



Inguna Skadiņa un Andrejs Spektors

Latviešu valoda datorā — pētījumi, resursi, tehnoloģijas

Zinātniskā konference "Matemātika un informātika pēc 50..."
2009. gada 9. novembris

```
01011100110110101 101000110100
10110100111010100 1011011010111001
110111110111010001 101100010011011100
101001 100110 0011001
011001 100110 010001
010111 1010100 1101110
011101 0100000 1100001 0111000
001110010011000010 1110011 0110100
10010000001001100010 1010101 0011010
1001001 01001001 1101000 1000110
0101101 0111010 0011011
0010110 1001101 0100111
1010011 0110010 1100001
1010001 0000001 1101100
0010000 001110 0100110 0010001
1010100 1110001 1001100 1000110
00000011 10010011 11100110101001110110
000010110111101100 1101100011110010
11101100101011 110110100100
```



Pirmsākumi

- Vēlme lietot datoru valodas apstrādei ir tikpat veca kā pirmais dators
- 60. gadu sākumā A.Gobzemis, V.Gorobecs, V.Juriks, T.Jakubaite izstrādā tiešās mašīntulkošanas sistēmu zinātnisko rakstu tulkošanai no krievu valodas latviešu valodā
- Valodas analīze, izmantojot kvantitatīvās metodes, uzsākta 60.gados (T.Jakubaite, S.Kļaviņa, A.Lorencs, Z.Nesaule)



Pirmsākumi

- Plašs kvantitatīvo metožu lietojums pētniecībā vērojams 70.gados (V.Drīzule, T.Jakubaite, S.Kļaviņa, V.Kuzina, N.Mecs, M.Oša, A.Rubīna, N.Sika, E.Soida, R.Zarovska)
- Veikta valodas kvantitatīvā analīze (vārdu krājuma sastāvs dažādu funkcionālo stilu tekstiem, vārdšķiru statistiskais raksturojums u.c.)
- Izstrādāta “Latviešu valodas biežuma vārdnīca” un “Latviešu valodas inversā biežuma vārdnīca”



Mākslīgā intelekta laboratorijas rašanās

- Tautasdziesmu tekstu ievade datorā
- 1988. gadā MII Lietišķās informātikas daļa uzsāk datorizētu latviešu valodas pētniecību:
 - Pētītas latviešu valodas kodēšanas iespējas, 1990. gadā izveidots kodu tabulas standarts (G. Lučkins un L. Vasermans)
 - Seno tekstu pētniecība (M.Baltiņa, E.Milčonoka, A.Ozoliņa)
 - Mūsdienu latviešu valodas pētniecība (A.Spektors, I.Āboliņa, R.Čevere, K.Dancis, I.Greitāne, B.Krauze-Krūze, U.Sarkans)
- 1992. gadā izveidota Mākslīgā intelekta laboratorija, kuru no tās pirmsākumiem līdz mūsdienām vada Dr. fiz. Andrejs Spektors

AI
lab



Galvenie darba virzieni

- **Korpuslingvistika** (M.Baltiņa, A.Spektors, E.Andronova, K.Levāne-Petrova, G.Nešpore, V. Krugļevskis, S. Reinsone)
- **Mašīnlasāmu vārdnīcu izveide** (A. Spektors, E.Andronova, G.Nešpore, N.Grūzītis, I. Ilziņa, V. Rostoks)
- **Programmriki valodas automatizētai apstrādei** (N. Grūzītis, I. Skadiņa, P. Paikens, M. Virza, V. Krugļevskis, I. Poikāns)
- **Mašīntulkošana** (I.Skadiņa, E. Andronova, E. Brālītis, M. Virza)
- **Semantiskā tīmekļa tehnoloģijas** (G.Bārzdiņš, A. Spektors, G.Nešpore, N.Grūzītis, B. Valkovska, J. Džeriņš)
- **Runas tehnoloģijas** (I.Auziņa, M. Pinnis, R. Berzinskis, G. Rābante)
- **Mācīblīdzekļi** (A.Spektors, E.Andronova, K.Levāne-Petrova, G. Nešpore, I.Auziņa, N. Grūzītis)

AI
lab



Korpuslingvistika

- Tekstu uzkrāšana sāka 80. gadu beigās un 90. gados, kad tekstus ievadīja datorā. Kopš 90. gadu vidus teksti tiek skenēti un pēc tam pārbaudīti, skenēšanas precizitāte ir 98.5%
- Nozīmīgākās tekstu kolekcijas:
 - Latviešu literatūras klasika (3, 5 milj. vārdlietojumu)
 - Latviešu teiku un pasaku, ticējumu un sakāmvārdu datorfonds
 - “Rīgas Balss” latviešu un krievu valodās
- Vairāki paralēlie tekstu korpusi:
 - Sastatīts Orvela ”1984”
 - Platona “Valsts” tulkojums
 - Sadarbībā ar TTC veidots juridisko tekstu paralēlais korpus

AI
lab



Seno tekstu korpuss

- Seno tekstu ievadīšana uzsākta 80. gadu un turpinās līdz pat mūsdienām:
 - **1989.g.** — M. Baltiņa ierodas LU MII, lai ievadītu Ījaba grāmatu datorā
 - **1992.–1994. g.** — M. Baltiņas vadībā datorā sāka ievadīt visu E.Glika vadībā tulkoto Bībeli
 -
 - Kopš **2002. g.** Seno tekstu korpuss (<http://www.korpuss.lv/senie/>)

to layke / kad Kyrenius Semmesoge exkan Syrten
by. Vnd Ickwens nogana / ka the höw mheslotes
lichte / nckwens exkan souwe Pille.

Ead nogan arridtezan Joseph no Galilea / aran
tās Pilles Nazareth / exkan to Judde semme / vs to
Pille David / kattrā tur dhēween Bethlehem / Eas
pectez ka tas no to Namme vnd Kadde David by /
vs to / ka tās höw mheslote lichte / ar Maria souwe
saloulata Gaspasche / kattrā tur apgrutenata by.
Vnd kad the tur path by / nātee tas laix ka thay pe
ezimpt by. Vnd tha peezimme souwe pirmone Dhe
le / vnd tinne to exkan autims / vnd lichte to exkan we
ne Szille. Aesto themis nhe by ezittur neewena
weta exkan to Maiawete.

Vnd tur by Ganne exkan to patte wete / tur
prettlybe wuerßon to louke / te ganne py to Nakte
souwes lopes. Vnd Rouge / tas Engels tha fun
ge apspyden tōs / vnd the issabhas höw lote / Vnd
tas Engels haten vs tems / nhe nßabhates yums /



Seno tekstu korpus

- Pašlaik korpus ietver 16. gs., 17. gs., kā arī vairākus 18. gs. sākuma tekstus latviešu valodā, kopā ap 1 miljonu vārdlietojumu
- Izstrādāti korpusam nepieciešamie meklēšanas, statistikas un konkordances rīki
- **2002. g.** izveidota darba grupa Seno tekstu korpusa un “Latviešu valodas vēsturiskās vārdnīcas (16.-18.gs.)” izveidei (<http://www.tezaurs.lv/lvvv/>)

SENIE

latviešu valodas seno tekstu korpuss

Konkordances rezultāts

KONKORDANCE

SĀKUMLAPA

Vārdformas šablons: **naud%**

Avots: SENIE

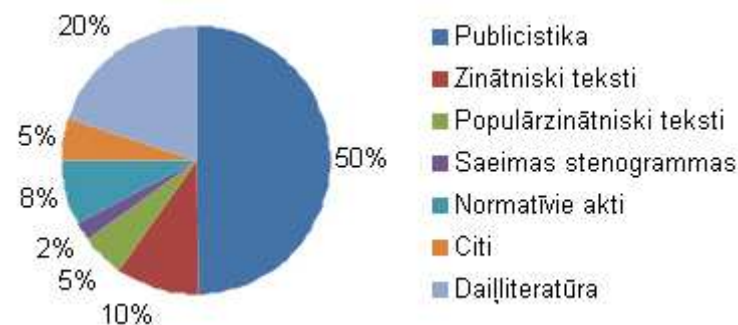
Statistika: vārdformas - 26, vārdlietojumi - 166

pirzeeta beß	<u>Naudas</u>	und wellte tick lieds Peemu ka Wienu Nahzeeta
nhe ka dauds	<u>Naudas</u>	par Manntu kraht
tō wairs	<u>Naudas</u>	nhe gir Macka` us Addošchanu jajämm
šeeta nahzeeta šchurr py Vhdeni und juhß kattreem	<u>Naudas</u>	nhe gir nahzeeta šchurr pirzeeteß und ehdeeta
nahzeeta šchurr pirzeeteß beß	<u>Naudas</u>	und wellte
wairahk	<u>Naudas</u>	dabbuit warr Winjam buhß py Christi Eenaidneekeem
Širrdi us	<u>Naudas</u>	dšiedammeß Buhß tah Mahzekļam
Mehteli kaß jauns buhdams dauds	<u>Naudas</u>	maxajis bett nu
m Zilwehkam 80 Dalderus tas irr 40 Dalderus labbas	<u>Naudas</u>	tur klaht tam Wihrišchkam šešchi reiš zaur Rihkšt
Zilwehkam 80 Dahlderus tas irr 40 Dahlderus labbas	<u>Naudas</u>	tur klaht tam Wihrišchkam šešchi reiš zaur Rihkšt
pirzetees Ghuddribu kad juhs beß	<u>Naudas</u>	to warraht dabbuit
mannas	<u>naudas</u>	sinnaischi kahd zits
Tam tahs	<u>naudas</u>	suit
ohtra` bahršta Bett nu tam Judas wišši	<u>Naudas=Ghabbali</u>	par Wirrwehm ißdohdahß par ko wings pirmahk pree
ja taß	<u>Naudas=ghabbalings</u>	winjai pašuhd redši ka winja
		Schähl gir dhe ka kahdam Baggatam kad



Mūsdienu latviešu valodas korpuss

- Pašlaik uzkrāti teksti ar vairāk nekā 100 milj. vārdlietojumu
- Izveidots līdzsvarots 3, 5 miljonus vārdlietojumu liels mūsdienu latviešu valodas tekstu korpuss (www.korpuss.lv)



- Morfoloģiski marķēts korpuss (~30 tūkst. vārdlietojumu)
- Saeimas stenogrammu korpuss
- Latvijas tīmekļa korpuss



launs vaicājums

nosaukums:

miljons-1.0

lanet . lv 2002 . gada karstā un sausā **vasara** mūs atstāja bez sēnēm vasarā un arī
 . " paskaidroja Kaspars , kuram tā gada **vasara** pagāja atbilstoša auto meklējumos .
 — — diena , nakts , diena , nakts , **vasara** , ziema — — mums tikai jāslīd pa
 nekā būtiska sakāma nav . * * * Bija **vasara** , bija karsts , un bija daudz laika . Katru
 Varbūt . Tam nebūtu tā jābeidzas . Bija **vasara** , bija karsts , es gāju pa ielu , ātri
 vannā ūdeni . Laukā ir silts laiks , **vasara** pilnbriedā . Gatis laiž vannā ūdeni
 visviens , zivtiņās iet — ziema tā vai **vasara** — zaļi svārki , tie paša brezenta
 līdz pavasarim . Un vecāki uzskatīja , ka **vasara** man — drošs paliek nedrošs — jāpavada
 izslēdzu projektoru un dodos laukā . Ir **vasara** . Saule spīd cauri koku galotnēm spoži
 kreklī ar trim podziņām . Ir 1974 . gada **vasara** . No jūras vējš atnes krievu raidstaciju
 un un dodos uz staciju . VASARA Ziemeļu **vasara** allaž ir īsa un nepastāvīga , allaž
 atvadas . — — Tu jautāji , kāda būs **vasara** . Tā tu jautāji dienā , kad zeme , pie
 malā klausās un zina — šī ir miera **vasara** . Bezšaubu laiks . Puikas brūnā mugura
 par daudz palaist muti un lekties jā , **vasara** tikko sākusies , citi plāno atvaļinājumus

Atrasto vārdlietojumu skaits: 24

> Query : "vasara"

Parādīts: viss/24 Rindiņa: 10 Iezīmēts: 1

launs vaicājums

nosaukums:

ledus

nevarēsi meitene atkārtoja . varēšu gan ,
 kabatas rēgojās adītas cepures stūris .
 par evenku , mongoli vai ķīnieti , taču
 , tev būs baīl . - man nemaz nav baīl ,
 irietis . meitene bija galvas tiesu garāka .
 plāns . no āgenskalna liča kreisās malas
 esot bezpajumtnieki un žurkas . neviens
 grauzās sodrēji , putekli un dubli .
 tomēr pārāk lēni . severīnam likās , ka
 kisjai trūka drosmes . pietika ar to , ka
 meitene tur stāv un priecājas , ka izdevies
 pārskriet pāri licim . varbūt tagad

viņš/viņš/Pp3msn
viņš/viņš/Pp3msn
viņi/viņš/Pp3npl
viņš/viņš/Pp3msn
viņai/viņa/Pp3fsd
viņos/viņš/Pp3mpl
viņiem/viņš/Pp3mpd
viņa/viņš/Pp3msg
viņš/viņš/Pp3msn
viņš/viņš/Pp3msn
viņu/viņš/Pp3msa
viņa/viņa/Pp3fsn

iebilda . kalsnējs , melnmatains zēns
 to pikti iestūķēja dziļāk (nemaz
 sarunājās latviski . - nevarēsi , tev
 sabozās . - ir gan , es redzu . - pašai
 mugurā bija tumši zils mētelītis ar
 vienaldzīgi noraudzījās trīsdesmitstāvi
 nepievērsa uzmanību : nedz apkārtējā
 ceļu krustoja vairākas plaisas , gar
 skrien uz vietas , taču atskatīties un
 zināja meitene tur stāv un priecājas
 piedabūt pārskriet pāri licim . varbūt
 jau sāka nožēlot severīns taču var



Mašīnlasāmu vārdnīcu izveide

- Esošo latviešu valodas vārdnīcu digitalizēšana un konvertēšana mašīnlasāmā formā
- Mērķis – izveidot visaptverošu vārdnīcu datubāzi, kas ietvertu pēc iespējas vairāk latviešu valodas vārdu un to nozīmju aprakstu, ko var izmantot gan galalietotāji, gan automatizēti valodas apstrādes rīki



Internetā publicētās vārdnīcas

- Skaidrojošā vārdnīca (~140 000 šķirkļu, <http://www.tezaurs.lv/sv/>)
- Latviešu literārās valodas vārdnīca (~64 000 šķirkļu, <http://www.tezaurs.lv/lvv/>)
- Daudzvalodu terminoloģijas datu bāze TTC vajadzībām (~115 000 šķirkļu)
- Latviešu valodas vārdnīca (~30 000 šķirkļu, <http://www.tezaurs.lv/lvv/>)
- Mīlenbaha-Endzelīna vārdnīca (~77 000 + ~55 000 šķirkļu, <http://www.ailab.lv/mev/>)



LLVV

http://www.tezaurs.lv/llvv/ Google

LETONIKA Vārds:

~64 000 šķirķļu
beta versija
22.10.2009.

meklēt vārdus, kas sākas ar...
 piedāvāt arī ortogrāfiski līdzīgos vārdus
 nerādīt citātus

Latviešu literārās valodas vārdnīca
Projektu finansē Valsts Pētījumu programma "Letonika"
Iestrādes atbalstījis Valsts Kultūrkapitāla fonds
© LU [Matemātikas un informātikas institūts](#)
© LU Latviešu valodas institūts
© Apgāds "Zinātne"

Mobilā versija: www.tezaurs.lv/wap

[plāns¹](#)
[plāns²](#)
[plāns³](#)

[plānsaimniecība](#)
[plānsienū](#)

Ortogrāfiski līdzīgie vārdi:

[lāns](#)
[pāns¹](#)
[pāns²](#)
[klāns](#)
[klāns](#)
[plakš](#)
[plakš](#)
[plašs](#)
[ulāns](#)
[plats](#)
[plāts](#)
[plands](#)

plāns³ [pla~ns] -a, v.; apv.

Klons. Arī grīda.

Rijas plāns.

Virtuves plāns.

Kambarītim nebij grīdas, tāpat kā istabai; bet plāns te izskatījās gludāks, tīrāks, sausāks, jo vistas iekšā. *Jaunsudrabiņš 3, 17.*

Klukšķe [vista] sēdēja un sēdēja savā cisu grozā tumšajā kaktā zem gultas, tikai reizū reizēm iznāca plānā drusku iekost un padzert, tad steidzās atkal atpakaļ perēt. *Birznieks-Uptis IV, 95.*

Ja man vēl šodien istabai kakti jāizslauka tikpat tīri kā plāna vidus, tad to man tika mācījusi pama Klints 1, 33.

Pirmajā pāri plāna vidū izgāja Subergs ar Barbu, veikli uzņemot dejas soli.. *Dorbe 6, 13.*

Jautājumiem, komentāriem, atsauksmēm: [tezaurs @ ailab punkts lv](mailto:tezaurs@ailab.punkts.lv)

Pamanījāt šajā šķirķlī kļūdu?

Display a menu



www.tezaurs.lv/llvv, www.tezaurs.lv/wap



Programmāriki valodas apstrādei

- Morfoloģiskās analīzes un vārda formu ģenerēšanas programmas
- Automātiskas morfoloģiskās marķēšanas rīki

pirmā <Mosfsny> pirmā

nodaļa <Ncfsn4> nodaļa

Atmosfēra <Ncfsn4> atmosfēra

tu <Pp2nsn> tu

to <Pdnmsa> tas

nevarēsi <Vmnipti32say> nevarēt

meitene <Ncfsn5> meitene

atkārtoja <Vmnipti23san> atkārtot

.

AI
lab



Programmāriki valodas apstrādei

- Morfēmiskās analīzes rīki
- Rīki terminoloģijas izguvei
- Sintaktiskās analīzes modeļi un rīki



Mašīntulkošana

- 1994.g. LU MII Mākslīgā intelekta laboratorijā tiek uzsākta interlingvas tulkošanas sistēmas modeļa LATRA izstrāde

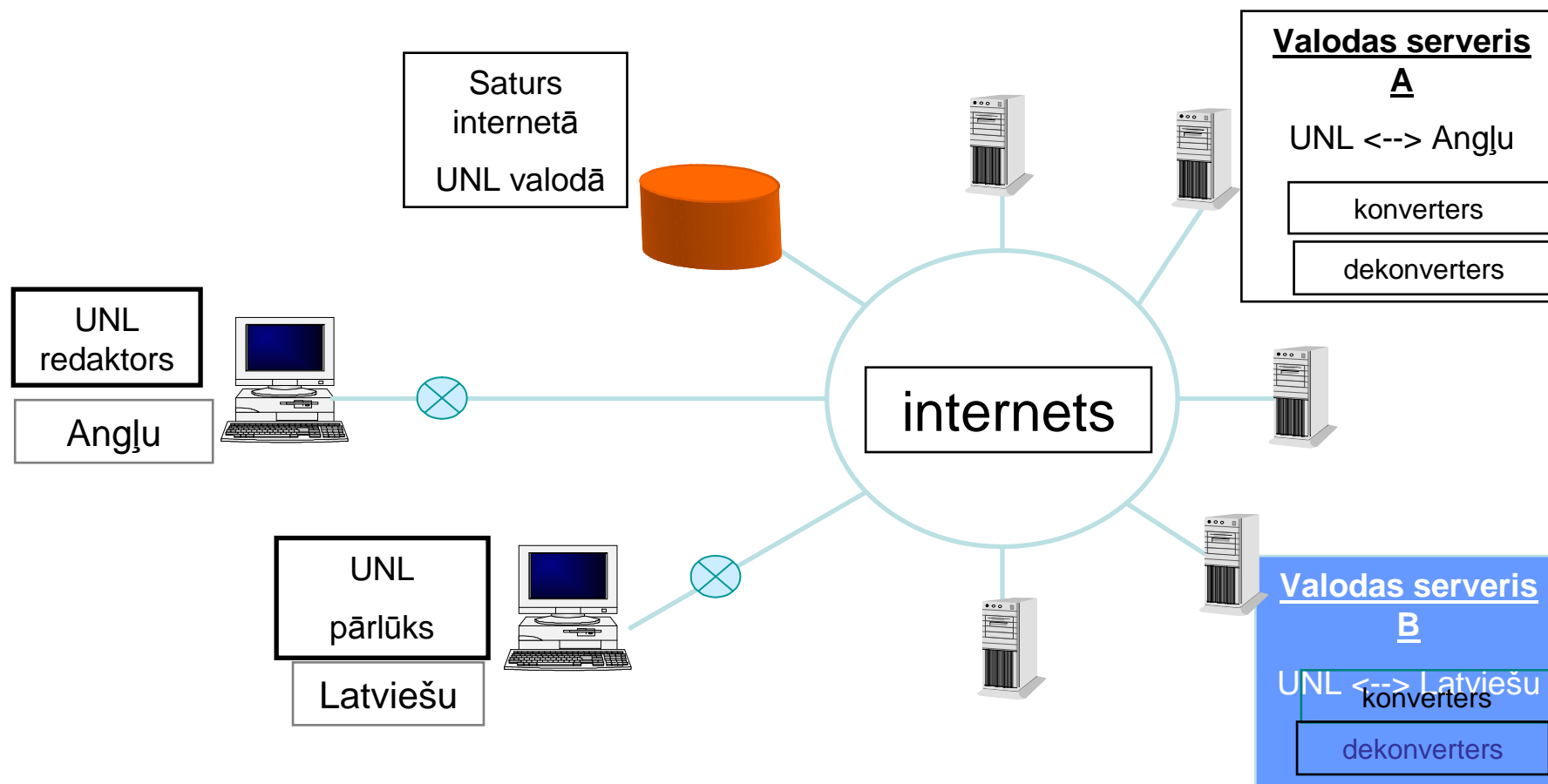
Gada sākumu iezīmēja arī struktūrkapitāla palielināšanas maratons, kuru ir spiesta uzsākt komercbanku lielākā daļa.

The marathon of the increase of statutory capital which the most of the commercial bank must begin characterised the beginning of the year also.

- 1997. g. LU MII Mākslīgā intelekta laboratorija iesaistās Apvienoto Nāciju universitātes projektā *Universal Networking Language*

AI
lab

Universālā tīkla valoda





Statistiskā mašīntulkošana

Pētījumi statistiskajā mašīntulkošanā uzsākti 2005. gadā LZP projekta „Statistisko metožu izvērtējums angļu-latviešu tulkošanas sistēmā” ietvaros un turpinās projektā “Faktorēto metožu lietojums angļu-latviešu statistiskajā mašīntulkošanas sistēmā”

Tulkojamais angļu val. teksts(max 400 simboli): This appropriation is intended to cover vehicle maintenance and operating costs and costs relating to the use of public transport.	Tulkot!	Tulkojums latviešu valodā: Šis apropriācijas ir paredzēts , lai segtu transportlīdzekļu tehniskās apkopes un ekspluatācijas izmaksas un izmaksas , kas attiecas uz sabiedriskā transporta lietošana .
--	----------------	---

eksperimenti.ailab.lv/smt






Runātās valodas resursi

- “ONOMASTICA-COPERNICUS” (1995-1997) projekta ietvaros transkribēti ~250 000 īpašvārdu
- 2001. gadā uzsākta runas korpusa izveide:
 - 15 runātāju ierunātas frāzes (1 300)
 - 8 stundu garš semināra ieraksts
 - 50 runātāju ierunāts teksts (1000 vārdlietojumu)
 - Teksti ir transkribēti (t.i., pierakstīti mašīnlasāmā formā) un segmentēti, izmantojot *Transcriber* un *WaveSurfer* programmatūru



Runas tehnoloģijas

- Izstrādāta grafēmu-fonēmu pārveides programmatūra un likumi. Automātiskās transkripcijas precizitāte ir 92%
- Izstrādāts eksperimentāls runas sintezators 
- Ir izstrādāts eksperimentāls runas analīzes modulis, kas spēj atpazīt aptuveni 80 vārdus

AI
lab



Latviešu valodas mācīblīdzekļi

<http://valoda.ailab.lv/latval/>

- Vispārējas ziņas par latviešu valodu
- Latviešu valodas multimediju mācīblīdzeklis iesācējiem
- Latviešu valoda sākumskolai
- Latviešu valoda pamatskolai
- Latviešu valoda vidusskolai
- Latviešu valoda bilingvālai apmācībai
- Latviešu valoda bērniem ar dzirdes traucējumiem
- Latviešu valodas vārdu analizators un sintezators



VALODA

SKAŅU MĀCĪBA

MORFOLOĢIJA

TENZĪMA MĀCĪBA

LEKSIKOĢIJA





CLARIN - Vienota valodas resursu un tehnoloģiju infrastruktūra

- Izveidot integrētu, sadarbību veicinošu pētniecības infrastruktūru, kas ļautu viegli piekļūt un izmantot valodas resursus un tehnoloģijas
- Novērst pašreizējo sadrumstalotību un piedāvāt stabilu, pastāvīgu un paplašināmu infrastruktūru

