

Tarptautinis mokslinis seminaras  
*A. Bezenbergeris - etnografas ir akademinės baltistikos pradininkas*  
(Juodkrantė, Neringa 2011 09 9–10)

## Senosios lietuvių kalbos referencinis korpusas (SLIEKKAS). Lingvistinė anotacija

Mieli kolegos ir bičiuliai! Esu Jums labai dėkinga už pakvietimą į Adalbertui Bezenbergeriui skirtą seminarą.

Pirmasis publikavęs lietuvių kalbos senuosius tekstus ir pirmasis pradėjęs analizuoti senąją lietuvių kalbą, Bezenbergeris iš esmės padėjo pagrindus visiems tolesniems XVI–XVII amžiaus lietuvių kalbos tyrimams. Deja, dar negalime teigti, kad tradiciškai senąją lietuvių kalbą vadinamą laikotarpį jau esame ištyrę išsamiai, nuodugnai ir sistemingai. Norisi tikėtis, kad iki 2077 metų, kai sukaks 200 metų nuo Bezenbergerio *Beiträge zur Geschichte der litauischen Sprache (Auf Grund litauischer Texte des XVI. und XVII. Jahrhunderts, Göttingen: Robert Peppmüller, 1877)*, šių dviejų šimtmečių lietuvių kalba bus detaliai ir sistemingai aprašyta, turėsime lietuvių kalbos istorinį žodyną ir senosios lietuvių kalbos gramatiką.

Siekdami sistemingo, o ne paremto atsitiktiniais pavieniais duomenimis tyrimo, pirmiausia turime nuspręsti, kaip parengti tokią sistemą, kuri leistų mums rasti patikimus atsakymus į įvairius lingvistinius klausimus. Pirmiausia, žinoma, reikia turėti visus tiriamo laikotarpio tekstus. Čia neaplenkiama yra Lietuvių kalbos institute rengiama bei nuolat pildoma *Senųjų raštų duomenų bazė*, apimanti tekstus nuo 1573 iki 1816 metų. Ji yra didžiausias ir tekstų dokumentinio tekstų perteikimo atžvilgiu patikimiausias senosios lietuvių kalbos tekstynas. Kitas etapas yra suskaitmenintus tekstus parengti tolesnei analizei. Šiuo atžvilgiu pavyzdiniai yra Frankfurto prie Maino Goethe's universiteto TITUS (*Thesaurus Indogermanischer Text- und Sprachmaterialien*) duomenų bazėje esantys lietuviški tekstai. Šioje bazėje įvairių tekstų žodžių formos yra pagal tam tikrą indeksacijos sistemą susietos vienos su kitomis ir iš dalies su vertimo šaltiniais.

Modernūs kalbos istorijos tyrimai reikalauja ne tik priėjimo prie patikimų empirinių duomenų, bet ir kokybiškos tų duomenų paieškos galimybės pagal įvairius lingvistinius kriterijus. Taigi suskaitmenintus ir filologiškai užtikrintus tekstus, kaip empirinių rašytinių senosios lietuvių kalbos faktų rinkinį, reikia paruošti tolesniems moksliniams tyrimams, visų pirma lingvistiniams, bet taip pat ir istoriniams plačiaja prasme.

**Senosios LIETUVIŲ Kalbos KorpusAS** ⇒ **SLIEKKAS** ir siekia parengti tokio tyrimo mokslinius bei techninius pagrindus. Referencinio senosios lietuvių kalbos korpuso tikslas yra reprezentuoti šio laikotarpio kalbą kaip visumą bei pateikti išsamią informaciją apie kalbą, jos variantus, žodyną ir gramatiką. SLIEKKAS turi padėti ir

pradėti įgyvendinti du didžiausius diachroninės lituanistikos darbus: senosios lietuvių kalbos gramatikos ir istorinio lietuvių kalbos žodyno rengimą.

SLIEKKAS susideda iš trijų pagrindinių dalių: iš metainformacijos apie tekstus, iš pačių tekstų ir iš jų lingvistinės anotacijos. Visas tris dalis kartu ar kiekvieną atskirai turi būti galima tirti pagal pasirenkamus paieškos kriterijus:

metaduomenys	paieškos įrankiai
rašytiniai tekstai (empiriniai kalbos faktai)	
lingvistinė anotacija	

SLIEKKAS yra daugiapakopės (vad. *Stand-off*) architektūros struktūra. Ją sudaro įvairūs metalingvistinės, teksto ir lingvistinės informacijos anotaciniai sluoksniai. Pats tekstas ir kiekvienas jo anotacinis sluoksnis faktiškai yra atskiras ir integralus dokumentas. Tokią architektūrą įgalina sukurti XML (*Extensible Markup Language*) duomenų struktūrų kalbos formatas. Kiekvieną sluoksnį galima plėsti ir pildyti nepažeidžiant kitų sluoksnių, tačiau visi sluoksniai yra tarpusavyje sinchronizuoti ir susieti. Toks formatas leidžia pagal įvairius kriterijus koduoti pirminius teksto duomenis ir pateikti visas reikalingas tekstologines bei lingvistines teksto anotacijas. Paieškos kriterijus galima modeliuoti pagal kiekvieną anotacijos sluoksnį ir pagal kelių sluoksnių kombinacijas. Tekstai su daugiapakopėmis anotacijomis perauga „tekstynų“, t. y. tekstų rinkiniu, vadinamas ribas. Todėl ne iš išpaikimo, bet sąmoningai vartojamas terminas „korpusas“.

Čia nekalbėsiu apie korpuso techninį konceptą, apie anotacines programas ir apie vizualizacijos bei paieškos sistemas. Šita korpuso pusė yra ne tiek lingvistikos, kiek lingvistinės informatikos dalykas. Taip pat nedetalizuosiu tekstologinės anotacijos, kuri tiesiogiai susijusi su rankraščio ar spaudinio kritiniu parengimu ir atspindi arba atstoja kritinį tekstologinį aparatą.

Šiuo atveju mums svarbiausias yra vieningais principais paremtos lingvistinės anotacijos sistemos kūrimas. Kaip ir rašant gramatiką ar sudarant žodyną, lingvistinė anotacija yra pagrįsta hierarchijos principu. Lingvistinė informacija pateikiama sluoksniais, pradedant nuo stambiausių kategorijų ir detalizuojant iki smulkesnių fleksinės morfologijos požymių.

LEMATIZAVIMAS, GLOSAVIMAS. Kiekviena korpuso teksto „žodžio forma“ (šis terminas, vad. *Token*, apima ne tik atskirų žodžių formas, bet ir skaitmenis, skyrybos ženklus, grafinius simbolius bei santrumpas) yra lematizuojama. Pirmiausia pateikiama standartizuota žodžio forma, tokia, kokia ji būtų atstatyta istoriniame žodyne, ir antraštinis žodis, lema, t. y. pagrindinė forma istoriniame žodyne. Toliau antraštinis žodis (lema) glosuojamas, t. y. užrašomas dabartine bendrine kalba ir pateikiamos jo reikšmės. Taip pat glosuojama konkreti pavartota žodžio forma, pateikiama jos reikšmė konkrečiame aprašomame kontekste. Kaip tai atrodo, matyti pirmoje lentelėje:

1 lentelė. Standartizuota forma, lematizavimas, glosavimas, vertimas

<i>Anotacijos sluoksniai</i>	<i>Reikšmė</i>	<i>Pavyzdys</i>	
<b>Rankraštis, spaudinys</b>	forma rankraštyje ar spaudinyje	fchwaifdes	Gregorius
<i>... (tekstologinės anotacijos sluoksniai)</i>			
<b>Standartizuota forma</b>	standartizuota forma (atstatyta kaitoma forma istoriniame žodyne)	žvaizdės	Gregorius
<b>Lema</b>	antraštinis žodis (pagrindinė forma istoriniame žodyne)	žvaizdė	Grigalius
<b>Lemos glosavimas</b>	antraštinio žodžio glosavimas: užrašymas bendrine kalba ir reikšmės, vertimas (skolinių, tarmybių ir nebevertojamų, dabartiniam vartotojui nebeatpažįstamų žodžių, kitomis kalbomis parašytų žodžių)	žvaigždė	Grigalius
<b>Konkrečios formos glosavimas</b>	konkrečios žodžio formos glosavimas: užrašymas bendrine kalba ir reikšmė (vertimas) konkrečiame kontekste	Aušrinė žvaigždė, Veneros planeta	Grigalius Didysis

LINGVISTINĖ ANOTACIJA: *POST (Part of Speech-Tagging)*. Tolesni korpuso sluoksniai skirti lingvistinei anotacijai. SLIEKKO lingvistinė anotacija iš esmės yra morfologinė. Morfologinė anotacija apima kalbos dalių (gramatinių klasių) žymėjimą (vad. *Part of Speech-Tagging, POS*) ir morfologinius žodžių formų aprašus. Anotuojama hierarchijos principu, pirmiausia lema, tada konkreti pavartota žodžio forma. Visas SLIEKKO pažymų kompleksas (vad. *Tagset*) suskirstytas taip pat hierarchiškai, pradedant bendra informacija apie kalbos dalį ir toliau skirstant į smulkesnius morfologinius poklasius.

Lemos anotuojamos pagal 10 gramatinių klasių.

2 lentelė. Kalbos dalys (gramatinės klasės) ir jų pažymos

1. <b>ADJ</b> – būdvardis (adiectivum)	7. <b>N</b> – daiktavardis (nomen): <b>NA</b> bendrinis (apeliatyvinis), <b>NT</b> tikrinis
2. <b>ADV</b> –rieveiksmis (adverbium)	8. <b>P</b> – įvardis (pronomen): <b>PD</b> demonstratyvinis, <b>PI</b> indefinitinis, <b>PK</b> klausiamasis (interogatyvinis), <b>PPER</b> asmeninis
3. <b>AP</b> – prielinksnis, polinksnis (adpositio)	9. <b>PTK</b> – dalelytė (particula)
4. <b>CARD</b> – kiekinis (cardinale) skaitvardis, <b>ADJO</b> kelintinis (ordinale) skaitvardis	10. <b>V</b> – veiksmazodis (verbum): <b>VA</b> pagalbinis, <b>VV</b> pagrindinis
5. <b>ITJ</b> – jaustukas (interjectio), ištiktukas (onomatopoeticum)	
6. <b>KO</b> – jungtukas (conjunctio)	

Kiekvienos žodžio formos gramatinė klasė anotuojama dvejopai: pagal lemą (antraštinio žodžio gramatinė klasė) ir pagal konkrečią pavartotą žodžio formą (konkrečios žodžio formos pavartojimas). Tokiu būdu lemos pagrindinė gramatinė klasė atskiriama nuo konkrečios formos vartojimo. Šis skyrimas padeda:

- žymėti gramatinių klasių kitimus, pvz., būdvardžių daiktavardėjimą, dalyvių būdvardėjimą bei daiktavardėjimą; daiktavardžių irrieveiksmių virtimą adpozicija (prielinksnium ar polinksnium), daiktavardžių, būdvardžių, dalyviųrieveiksmėjimą ir pan.;

- apibūdinti ambivalentiškus atvejus, kai neaišku, kuriai kalbos daliai priskirti konkrečią žodžio formą, pvz., sustabarėjusių linksnių prieveiksmėjimo atveju konkreti pavartota forma yra prieveiksmis *atėjo vidunakčiu* (ADV), o pati lema *vidunaktis* – daiktavardis (NA);
- adpozicijos (prielinksnio) klasėje skirti prielinksnius (prepozicijas) ir polinksnius (postpozicijas), pvz.: *link* (APPR) *namų* ir *namų link* (APPO);
- pateikti aiškia neasmenuojamų veiksmažodžio formų definiciją, pvz., *kalbėdamas tarė* konkreti pavartota forma yra predikatyvinis pusdalyvis/semiparticipijas (VVSP), o lema *kalbėti* – veiksmažodis (VV). Toks lemos ir konkrečios formos skyrimas svarbus aprašant bendratį, dalyvius, gerundyvus (reikiamybės dalyvius), pusdalyvius (semiparticipijus), padalyvius (gerundijus) ir supiną.

Kiekviena iš lemos gramatinių klasių turi įvairaus detalumo subklasifikaciją, kuri priklauso nuo konkrečios pavartotos žodžio formos morfologijos.

Tais atvejais, kada morfologinė išraiška kinta priklausomai nuo sintaksinės funkcijos, anotacija apima ir sintaksinės morfologijos lygmenį. Morfosintaksinis anotavimas leidžia išryškinti skirtybes, kurios neatsispindi grynai morfologiniame apraše. Tai svarbu šiais atvejais:

- o įvairuojanti būtinojo (vardažodinis tarinys) ir laisvojo (tarininis pažyminys) predikatyvo morfologinė forma: vardininkas arba įnagininkas, pvz., *Jonas buvo pranašas/pranašu*. Konkrečių formų morfologinėse pažymose skiriamas predikatyvinis daiktavardis, būdvardis, skaitvardis, įvardis, dalyvis (+gerundyvas), pusdalyvis;
- o pagalbinio veiksmažodžio ir pagrindinio veiksmažodžio išskyrimas reikalingas aprašant tokius veiksmažodinius kompleksus kaip sudėtinės laikų formas. Todėl tikslinga lemos morfologinėse pažymose veiksmažodį skirti į pagrindinį (VV) ir pagalbinį (VA). Konkrečių formų morfologinėse pažymose svarbu skirti atributinę ir predikatyvinę bei adverbialinę neasmenuojamų veiksmažodžio formų vartoseną.

Plg. 3 lentelę:

3 lentelė. Lemos ir konkrečios žodžio formos POS

<i>Anotacijos sluoksniai</i>	<i>Reikšmė</i>	<i>Pavyzdys</i>
<b>M1a lema</b>	antraštinio žodžio (lemos) kalbos dalis	[Jonas tapo] <i>pranašu</i> NA (bendrasis daiktavardis)
<b>M1b konkreti žodžio forma</b>	konkrečios žodžio formos morfologinė (šiuo atveju morfosintaksinė) charakteristika	NAP (bendrasis daiktavardis, predikatyvinis)

MORFOLOGINĖ ANOTACIJA. Tolesnė morfologinė anotacija taip pat yra hierarchinė. Ji susideda iš trijų sluoksnių. Pirmas ir antras sluoksnis žymi tokias nekintamas, inherentines kalbos dalių morfologines kategorijas kaip kaitybos klasė (kaitybos kamienas) ir giminė (giminėmis nekaitomų klasių: daiktavardžių ir daiktavardiškų skaitvardžių). Kaip ir anotuojant kalbos dalis (POS), apibūdinamos pirmiausia lemos (antraštinio žodžio), o tada konkrečios žodžio formos kategorijos. Toks lemos ir žodžio formos skyrimas leidžia fiksuoti kaitybos klasės ir giminės įvairavimą, pvz.:

- lema – *amžius* (ju\_Masc) ir tekste vartojamos įvairių kamienų formos *amžias* (ja\_Masc), *amžas* (a\_Masc), *amžis* (i\_Masc);
- lema – *gražus* pagrindinė būdvardžio forma nelyginamuoju laipsniu (*u,jo* kamienas, vyriškoji ir moteriškoji giminė), o konkreti žodžio forma *gražiausias* aukščiausiuoju laipsniu (*ja* kamienas);
- lema – *neklaužada* yra bendrosios giminės (Com), o konkreti žodžio forma junginyje *baisus* *neklaužada* yra vyriškosios giminės (Masc).

Trečias sluoksnis žymi kintamus fleksinės morfologijos požymius: giminę (giminėmis kaitomų klasių), skaičių, linksnį, laipsnį, asmenį, nuosaką, laiką, diatezę (rūšį), emfazę (apibrėžtumą, kitaip, įvardžiuotines formas). Plg. 4 lentelę:

4 lentelė. Lemos ir konkrečios žodžio formos morfologinė anotacija

<i>Anotacijos sluoksniai</i>	<i>Reikšmė</i>	<i>Pavyzdys</i>		
<b>M2a lemos fleksija</b>	antraštinio žodžio (lemos) nekintamos morfologinės kategorijos (kaitybos klasė ir giminėmis nekaitomų klasių giminė)	<i>gyvename</i> ti-Inf	<i>gražiausiame</i> u,jo	<i>amže</i> ju_Masc
<b>M2b konkrečios žodžio formos fleksija</b>	konkrečios žodžio formos nekintamos morfologinės kategorijos (kaitybos klasė ir giminėmis nekaitomų klasių giminė)	a	ja	a_Masc
<b>M2c konkrečios žodžio formos fleksinė morfologija</b>	konkrečios žodžio formos fleksinės morfologijos požymiai	Act_Ind_Pres_Pl_1	Sup_Masc_Sg_Loc	Sg_Loc

Daugiasluoksnią hierarchinę korpuso anotacijos struktūrą ne tik įgalina detaliai aprašyti kiekvieno teksto morfologinę sandarą, bet ir gali būti diachroninio žodyno bei senosios lietuvių kalbos gramatikos pagrindas. Toliau galima ir reikia numatyti sintaksinės anotacijos sluoksnius, eksplisitinį ir implicitinį citatų žymėjimą, susiejimą su tiksliai nustatytais ir lingvistiškai anotuotais vertimo šaltiniais. Bet tai jau atskiro pranešimo reikalinga tema.

SLIEKKO iniciatyvoje ir parengiamuosiuose darbuose šiuo metu bendradarbiauja Frankfurto prie Maino Goethe's universitetas (Vokietija), Lietuvių kalbos institutas (Vilnius) ir Pisos universitetas (Italija). 2010 metais SLIEKKA rėmė Lietuvos Respublikos švietimo ir mokslo ministerija pagal programą „Mokslo ir studijų modernizavimas“. Viliamės, kad tolesni korpuso darbai bus remiami bendromis tarptautinėmis jėgomis.

Ačiū Jums už kantrybę ir dėmesį! Būsiu labai dėkinga už Jūsų pastabas, komentarus ir pasvarstymus, kuriuos labai prašyčiau man atsiųsti elektroniniu paštu.