

Samisk språkteknologi ved Universitetet i Tromsø

Disambiguering og syntaktisk analyse(1) ◊ Samisk stavekontroll(2)

(1): Lene Antonsen, Saara Huhamanniemi, Marit Julien, Ilona Kivinen, Trond Trosterud, Linda Wiechetek <http://giellatekno.uit.no>

2): Børre Gaup, Sjur Moshagen, Thomas Omma, Maaren Palismaa, Tomi Pieski (Sametinget i Noreg) <http://divvun.no>



SÁMEDIGGI SAMETINGET



Resultat – samisk

Grammatisk analyse

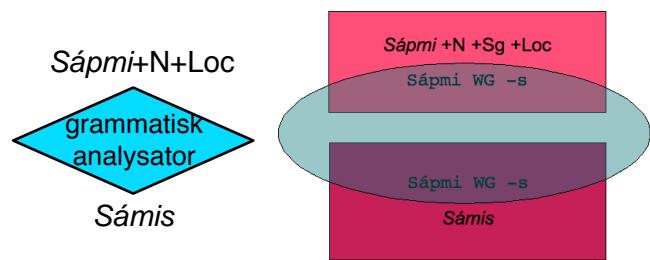
Vi har laga ein grammatiskeksanalyse for nord- og lulesamisk tekstu (ordboying og -avleiring, lydovergangar), og ein setningsanalyse for nord- og lulesamisk (den nordsamiske har presisjon = 97,5 % og recall = 94,4 %).

Vi har laga eit syntaktisk annotert korpus for nordsamisk, som vi har gjort tilgjengeleg for språkforskarar over internett.

Stavekontroll

Den grammatiskeksanalyzator er grunnstammen i Sametinget sin stavekontroll for nord- og lulesamisk, som blir publisert i slutten av 2007 (betaversjonen ferdig i februar 2007).

Programmet inneholder også eit orddeleprogram.



Tekst frå den samiske avisa *Min Áigi*, i det samiske retteprogrammet
(korrekte former: čeabehii, gávnaháhit, boadásii)

Ii juohke bohccø xeabehii
GPS-čeabetbáttid dahje dan teknika maid gavnaháhit buoremussan ii eisege leat ulbmil bidjet juohke bohccui.
– Bálgoosiid siste leat earláagan vuogádagat ja várra dat boadašii leat siiddaid anus.

A n d r e s p r á k

Gronlandsk stavekontroll

Utan morfologiske automater er det umogleg å analysere grønlandsk, eller lage stavekontrolla: dei fleste orda på grønlandsk tilsvarer ei heil setning på norsk. I oktober 2006 vart det første korrekturprogrammet for grønlandsk lansert, basert på infrastrukturen til dei samiske analyzatorane.

Ordbøker for komi og kvensk

Tida då ordboks- og språkteknoiprojektet var uavhengige av kvarandre er over. Med å koble ordbok og grammattikk saman får vi eit kraftig analyseverktøy.

Med utgangspunkt i arbeidet vårt med samisk terminologi og analyse har vi gjort ei komi-ordbok om til ein integrert ordbok og analyzator, og samtidig til ein modell for arbeid med maskinomsetjing i framtida. Vi har også laga eit oppsett for arbeid med kvensk leksikografi.

Færøysk analyzator

Vi har utarbeidd ein betaversjon for ein færøysk analyzator, som kan danne utgangspunkt for den færøyske delen av eit dansk-færøysk maskinomsettingsprogram.

"<Ataasinngorpat>"
"ataasinngor" V Int 3Pl @PRED
"ataasinngor" V Ind 2Sg 3Sg @PRED
"ataaseq" NNGUR V Con 3Sg @PRED
"ataaseq" NNGUR V Int 3Pl @PRED
"ataasinngor" V Con 3Sg @PRED
"ataaseq" NNGUR V Ind 2Sg 3Sg @PRED
"<16.10.2006>"
"16.10.2006" Num Abs @SUBJ
"
". ." CLB <<<
"<Kukkuniaat>"
"kukkuniaat" N Abs Sg @SUBJ
"Kukkuniaat" N Abs Sg @SUBJ
"kukkuniaat" N Abs Sg @OBJ
"kukkuniaat" N Abs Sg @OBJ
"<siulleq>"
"siulleq" N Abs Sg @SUBJ
"siulleq" N Abs Sg @OBJ
"<aaneqarsinnaangorluni>"
"aaneqarsinnaangorluni" ?
"<piariissaag>"
"piareer" SSA V Ind 3Sg @PRED
"pia" RISSA V Ind 3Sg @PRED
"
". ." CLB <<<
"<Aallerneq>"
"aaller" NIQ N Abs Sg @SUBJ
"aallerneq" N Abs Sg @SUBJ
"aallerneq" N Abs Sg @OBJ
"aaller" NIQ N Abs Sg @OBJ
"<akeqanngilag>"
"akik" QAR NNGIT V Ind 3Sg @PRED
"aki" QAR NNGIT V Ind 3Sg @PRED
"
". ." CLB <<<
"<Uani>"
"uv" Adv Lok @ADVL
"<aallersoqarsinnaavoq>"
"aallersoq" QAR SINNAA V Ind 3Sg @PRED
"aaller" TUQ QAR SINNAA V Ind 3Sg @PRED
"
". ." CLB <<<

A n a l y s e m e t o d a r

Vi har brukt endelige tilstandautomater for morfologisk analyse. Til morfofonologien bruker vi tonivåmorfolologi. Analyzatorane våre er kompliert med komplitorar frå Xerox, (lexc, twlc og xfst, jf. <http://www.fsmbook.com>).

Til syntaktisk analyse har vi brukt ein grammatiskeks basert syntaktisk analyzator (føringsgrammatikk), vislcg frå Syddansk universitet, jf. <http://visl.sdu.dk>.

Automatisk analyse av 1. Mosebok, v. 1-2:

"<Álggus>"
"álgu" N Sg Loc @ADV
"<Ipmil>"
"ipmil" N Sg Nom @SUBJ
"<sivdnidii>"
"sivdnidit" V TV Ind Prt Sg3 @+FMAINV
"<almmi>"
"albmi" N Sg Acc @OBJ
"<ja>"
"ja" CC @CC-NP
"<eatnama>"
"eana" N Sg Acc @OBJ
"
". ." CLB <<<
"<Eanan>"
"eanan" N Sg Nom @SUBJ
"<lei>"
"leat" V IV Ind Prt Sg3 @+FMAINV
"<ávdin>"
"ávdin" A Sg Nom @SPRED
"<ja>"
"ja" CC @CC-NP
"<guorus>"
"guorus" A Sg Nom @SPRED
"
". ." CLB <<<
"<Seavdnijatvuoha>"
"seavdnijatvuoha" N Sg Nom @SUBJ
"<govčai>"
"gokčat" V TV Ind Prt Sg3 @+FMAINV
"<čiekjala>"
"čiekjala" N Sg Acc @OBJ
"
". ." CLB <<<
"<Ipmila>"
"ipmil" N Sg Gen @GN
"<Vuogjá>"
"Vuogjá" N Prop Plc Sg Nom @SUBJ
"<sattáhalai>"
"sattáhallat" V IV Ind Prt Sg3 @+FMAINV
"<čáziid>"
"čáhcí" N Pl Gen @GP
"<bajábealde>"
"bajábealde" Po @ADV
"
". ." CLB <<<
"<De>"
"de" Adv @ADVL
"<Ipmil>"
"ipmil" N Sg Nom @SUBJ
"<celkkii>"
"cealkit" V TV Ind Prt Sg3 @+FMAINV
"
". :" CLB
"
". ." PUNCT RIGHT
"<Šaddos>"
"šaddat" V IV Imprt Prs Sg3 @+FMAINV
"<čuovga>"
"čuovga" N Sg Nom @SUBJ
"
". !>"
"excl" CLB <<<
"
". ." PUNCT RIGHT
"<Ja>"
"ja" CC @CC-VP
"<čuovga>"
"čuovga" N Sg Nom @SUBJ
"<šattai>"
"šaddat" V IV Ind Prt Sg3 @+FMAINV
"
". ." CLB <<<

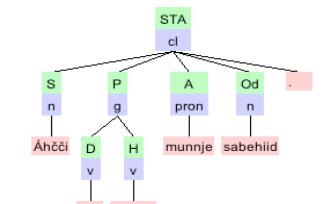
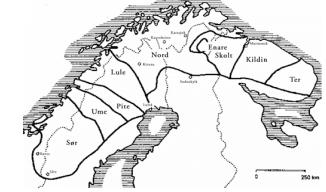
F r a m t i d s p e r s p e k t i v

Utvide til andre samiske språk

Vi kan bruke ein grammatiskeksanalyzator for eitt språk som utgangspunkt for tilsvarende analyzatorar for nært beslektte språk. Det er seks samiske skriftspråk, vi har ein analyzator for to av dei, er i ferd med å lage den tredje, og ser på dei tre siste. Delar av regeloppsettet kan brukast på nytt.

Pedagogisk programvare

Samiskundervisninga blir ofte organisert som desentralisert nettundervisning. Samiske verb blir boygd i eintal, total og fleirtal i tre personar, ordet for «ikkje» er eit verb. Å lære dette mønsteret er ein sentral del av nybyjaroplæringa, automaten vår gir det mogleg å lage elevstyrte, interaktive spel som for det til. Vi kan også la analyzerte setningar gå inn i grammatikkspill, som her:



Informasjonssøk på samisk (Kva er vitsen med å lage noko du veit du ikke finn att?)

For å søke på ordformene språk, språket, språkene, språka kan vi skrive språk*. For det tilsvarende samiske giella, er ca. 60 % av belegga i løpende tekstu i sakkalla svakt stadium (giela), mens ca. 10 % står i illativ (gillii). Berre 30 % står i sterkt stadium, dvs. at ved å søke etter giella* går vi glipp av 70 % av dei relevante formene. Automaten vår kjemmer att samlede boyingsformer av eit quart samisk ord, og kan danne stammen i avanserte informasjonssøkeprogram.

586 Geavatlacčat mearkkaša dát ahte sámeigella ii leat doaibmi **giella** diehtojuohkinteknologija oktavuodain .

680 Oslo universitetas galgá sikhkarit ain leat oahpahus suoma-ugralaš **gielain** boahtteáiggis .

681 Dán áássis leat sáhka sámeigiel fáladaga heaittheamis ii ge suoma-ugralaš **gielaid** heaittheamis oppalačat .

936 Dálá **lemma: giella** pos: N ja ku syn: @OBJ **giela** steathta sáttá leat sikhé áittan ja vejolašvuohant sámi **giela** syn: @OBJ **gielain** ja ovdideamis .

950 Kult **number: Pl** sihl **case: Acc** áilut lohkamuša girjjálašvuoda mas sámi mánát ja nuorat seammás várjavilcii ja ovdidivcii sámi **giela** .

Sámediggi oaidná ahte kulturráddi lea ruðaid prioriteren biergasiidda mat mánáigárddiin galget geavahuvot sámi **giela** ja kultuvra gaskustearpmái

S y n t a k t i s k e r e g l a r

"<Mii>"
"M" Num Ill
"mii" Pron Interr Sg Nom
"mii" Pron Rel Sg Nom
"mun" Pron Pers Pl1 Nom ← Fleitalspronomen fordi verbet står i fleitral
"
". ." CLB
"
"ii" V IV Neg Ind Pl1 ← verbet står i fleitral
"
"<leat>"
"leat" V IV Ind Prs Pl3
"leat" V IV Ind Prs Sg2
"leat" V Inf
"leat" V IV Ind Prs ConNeg ← leat står etter nektignsverbet, det er altså ikkje Inf eller personboygd form
"leat" V IV Ind Prs Pl1
"
"dat" Pron Pers Sg3 Acc
"dat" Pron Dem Sg Acc
"dat" Pron Pers Sg3 Gen
"dat" Pron Dem Sg Gen ← dan er ikkje demonstrativ, sidan det ikkje står foran adjektiv eller substantiv. Ikkje genitiv, sidan det ikkje står etter talord eller preposisjon, og ikkje foran substantiv, adjektiv eller postposisjon
"
"muitalan" "muitalit" V* TV Actio Der/eapni N Sg Gen
"muitalit" V TV Actio Acc
"muitalit" V TV Actio Gen
"muitalit" V TV Ind Prs Sgl
"muitalit" V TV PrfPrc
"muitalit" V TV Actio Nom ← muitalan er partisipp, sidan det står som komplement til nektingsverb
"
". ." CLB
"
"<Mii>"
"mun" Pron Pers Pl1 Nom @SUBJ
"
". ." CLB
"
"ii" V IV Neg Ind Pl1 @FAUXV
"
"<leat>"
"leat" V IV Ind Prs ConNeg @FAUXV
"
"<dan>"
"dat" Pron Pers Sg3 Acc @OBJ
"
"muitalit" "muitalit" V TV PrfPrc @FMAINV
"
". ." CLB
"
"<Mii>"
"mun" Pron Pers Pl1 Nom @SUBJ
"
". ." CLB
"
"ii" V IV Neg Ind Pl1 @FAUXV
"
"<leat>"
"leat" V IV Ind Prs ConNeg @FAUXV
"
"<dan>"
"dat" Pron Pers Sg3 Acc @OBJ
"
"muitalit" "muitalit" V TV PrfPrc @FMAINV
"
". ." CLB
"
". ." CLB