *nuo* (*no*), *izmantu* (*izmanto*), *nuomaksā* (*nomaksā*), *kaut kuo parduat* (*kaut ko pārdot*). Turklāt šo darbu autoru dzimtā valoda ir lietuviešu, kurā divskani [uo] raksta ar burtu kopu *uo*.

Bez iepriekš minētajām kļūdām vērojamas arī citas:

- latīņu alfabēta un kirilicas rakstzīmju jaukšana: *mu* (*tu*), *līdzu* (*līdzi*);
- kopā vai šķirti rakstāmi vārdi: *Uzredzēšanos!* (*Uz redzēšanos!*), *takā* (*tā kā*);
- lieki burti: *rakstinieku* (*rakstnieku*), *kollekcionēt* (*kolekcionēt*), *dzintaras* (*dzintara*), *gruppas* (*grupas*), *vissvairāk* (*visvairāk*);
- izlaisti burti: *kau* (*kaut*), *vis* (*viss*), *pieju* (*pieeju*), *apuseju* (*abpusēju*), *ļot* (*ļoti*), *nopirk* (*nopirka*), *ļau* (*ļauj*);
- pārstatīti burti:
 - divskaņos *ei* un *ie*: *neveina* (*neviena*), *atteicas* (*attiecas*),
 - līdzskaņu savienojumos: *dāzrs* (*dārzs*), *kartā* (*katrā*), *ārtsi* (*ārsti*);
- lielo sākumburtu lietojums un īpašvārdu rakstība.

Runājot par lielo sākumburtu lietojumu, visbiežāk valodas apguvēji kļūdījušies vēstulē (viens no A un B līmeņa pārbaudes darbu uzdevumiem bija uzrakstīt vēstuli) nerakstot vietniekvārdu *Tu*, *Jūs* un *Tavs* formas ar lielo sākumburtu. Reizēm tekstos teikuma pirmais vārds rakstīts ar mazo burtu, piemēram, *man ir labi*.

Secinājumi

Pareizrakstības kļūdas atspoguļo valodas apguvēja spēju izrunāt un saklausīt latviešu valodas skaņas un skaņu savienojumus.

Analizējot kļūdas, jāsecina, ka visos valodas apguves līmeņos visizplatītākās ir pareizrakstības kļūdas (37% no visām kļūdām). Savukārt pareizrakstības kļūdu vidū visbiežāk vērojams kļūdainis diakritisko zīmju lietojums, arī diakritisko zīmju neesamība, kas saistīta ar 1) īso un garo patskaņu šķīrumu rakstos, 2) palatālo līdzskaņu *ļ, ņ, ķ, ģ*, afrikātu *č, dž*, apikāli alveolāros priekšējo mēleņu *š, ž* atveidi rakstos.

Kā jau rakstā minēts – lielākajai daļai valodas apguvēju dzimtā valoda ir viena no austrumslāvu valodām, līdz ar to pareizrakstības kļūdas galvenokārt saistītas ar viņu dzimtās valodas – krievu, ukraiņu vai baltkrievu valodas –

fonētiski fonoloģiskās sistēmas atšķirību no latviešu valodas, kurā tiek šķirti garie un īsie patskaņi, ir palatālie līdzskaņi, kā arī pirmās zilbes uzsvars. Veidojot latviešu valodas apguvēju korpusu, kurā tiek iekļauti to ārzemju studentu darbi, kas Latvijas augstskolās apgūst latviešu valodu, jau sākotnējā datu analīze parāda, ka pareizrakstības kļūdas atšķiras.

Kļūdu analīze parāda latviešu valodas apguves problemātiku un ir mācību un metodisko materiālu izveides pamats.

Literatūra

Darģis et al. 2018 – **Darģis, Roberts, Auziņa, Ilze, Levāne-Petrova, Kristīne**. The Use of Text Alignment in Semi-Automatic Error Analysis: Use Case in the Development of the Corpus of the Latvian Language Learners. *Proceedings of the 10th International Conference on Language Resources and Evaluation (LREC 2018)*, 2018 Miyazaki, Japan: European Language Resources Association (ELRA)

Ellis 1994 – **Ellis, Rod**. *The Study of Second Language Acquisition*. Oxford University Press, 1994.

Nītiņa, Daina. Grigorjevs, Juris 2013 – **Nītiņa, Daina. Grigorjevs, Juris** (red.). *Latviešu valodas gramatika*. Rīga: LU Latviešu valodas institūts, 2013.

Strautiņa, Vaira, Šulce, Dzintra 2009 – **Strautiņa, Vaira, Šulce, Dzintra**. *Latviešu valodas pareizrūna un pareizrakstība*. Rīga : RaKa, 2009.

Vanags 2019 – **Vanags, Pēteris**. Latviešu valoda. *Nacionālā enciklopēdija*. <https://enciklopedija.lv/skirklis/9891> (skatīts 06.05.2019)

Znotiņa 2017 – **Znotiņa, Inga**. Otrās baltu valodas apguvēju korpuss: izveides metodoloģija un lietojuma iespējas: promocijas darbs filoloģijas doktora grāda iegūšanai valodniecības zinātņu nozares lietišķās valodniecības apakšnozarē. Liepāja: Liepājas Universitāte, LiePA, 2017.