

Europos kalbų išteklių konsorciumo seminaras

Diskusija Kalbos duomenys ir skaitmeniniai ištekliai Lietuvoje

Andrius UTKA
Jolanta ZABARSKAITĖ
Audrius VALOTKA
Audra IVANAUSKIENĖ

Europos kalbų išteklių konsorciui aktualūs lietuvių kalbos duomenys ir skaitmeniniai ištekliai

Andrius UTKA



1. Kokie lietuvių kalbos ištekliai yra aktualūs Europos kalbų konsorciui?
2. Kokie lietuvių kalbos ištekliai nėra aktualūs?
3. Kaip atsiranda nauji kalbos ištekliai?
4. Kokios Lietuvos institucijos šiais kalbos ištekliais disponuoja?
5. Kaip juos galima gauti?



1. Lygiagretūs duomenys, (lietuvių kalba su viena iš 24 oficialių Europos Sąjungos kalbų):
 - Lygiagretūs tekstynai,
 - Vertimo atmintys,
 - Žodynai ir leksikonai.
2. Vienkalbiai lietuvių kalbos duomenys:
 - Bendrieji tekstynai,
 - Specialieji tekstynai (tam tikros temos).



- Lygiagretūs duomenys su **ne oficialiomis Europos Sąjungos kalbomis**.
- **Nedabartinės** lietuvių kalbos duomenys.
- Įvairūs **vienkalbiai** lietuvių kalbos **žodynai**.
- Mažiau reikalingi ir **grožinės literatūros** tekstynai.



Lygiagretūs ištekliai

- Vertimo atmintys kaupiasi verčiant su vertimo atminčių sistemomis (pvz., Trados),
- Lygiagrečius tekstynus kaupiant reikia ne tik gauti tekstus, bet ir juos sulygiuoti naudojant specialius lygiagretinimo įrankius,
- Kuriant žodynus paprastai atliekamas nelengvas leksikografinis darbas.

Vienkalbiai ištekliai

- Juos kaupti lengviau, nei lygiagrečius:
turint tinkamas kopijavimo teises tekstus galima automatiškai siurbti iš interneto arba gauti iš leidyklų.

Kokios Lietuvos institucijos šiais kalbos ištekliais disponuoja?



- Vilniaus universitetas (sukūręs pirmąją statistinio mašininio vertimo sistemą),
 - lietuvių/anglų ir lietuvių/prancūzų lygiagretūs tekstynai,
 - dvikalbiai žodynai.
- Vytauto Didžiojo universitetas (sukūręs pirmąją anglų-lietuvių internetinę taisyklinę mašininio vertimo sistemą)
 - anglų-lietuvių, latvių-lietuvių, vokiečių-lietuvių lygiagretieji tekstynai
 - dabartinės lietuvių kalbos tekstynas (200 mln. žodžių) ir Bendrasis internetinis tekstynas (1 mlrd. žodžių).
- Lietuvių kalbos institutas
 - dvikalbiai žodynai (lietuvių/anglų, latvių, lenkų, vengrų, vokiečių)
- Valstybinė lietuvių kalbos komisija
 - terminų bankas (terminai, kurie turi vertimus)
- Vertimo biurai
 - įvairios vertimo atmintys



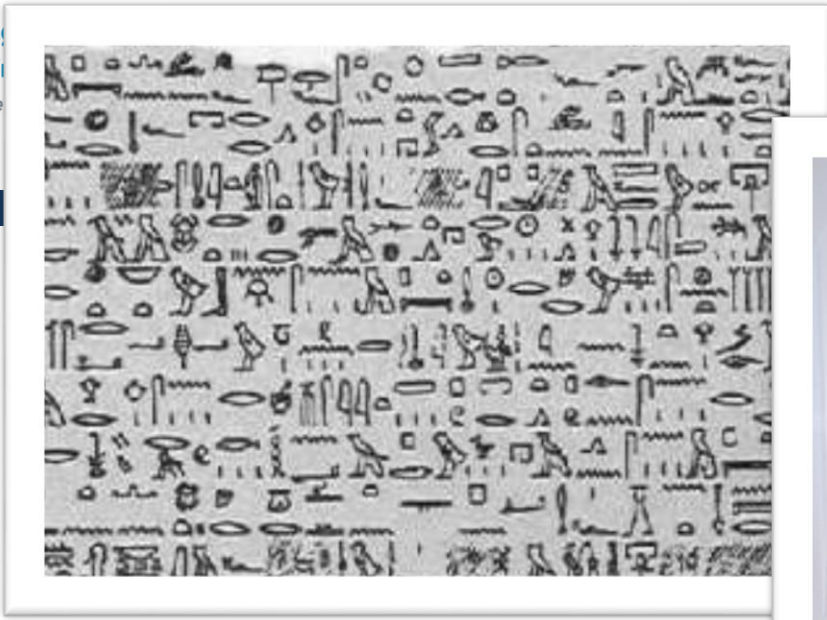
- Tiesiogiai tariantis su institucijomis:
 - nes daugelis išteklių nors ir yra prieinami internete, tačiau duomenų parsisiuntimas yra ribojamas licenciniais susitarimais.
- Per mokslinių tyrimų infrastruktūras (MTI):
 - META-SHARE (metashare.tilde.com),
 - CLARIN (šiuo metu kūrimo fazėje) (www.clarin-lt.lt)
(kuriose licencinės sutartys sprendžiamos sistemiškai)
- Tariantis su vertimo biurais:
 - turint omeny, kad daugelis vertimo tekstų yra privačių klientų nuosavybė su konfidencialumo apribojimais.

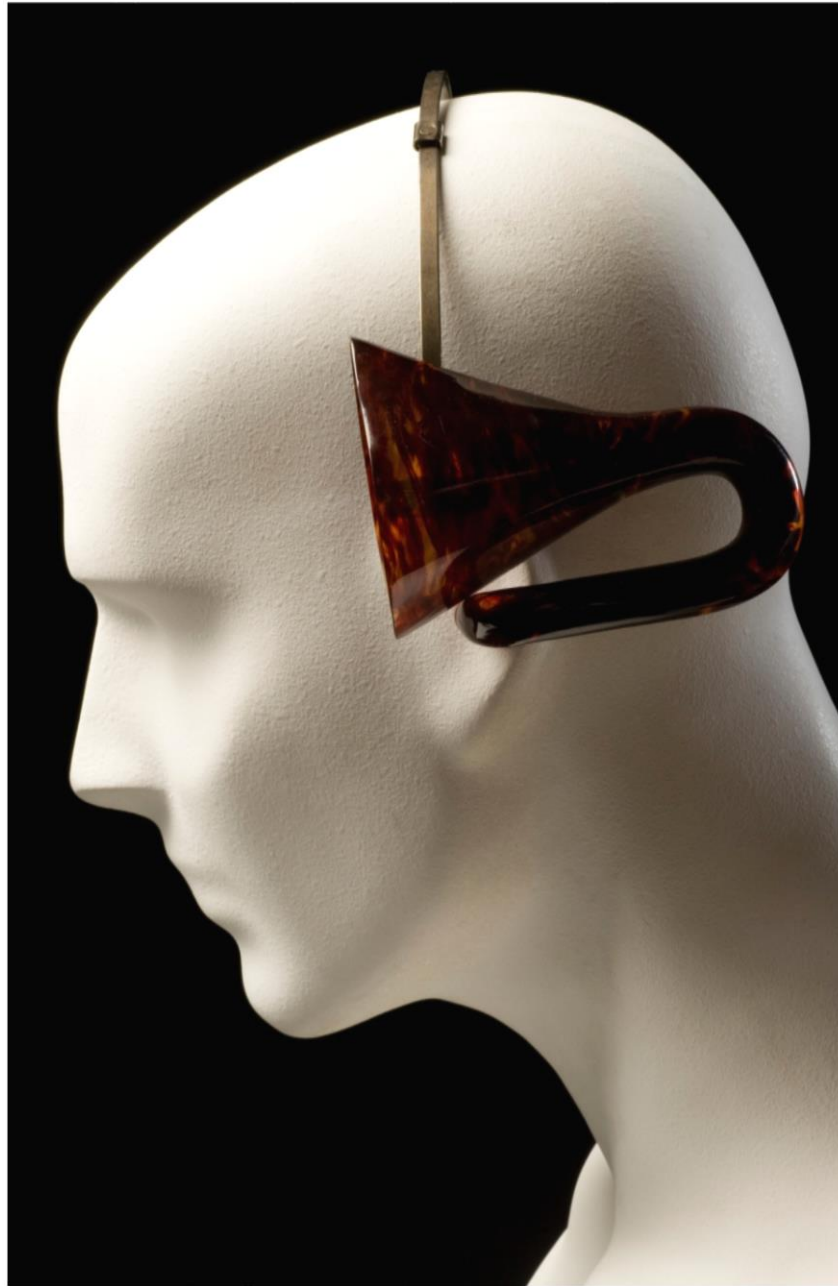


Audrius VALOTKA

◦ VILNIAUS UNIVERSITETO FILOLOGIJOS FAKULTETAS

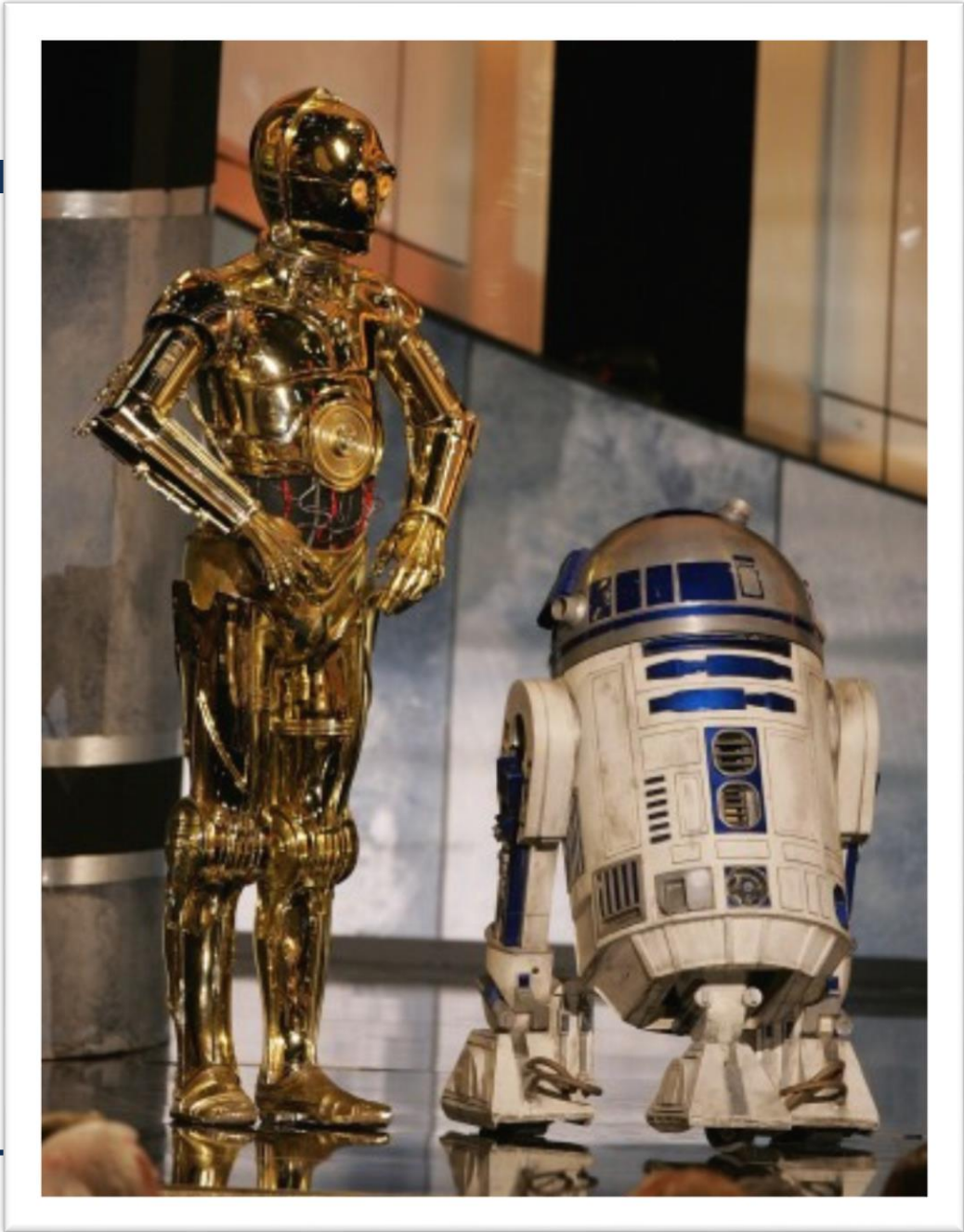






Kalba A (tekstas) - Kalba B (tekstas)

Šneka A – Kalba A – Kalba B –
Šneka B









• **AČIŪ UŽ DĖMESĮ**



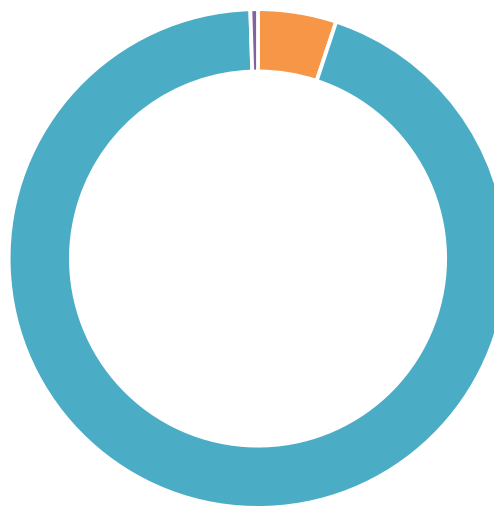
Europos kalbų išteklių konsorciumo seminaras

Diskusija Kalbos duomenys ir skaitmeniniai ištekliai Lietuvoje

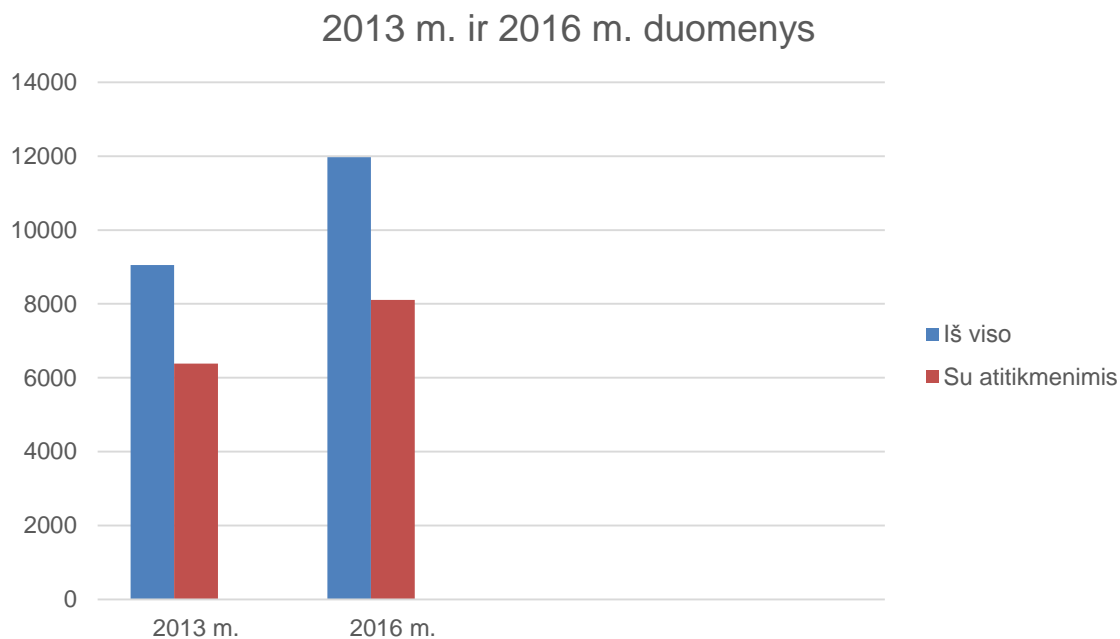
Audra IVANAUSKIENĖ



Aprobuota – 11971 termino straipsnis
Teiktina – 222834 terminų straipsniai
Neteiktina – 1160 terminų straipsnių



■ Aprobuoti ■ Teiktini ■ Neteiktini



Europos kalbų išteklių konsorciumo seminaras

Diskusija Kalbos duomenys ir skaitmeniniai ištekliai Lietuvoje

Jolanta ZABARSKAITĖ



- LKI valdomų ir kuriamų išteklių tipai
- LKI išteklių ypatybės
- Vartotojai
- Išteklių kūrimo ir tvarkymo strateginės kryptys, pridėtinė vertė ir sklaida
- Ateities perspektyvos



- LKI mokslinės ir taikomosios veiklos, kaupiant, skaitmeninant ir kuriant skaitmeninius lietuvių kalbos išteklius, metodologinis pagrindas yra požiūris, kad:
- kai postmodernaus pasaulio raidos būtina sąlyga yra žinios ir kūrybiškumas (kurių raiškos yra kalbinės), tinkamai suvokta, organizuota ir naudojama kalba tampa ne tik praktine informacijos keitimosi priemone daugiakalbėje aplinkoje, bet ir kultūrinės, ekonominės bei socialinės vertės kūrimo priemone ir medžiaga.



LKI kaupiamų skaitmeninių išteklių tipai:

- skaitmenizuoti vienakalbiai žodynai (DŽ, LKŽ, BŽ, sinonimų, antonimų, palyginimų, frazeologijos žodynai ir t. t.),
- skaitmenizuoti dvikalbiai žodynai (lietuvių-vengrų, baltarusių-lietuvių),
- specializuotos duomenų bazės (pavardžių, vietovardžių, naujažodžių, naujųjų skolinių ir kt.),
- specializuoti tarminiai ištekliai (geolingvistinių duomenų daugiafunkcinė integruota bazė, tarmių archyvo informacinė duomenų bazė ir kt.),
- specializuoti tekstynai (sąkytinės žiniasklaidos tekstynas 1960–2010, senųjų raštų anototas tekstynas SLIEKKAS, senųjų raštų duomenų bazė ir kt.),
- geoinformacinė Lietuvos vietovardžių duomenų bazė,
- naujažodžių tartuvas it t.t.



- LKI skaitmeninių išteklių ypatybės:
- struktūrizuoti,
- su labai detalia atributine paieška,
- integruojantys įvairaus formato medijas (įvairialypis turinys),
- nuolat atnaujinami atsižvelgiant į naujausius kalbos pokyčius,
- integruoti.



- <http://www.prusija.lki.lt/>
- <http://lkiis.lki.lt/>
- <http://titus.uni-frankfurt.de/sliekkas/>



Plėtros strategija:

- tolimesnis visų LKI išteklių integravimas į bendrą infrastruktūrą,
- prasminių sąsajų sistemos sukūrimas,
- paieškų galimybių didinimas,
- paieškų vizualizavimas,
- paslaugų vartotojui plėtra ir įvairovė,
- integracija į viešojo administravimo paslaugas,
- integracija į daugiakalbes infrastruktūras.