

Manne dihtor galgá máhttit grammatihka?

LENE ANTONSEN – TROND TROSTERUD

Ođđaáigásaš servodagas ii oktage giella birge čállingiellan giellateknologiija haga. Sámegillii leat ráhkaduvvon grammatihkalaš analysáhtorat, vaiko erenoamážit stuora giellaide lea dábálaš geavahit statistihkalaš analysáhtoriid. Artihkal gieđahallá man muttus sámi giellateknologiija lea dál, mo sámeigiella earrána dárogielas ja eŋgelasgielas, ja mo dat erohusat váikkuhit giellateknologalaš čovdosiidda.

Vuosttaš kapihttalis čilgejetne lokaliserema, dahje vuodđoinfrastruktuvrra mii dahká vejolažžan čállit sámegillii dihtoriin. Dasto guorahalle mo sátneanalysas golbma iešguđetlágan lahkonanvuogi heivejit sámegillii (kap. 2.). Goalmmát kapihttalis gehčce grammatihkalaš cealkkaanalysa, ja de čuovvovaš kapihttalis makkár geavtlaš prográmmaide sáhtta analysáhtoriid atnit odne juo. Maŋimuš fáddá (kap. 5) lea mo giellateknologiija váikkuha sámeigiela geavaheapmái min servodagas. Guovddážin lea ahte dakkár reaiddut šaddet dehálažžan jus galgá leat vejolaš geavahit giella omd. hálddahušlaš giellan.

1. Lokaliseren

Dihtormáilmmis tearbma *lokaliseren* máksá buot maid dihtoris ferte heivehit dihto riikii ja gillii. Sámeigiela lokaliseremii gullá áigediehtu, alfabehtalaš sorterren, boallobeavdi ja erenoamáš bustávat.

Buot sámeigielaid čállingielain gávdnojit bustávat mat eai gávdno riikkaid váldogielain, nugo boahtá ovdan tabeallas 1.

Tabella 2. Mo sánit Dieđut/dieđut ja Ođđasat/ođđasat leat čállojuvvon interneahτας 20.11.09.

Ohcantearbma	<i>dieđut</i>	<i>dieđut</i>	<i>dieđut</i>	<i>ođđasat</i>
Ohcandomena	samiskhs.no	uit.no	.no	olles interneahтта
Dárogieđa d	0,7 %	30,4 %	3,9 %	9,0 %
Islánddagiela ð	56,5 %	22,2 %	2,7 %	86,0 %
Sámegiela đ	42,8 %	47,4 %	93,4 %	5,0 %

Go rivttes bustávva lea geavahuvvon *dieđut*-sátnái *.no*-domenas 93,4 % dáhpáhusain, muhto *ođđasat*-sátnái dušše 5,0 %, de lea sivva ahte *dieđut* lea geavahuvvon sámegiela teavsttain, ja *Ođđasat* lea siterejuvvon TV-programma-namman (boasttu bustávain) dárogieđa teavsttain.

Nubbi lokaliserenváttisvuohta ii leat diehtemeahttunvuođas, muhto ieš dihtorvuogádagas. Sullii jagi 2000 maŋŋel buot dihtorvuogádagat dorjot sámegiela, muhto servodaga stuora registarat Norggas, nugo Álbmotregisttar ja Brønnøysund-registarat¹, geavahit velge boares dihtorstandárddaid, sámegiela bustávid haga. Álbmotregisttaris livččii dehálaš geavahit á (*Ánde, Márjá, jna.*), ja dat ii leat mihkkege čuolmmaid teknihkalaččat, muhto Álbmotregisttaris lea mearriduvvon geavahit dušše dáru, ruođa ja duiskka bustávid. Brønnøysund-registaris geavahuvvo ng. *Latin 1* -kodatabealla, sámegiela bustávid haga (earret *á*). Ádjána vel máŋga jagi ovdal go Brønnøysund-registariidda lea vejolaš registreret fitnodatnama nugo *Šillju Bistro*.

Lokaliseremii gullá maiddái vejolašvuohta čujuhit iešguđet gillii. Dál juohke sámegiela lea standárda ISO-koda² (lulli-, julev- ja davvisámegiela ISO-639-kodat leat *sma, smj, sme*). Microsoft-fitnodat ii geavat ISO-kodaid, muhto baicca iežas giellakodaid. Windows MS-Office-páhkás lea vejolaš čujuhit buot sámegielaide, muhto Macintosh MS-Office-páhkás eai gávdno giellakodat. Dan dihte lulli-, julev- ja davvisámegiela divvunprográmmaid ferte ohat slováhka-, euskara- (baska-) ja katalánagiela vuolde.

¹ Našunála dárkkistan- ja registrerenvuogádat, mii erenoamážit guoská fitnodagaide ja servviide.

² ISO = International Organization for Standardization

2. Sániid analyseren

Sátni, mii lea muhtun teavsttas oassin, ii leat dušše sátni, muhto sátni ovddasta dihto leksema, grammatihkalaš sáni ja gullá dihto sojahanparadigmii. Gávdnoidit mánggalágan giellateknologiiijat sátneanalysa várás – juohke teknologiiija heive dihto geavahussii ja dihto gielaide. Dás guorahalle máddagastima, statistihkalaš lahkonanvuogi ja morfologalaš analysáhtora.

2.1. Máddagastin

Lunddolaš gielain seammá leksemas sáhttet leat mánga sojahanhámi. Go ohcá dihto leksema teavsttas dahje dokumeanttain, de galgá leat vejolaš gávdnat maiddá sojahuvvon hámiid. Sátnehámit *gahpira*, *gahpiris*, *gahpiriin* leat buot čadnon *gahpir*-leksemii, seamma ládje go dárogiel sánit *hatt*, *hatten*, *hattar*, *hattane*, *hattens*, *hattanes* gullet oktii.

Dábálaš vuohki lea máddagastin (eng. *stemming* ‘máddaga dahkat’). Dárogillii máddagastinprográmma válldášii eret gehčosiid *-en*, *-ar*, *-ane*, *-ens*, *-anes*, ja dalle bázášii *hatt* máttan. Máddagastin doaibmá maiddá bures agglutinerejeaddji gielas, dego omd. turkkagielas, mas olles paradigmas lea seammá mátta (*şapka* ‘gahpir’: *şapka-nunm şapka-ya*, *şapka-yı*, *şapka-da*, *şapka-dan*, ...).

Sámegillii lea maiddá vejolaš ráhkadit máddagasti, muhto sámegiela morfologiiija lea earálágan go turkkagiela morfologiiija go paradigmas sáhttet dássemolsašumi ja diftonganjuolgama dihte leat mánga máddaga. Dán artihkkala oktavuodas moai letne ráhkadan sámegiela máddagasti, ja dainna analyseren substantiivvaid ja vearbbaid sojahanparadigmaid³. Substantiivvain leat máddagastima mañnel gaskamearálaččat 3,1 máddaga, ja vearbain 4,9 máddaga.

Gávdnan dihte man stuora čuovvumušat das leat teakstaanalysii, de letne iskan analyseret vuosttaš 100.000 sáni Norgga Almmolaš Čielggadeami 1994–19:s (NAČ). Morfologalaš analysain (gč. Kap. 2.3) teavsttas leat

³ Substantiivvat leat *baste*, *beana*, *boazu*, *fális*, *gahpir*, *geažus*, *gieddi*, *mánná*, *nieida*, *reñko*, *sabet*, *viessu*, ja vearbbat leat *viežžat*, *diehtit*, *čierrut*, *boradit*, *čohkkát*, *gillet*, *dingot*. Máddagasti lea internehtas: <https://victorio.uit.no/langtech/trunk/gt/sme/src/sme-stemmer.xfst>

5.128 substantiivaleksema. Máddagasti mielde seammá teavsttas leat 9.108 substantiivamáddaga. Dat mearkkaša ahte 44 % máddagiin eai čadnojuvvo rivttes leksemii.

Máddagastinproseassa lea vejolaš buoridit, muhto sámegiela leat dasa hástalussan. Sámegeielain leat rikkis paradigmát, maid lea bággu gieđahallat jus galgá gávdnat rivttes leksema. Sámegeielain leat siskkáldas morfofonologalaš proseassat, dássemolsašupmi, vokálarievdan ja diftonganjuolgan, maid dihte buhtes morfologalaš máddagastima maŋŋel juohke sánis leat máŋga máddaga. Jos ulbmil lea juohke leksemii gávdnat buot sojahanhámiid, de máddagastin ii heive sámegillii.

2.2. Statistihkalaš analysáhtor

Máŋga giellateknologalaš prográmma eŋgelasgiela várás leat ráhkaduvvon unnán grammatihkalaš analysaiguin. Čállindárkkistanprográmmaide leat čohkken buot iešguđetlágan sátnehámiid stuora dárkkistuvvon teakstačoakkáldagas, analysa haga, ja dien ládje prográmma dovda «buot» normatiiva čállinhámiid eŋgelasgillii.

Guokte beali leat dahkan dán vuogi dohkálaš vuohkin: Vuosttažettiin leat ollu eŋgelasgiel teavsttat digitálalaš hámis viežžan ládje. Dasa lassin lea eŋgelasgielas geafes morfologiija. Juohke leksemii gullet dušše moadde sátnehámi, iige oktage dáin leat ollu hárvvit go dat eará hámit. Eanaš gielain máilmmis lea mohkkát struktuvra go eŋgelasgielas – struktuvra mas leksemii gullet viiddis sojahanparadigmat máŋggainlogiin dahje máŋggainčuđiin sátnehámiin, ja muhtun sátnehámit sáhttet leat hárvvenáčat. Statistihkalaš lahkonanvuohki doaibmá buoremusat gielain main gávdnojit hui ollu digitálalaš teavsttat ja unnán morfologiija. Čájehan dihte manne lea nu, de letne iskan guovtti vearbba sojahanhámiid frekveansa eŋgelasgiel, dárogiel ja sámegiela odđa testameanttas. Nubbi vearba lea hui frekveanta, nubbi lea gaskamearálaš frekveanta (gč. tabealla 3).

Tabealla 3. Ođđa testameantta sturrodat golmma gillii – ja vearbbat maid frekveansa letne iskan.

Ođđa testameanta	galle sáni	hui frekveanta vearba	gaskamearálaš frekveanta vearba
davvisámegillii (OT)	138.706	<i>dadjat</i>	<i>bálvalit</i>
dárogillii (DNT)	137.148	<i>si</i>	<i>tjene</i>
engelasgillii (NT)	188.616	<i>say</i>	<i>serve</i>

Tabealla 4 čájeha gallii dárogiel ja engelasgiel veorbbaid sátnehámit gávdnojit ođđa testameanttas. Passiivahámit eai leat mielde danne go sámegeielas passiivaveorbbat gullet suorggideapmái eaige sojaheapmái, nugo engelasgeielas ja dárogeielas. Homonymat leat tabeallas mielde dušše oktii, nuppi ruktui lea merkejuvvon x .

Tabealla 4. Guovtti veorbba frekveansa dárogeiel ja engelasgiel ođđa testameanttas. x = homonyma eará sátnehámiin.

Dárogeiel ođđa testameanttas (DNT)					Engelasgiel ođđa testameanttas (NT)				
Preseansa	<i>sier</i>	467	<i>tjener</i>	26	Preseansa	<i>say</i>	422	<i>serve</i>	33
					Sg2	<i>sayest</i>	39	<i>servest</i>	2
					Sg3			<i>serveth</i>	9
Preterihtta	<i>sa</i>	1028	<i>tjente</i>	7	Preterihtta	<i>said</i>	1058	<i>served</i>	5
Imperatiiva	<i>si</i>	199	<i>tjen</i>	2	Imperatiiva	<i>say</i>	x	<i>serve</i>	x
Infinihtta hámit					Infinihtta hámit				
Perfeakta partisihppa	<i>sagt</i>	136	<i>tjent</i>	11	Perfeakta partisihppa	<i>said</i>	x	<i>served</i>	x
Infinitiiva	<i>si</i>	x	<i>tjene</i>	24	Infinitiiva	<i>say</i>	x	<i>serve</i>	x
					Gerunda	<i>saying</i>	408	<i>serving</i>	6

Engelasgeiela s-hámit *serves* ja *says* eai gávdno danne go dát veršuvdna Ođđa testameanttas (NT) lea čállojuvvon jagis 1611, ja das leat 2. ja 3. persovdna *-st* ja *-th*, eaige leat 3. persovdna s-hámit. Go ohcá dábálaš dán áigásaš teakstačoakkáldagas mas leat 85.000 sáni (bealli ođđa testameanttas), de *says* lea 4 gearddi ja *serves* lea 2 gearddi. Boađus lea ahte buot dán guovtti

vearbba sojahanhámit gávdnojit engelasgiel ođđa testameantta sturrosaš teakstačoakkáldagas. Seammá guoská dárogillii, go dáru ođđa testameanttas leat buot sojahanhámit.

Miellagiddevaš lea buohtastahttit dán bohtosa sámegeielain. Engelasgielas ja dárogielas eai leat go 2–3 sierralágan finihtta hámi, dasa lassin bohtet 1–2 infinihtta hámi. Davvisámegeielas leat juohke vearbba 42 sierralágan finihtta hámi ja dasa lassin 7 infinihtta hámi. Dás eat geahča suorggádusaid go juohke suorggádussii gulašii fas sojahanparadigma. Tabeallas 5 leat *dadjat*-vearbba sátnehámit mat gávdnojit ođđa testameanttas (OT).

Tabealla 5. *dadjat*-vearbba frekveansa davvisámegeiel ođđa testameanttas (OT). Go sátnehápmi ii gávdno, de ruktu lea suivejuvvon. x = homonyma eará sátnehámiin.

	Indikatiiva preseansa		Indikatiiva preterihtta		Konditionála		Potentiála		Imperatiiva	
Sg1	<i>dajan</i>	4	<i>dadjen</i>	0	<i>dajašin</i>	1	<i>dajažan</i>	0	<i>dadjon</i>	0
Sg2	<i>dajat</i>	5	<i>dadjet</i>	x	<i>dajašit</i>	0	<i>dajažat</i>	0	<i>daja</i>	x
Sg3	<i>dadjá</i>	42	<i>dajai</i>	204	<i>dajašii</i>	4	<i>dajaža</i>	0	<i>dadjos</i>	0
Du1	<i>dadje</i>	182	<i>dajaim</i>	0	<i>dajašim</i>	0	<i>dajažetne</i>	0	<i>daddju</i>	0
Du2	<i>dadjabeahhti</i>	1	<i>dajaide</i>	0	<i>dajašeidde</i>	0	<i>dajažeahppi</i>	0	<i>daddji</i>	4
Du3	<i>dadjaba</i>	0	<i>dajaiga</i>	7	<i>dajašigga</i>	0	<i>dajažeaba</i>	0	<i>dadjoska</i>	0
PI1	<i>dadjat</i>	45	<i>dajaimet</i>	1	<i>dajašimmet</i>	0	<i>dajažit</i>	0	<i>dadjot</i>	0
PI2	<i>dadjabehtet</i>	24	<i>dajaidet</i>	0	<i>dajašiddet</i>	2	<i>dajažehpet</i>	0	<i>daddjet</i>	0
PI3	<i>dadjet</i>	47	<i>dadje</i>	182	<i>dajašedje</i>	0	<i>dajažit</i>	x	<i>dadjoset</i>	0
Neg	<i>daja</i>	3			<i>dajaše</i>	1	<i>dajaš</i>	1	<i>daja</i>	x
	Infinitiiva		Aktio essiiva		Perf.part.		Vearba- abessiiva		Gerunda	
	<i>dadjat</i>	x	<i>dadjamin</i>	1	<i>dadjan</i>	14	<i>dajakeahhtá</i>	0	<i>dajadettiin</i>	0

Muhtun grammatihkalaš sániin leat guokte dahje eanet vejolaš sátnehámi. Tabeallain lea dalle mielde sátnehápmi mii gávdno guorahallojuvvon teavsttas.

Nugo boahdá ovdan tabeallas 5, de ođđa testameanttas gávdná dušše 22 sierralágan sojahanhámi *dadjat*-vearbba, ja dat dahká 44 % buot hámiin. Buoremusat gokčá indikatiivva preseansa, muhto dátge vuolleparadigma báhcá váilevažžan.

Tabealla 6 čájeha davvisámegiell *bálvalit*-vearbba bohtosiid. Ođđa testameanttas leat dušše 15 sierralágan sojahanhámi, mat dahket 31 % vejolaš sátnehámiin.

Tabealla 6. *bálvalit*-vearbba frekveansa davvisámegiell ođđa testameanttas (OT). Go sátnehápmi ii gávdno, de ruktu lea suivejuvvon. x = homonyma eará sátnehámiin.

	Indikatiiva preseansa		Indikatiiva preterihitta		Konditionála		Potentiála		Imperatiiva	
Sg1	<i>bálvalan</i>	13	<i>bálvalin</i>	0	<i>bálvalivččen</i>	0	<i>bálvaleaččan</i>	0	<i>bálvalehkon</i>	0
Sg2	<i>bálvalat</i>	0	<i>bálvalit</i>	x	<i>bálvalivččet</i>	0	<i>bálvaleaččat</i>	0	<i>bálval</i>	x
Sg3	<i>bálvala</i>	10	<i>bálvalii</i>	3	<i>bálvalivččii</i>	0	<i>bálvaleažžá</i>	0	<i>bálvalehkos</i>	2
Du1	<i>bálvaleme</i>	0	<i>bálvaleimme</i>	1	<i>bálvalivččiime</i>	0	<i>bálvaležže</i>	0	<i>bálvaleadnu</i>	0
Du2	<i>bálvaleahppi</i>	0	<i>bálvaleidde</i>	0	<i>bálvalivččiide</i>	0	<i>bálvaleažžabeahhti</i>	0	<i>bálvaleahkki</i>	0
Du3	<i>bálvaleaba</i>	0	<i>bálvaleigga</i>	0	<i>bálvalivččiiga</i>	0	<i>bálvaleažžaba</i>	0	<i>bálvalehkoska</i>	0
Pl1	<i>bálvalit</i>	0	<i>bálvaleimmet</i>	0	<i>bálvalivččiimet</i>	0	<i>bálvaleažžat</i>	0	<i>bálvalehkot</i>	1
Pl2	<i>bálvalehpet</i>	0	<i>bálvaleiddet</i>	1	<i>bálvalivččiidet</i>	0	<i>bálvaleažžabehet</i>	0	<i>bálvalehket</i>	2
Pl3	<i>bálvalit</i>	x	<i>bálvaledje</i>	3	<i>bálvalivčče</i>	0	<i>bálvaležžet</i>	0	<i>bálvalehkoset</i>	1
Neg	<i>bálval</i>	5			<i>bálvalivčče</i>	x	<i>bálvaleaš</i>	0	<i>bálval</i>	x
	Infinitiiva		Aktioessiiva		Perf.part.		Vearbaessiiva		Gerunda	
	<i>bálvalit</i>	38	<i>bálvaleamen</i>	2	<i>bálvalan</i>	x	<i>bálvalkeahhtá</i>	0	<i>bálvalettiin</i>	1

Sámegiell ođđa testameanttas leat 139.000 sáni, juoga mii lea unnán dán oktavuodas. Interneahtta dahká stuorit teakstačoakkáldaga, dál leat sullii 75 milj. sámegiell sáni interneahatas⁴ (skábmamánuš 2009). Tabeallas 7 leat merkejuvvon sátnehámit mat leat interneahatas. Dan dihte go muhtun sátnehámit leat homonymat eará gielaid sátnehámiiguin, de ii leat frekveansa mielde tabeallas. 48 sierralágan hámis gávdnojedje 20 interneahatas, ja dat dahká 42 %. Giellatekno (Romssa universitehta sámi giellateknologiija guovddáš) ja Norgga Sámedikki Divvun-joavkku fiillat mat leat olámuttos interneahata bokte, eai leat mielde tabeallas, danne go muhtun fiillain leat genererejuvvon aiddo dáid vearbbaid sátnehámit.

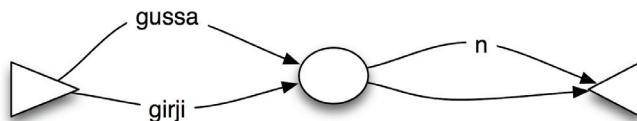
⁴ Interneahta sámegiella teakstameari lea árvoštallojuvvon ná: Giellatekno/Divvuma davvisámegiella teakstačoakkáldagas (6,997 milj. sáni) rehkenasttime man stuora proseantaoasi 34 dábalaš sátnehámi dahke (*bierrgu*, *boahibehtet*, *boahán*, *bohtet*, ...). Bijaimme vuoddu ahte sániid proseantaoassi lea seammá interneahatas, ja nu árvoštalaime interneahata teakstameari.

Tabella 7. bálvalit-vearba interneahntas (01.12.09). Go sátnehápmi ii gávdno, de ruktu lea suivejuvvon. +: sátnehápmi gávdno, 0: ii gávdno, x = homonyma eará sátnehámiin.

	Indikatiiva Preseansa		Indikatiiva Preterihntta		Konditionála		Potentiála		Imperatiiva	
Sg1	<i>bálvalan</i>	+	<i>bálvalin</i>	0	<i>bálvalivččen</i>	+	<i>bálvleaččan</i>	0	<i>bálvlehkcon</i>	0
Sg2	<i>bálvalat</i>	+	<i>bálvalit</i>	x	<i>bálvalivččet</i>	0	<i>bálvleaččat</i>	0	<i>bálval</i>	x
Sg3	<i>bálvala</i>	+	<i>bálvalii</i>	+	<i>bálvalivččii</i>	+	<i>bálvleažžá</i>	0	<i>bálvlehkcos</i>	+
Du1	<i>bálvaletne</i>	0	<i>bálvaleimme</i>	0	<i>bálvalivččiime</i>	0	<i>bálvležže</i>	0	<i>bálvleadnu</i>	0
Du2	<i>bálvleahppi</i>	0	<i>bálvaleidde</i>	0	<i>bálvalivččiide</i>	0	<i>bálvleažžabeahhti</i>	0	<i>bálvleahkki</i>	+
Du3	<i>bálvleaba</i>	+	<i>bálvaleigga</i>	0	<i>bálvalivččiiga</i>	0	<i>bálvleažžaba</i>	0	<i>bálvlehkcoska</i>	0
Pl1	<i>bálvalit</i>	+	<i>bálvaleimmet</i>	0	<i>bálvalivččiimet</i>	0	<i>bálvleažžat</i>	0	<i>bálvlehkot</i>	+
Pl2	<i>bálvlehpēt</i>	+	<i>bálvaleiddet</i>	0	<i>bálvalivččiidet</i>	0	<i>bálvleažžabehtet</i>	0	<i>bálvlehkēt</i>	+
Pl3	<i>bálvalit</i>	x	<i>bálvaledje</i>	+	<i>bálvalivčče</i>	+	<i>bálvležžet</i>	0	<i>bálvlehkcosēt</i>	0
Neg	<i>bálval</i>	+			<i>bálvalivčče</i>	x	<i>bálvlelaš</i>	0	<i>bálval</i>	x
	Infinitiiva		Aktio essiiva		Perf.part.		Vearbaabessiiva		Gerunda	
	<i>bálvalit</i>	x	<i>bálvleamen</i>	+	<i>bálvalan</i>	x	<i>bálvalkeahhtá</i>	+	<i>bálvalettiin</i>	+

2.3. Morfologalaš analysáhtor

Gielaide main leat unnán digítalalaš teavsttat ja ollu morfologiija, lea čoaiddus ráhkadit automáhtaid mat lasihit gehčosiid máddagiidda. Govvosi 1 automáhta addá sátnehámiid *gussa*, *gussan*, *girji*, *girjin*.

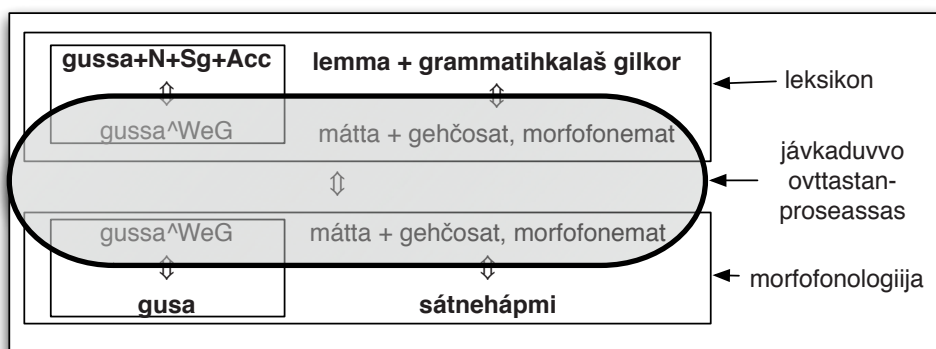


Govv. 1. Morfologalaš automáhta.

Grammatihkalaš transduser (Beesley – Karttunen 2003) lea dihto automáhtatiipa mas juohke sátnehámiis leat guokte ovddasteami. Sátnehámiis *gussan* lea maiddá ovddasteapmi «gussa+N+Ess». Sátnehápmi lea atta (enj. *input*) ja leksema + grammatihkalaš gilkorat (= grammatihkalaš sátni) fas buvttus (enj. *output*), dahje nuppe ládje. Dákkár transduser gohčoduvvo loahpalaš dilletransduserin – FST (enj. *finite state transducer*).

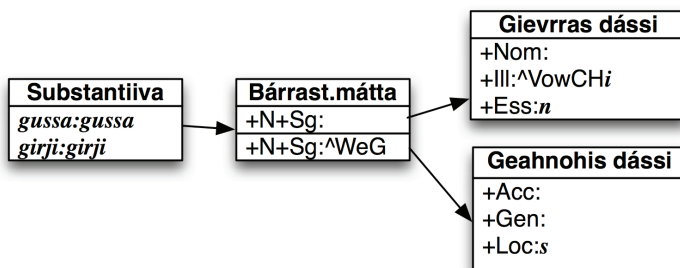
Turkkagielas lokatiivatransduser sátnái *kitab* ‘girji’ sáhttá leat transduserat [kitab:kitab] ja [+Sg+Loc:da], nu ahte atta lea *kitab+Sg+Loc*, ja buvttus *kitabda*. Nugo oinniimet kapihttalas 2.1. máddagastima birra, de dakkár transduser doaimmašii bures turkkagillii. Dat ii doaimmaše bures sámegillii, go jus atta lea *girji+Sg+Loc*, de buvttus livččii *girjis*.

Dan dihte sámegiela analysáhtoris leat guokte transdusera: Nubbi leksikonii ja geažusmorfologijai (leksikontransduser), ja nubbi siskkáldas morfofonologalaš proseassaide (morfofonologalaš transduser, gč. Moshagen – Sammallahhti – Trosterud 2004). Leksikontransdusera bajit dássi ovddasta grammatihkalaš sáni, ja vuolit dásis leat mátta, affivssat ja báhcahagat morfofonologalaš proseassaide maid boadus lea sátnehápmi. Govus 2 čájeha leksikon- ja morfofonologalaš transduseriid. Leksikontransdusera buvttus lea morfofonologalaš transdusera atta.



Govus 2. Analysáhtoris leat guokte transdusera, nuppis lea leksikon ja gehčosiid lasiheapmi, nuppis leat morfofonologalaš proseassat.

Govvosis 2 sátnehápmi *gusa* lea ovdamearkan. Go leksikon- ja morfofonologalaš transduserat leat ovttastuvvon, de gaskabuvttus jávká. Boadus lea transduser, man vuolit dásis lea *gusa* ja bajit dássi lea *gussa+N+Sg+Acc*. Transduseris bajit ja vuolit dásit sirrejuvvojit duppalčuoggáin (*gusa:gussa+N+Sg+Acc*), muhto loahpalaš buvttus lea dan haga, nugo govvosis 5.



Govus 3. Leksikontransduser. Báhcahagat leat \wedge WeG \rightarrow geahnohis dássi ja \wedge VowCHi \rightarrow vokála rievdan. Bajit dási segmeanttat leat duppalčuoggá gurutbealde, ja vuolit dási fas duppalčuoggá olgešbealde.

Leksikontransdusera siskkáldas struktuvra lea govvosis 3, mii čájeha ahte *gussa-* ja *girji-*leksemaid bajit dássái lasihuvvojit sátneluohtkájá kásusgilkorat. Vuolit dássái lasihuvvojit báhcahagat (eng. *trigger*) dássemolsašuddama ja soggevokála rievdamana várás, ja dasa lassin geažus.

Dássemolsašuddan ja vokálarievdan dahkko juogo fonologalaš konteavstta vuodul dahje lasihuvvon báhcahaga dihte. Govvosis 3 lea WeG geahnohis dási báhcahat ja VowCH lea soggevokála rievdamana báhcahat, mat morfofonologalaš transdusera njuolggadusain leat eaktun (gč. govvosa 4).

s s \rightarrow s	/	_	Vow*	WeG
r j \rightarrow r j j	/	_	Vow*	WeG
i \rightarrow á	/	_	VowCh	

Govus 4. Morfofonologalaš transduser. Konsonántaguovddáža molsašuddan lea «/»-symbola gurutbealde, ja eaktu fas olgešbealde. * mearkkaša «nolla dahje eanet».

Kompilašuvnnain jávkaduvvojit leksikon- ja morfofonologalaš transduseriid gaskadásat ovddasteamet (gč. govvosa 2). Boađus lea odđa morfologalaš transduser mainna sáhtta sihke analyseret sátnehámiid ja genereret sátnehámiid. Go analysere *gusa* ja *gussii* sátnehámiid, de oažžu leksema ja grammatihkalaš gilkoriid. Jus manná nuppe guvlui, de generere sátnehámiid *gusa* ja *gussii* (gč. govvosa 5).

gusa	gussa+N+Sg+Acc
gusa	gussa+N+Sg+Gen
gussii	gussa+N+Sg+Ill

Govus 5. Morfoloalaš analysáhtor: Olgeš bealde lea atta, gurut bealde fas buvttus.

Jus eai leat dárkilis ráddjejumit, de sáhtttá genereret hámiid mat eai gávdno. Dakkár riska ii leat statistihkalaš lahkonaivugiin.

Gillii gullet maiddái goallossániid, ja daid ii leat vejolaš leksikaliseret go gillii jámma ihtet ođđa goallosteamiid. Analysáhtoris lea dan dihte dynámalaš substantiiva+substantiiva goallosteapmi, mii maiddái generere goallossániid mat eai gávdno sámegielas, sihke semantihkalaččat ja čállinhámi ektui. Mañit čuovvumuša sáhtttá eastadit njuolggadusaiguin ja ráddjejumiiguin (Moshagen – Omma – Pieski 2009).

Goallossáni mearusoassi sáhtttá leat ol.nom, ol.gen dahje ml.gen. Vuogádaga ovdaválljen lea ol.nom. Leksikonii lea merkejuvvon dalle go leat eará vejolašvuoddat, gč. (1). Goallossáni mearusoasi kasusa sáhtttá stivret maiddái vuodđooasi bokte (2):

- (1) loddi ; +SgNomCmp SgGenCmp (omd. loddebivdu, lottečivga)
- (2) lávlun ; +SgNomLeft +SgGenLeft (omd. sálbmálávlun, sálmma-lávlun)

Goallossáni mearusoassi galgá muhtun goallossániin leat allegrohámis go lea ovttaidlogus, eará sániin fas dat ii galgga leat, dahje sáhtttá leat. Dát čovdojuvvo leksikonas iešguđetlágan joatkkaleksikonaid (bálgáid) bokte. Goallossániid mat eai heive dán vuogádahkii, sáhtttá leksikaliseret.

Giellaoahpalaččat sojahusa ja suorggideami erohus lea ahte sojahusa boađus lea sátnehápmi, ja suorggideami boađus lea ođđa leksema. Juohke substantiivamáddagis manná bálggis kásusmorfologijja čađa, ja máddagii lasihuvvojit kásusgehčosat. Dasa lassin juohke substantiivva ii-nominatiivamáttá joatká omd. diminutiiva-vuolleleksikonii. Diminutiiva-leksikon lasiha máddagii gehčosa (-š- dahje -ž-) ja diminutiivagilkora (+Der/Dimin). Suorggideami boađus sáhtttá maiddái leat ođđa sátneluohtká. Ovdamearkka dihte vearbamáddagat sáhttet joatkit Der/eapmi-leksikonii, nu ahte vearbba-

boahdá substantiiva, nugo *vuoruhit + eapmi => vuoruheapmi*. Teknikkalaččat suorggidahttin ja sojahus gieđahallojuvvojit seammá vuogi mielde.

Go galgá ráhkadit morfologalaš transdusera, de dárbbáša:

- Sátneleisttu digítálalaš hámis, mas juohke sátnái leat lasihuvvon dieđut mátta- ja sojahanluohká birra.
- Morfofonologalaš, dahje rievtti mielde morfografemihkalaš njuolggadusaid, vai sáhtá generaliseret jeavddalaš geažusmolsašumiid.
- Sorterejuvvon listtu mas leat rájálaš luohkáid sátnehámit ja spiehkastagat rabas sátneluohkáin.
- Buori referánsagrammatihka, gos lea vejolaš gávdnat jeavddalaš ja jeavddahis sojahanminstariid.

Giellatekno/Divvuma vásáhusaid vuodul árvvoštalle ahte morfologalaš transdusera ádjána 0,5–1,5 bargojagi ráhkadit, dan mielde man mánggabalaš morfologiija gielas lea. Go vuos lea ráhkadan dakkár analysáhtora, de sáhtá álkít kompíleret iešguđetlágan variánttaid, omd. deskriptiiva vs. normatiiva analysáhtora.

2.4. Čeahkkáigeassu

Dán kapihttalis leat čilgejuvvon golbma lahkonanvuogi sátneanalysere-mii. Gielaide main lea unnán morfologiija sáhtá geavahit máddagasti diehtoohcamii, ja statistihkalaš lahkonanvugiin ráhkadit sátnedárkkistan-programma. Sáme-gielaide ii goabbáge vuohki heive. Dan sajis leat válljen oahpahit grammatihka dihtorii nu ahte sáhtá analyseret ja genereret juohke teorehtalaččat vejolaš sojahanhámi. Lea stuora bargu dan dahkat, muhto Kap. 4 čájeha mo grammatihkalaš analysáhtoriin sáhtá ráhkadit mánga ávkkálaš geavaheaddjiprogramma.

3. Cealkagiid analyseren

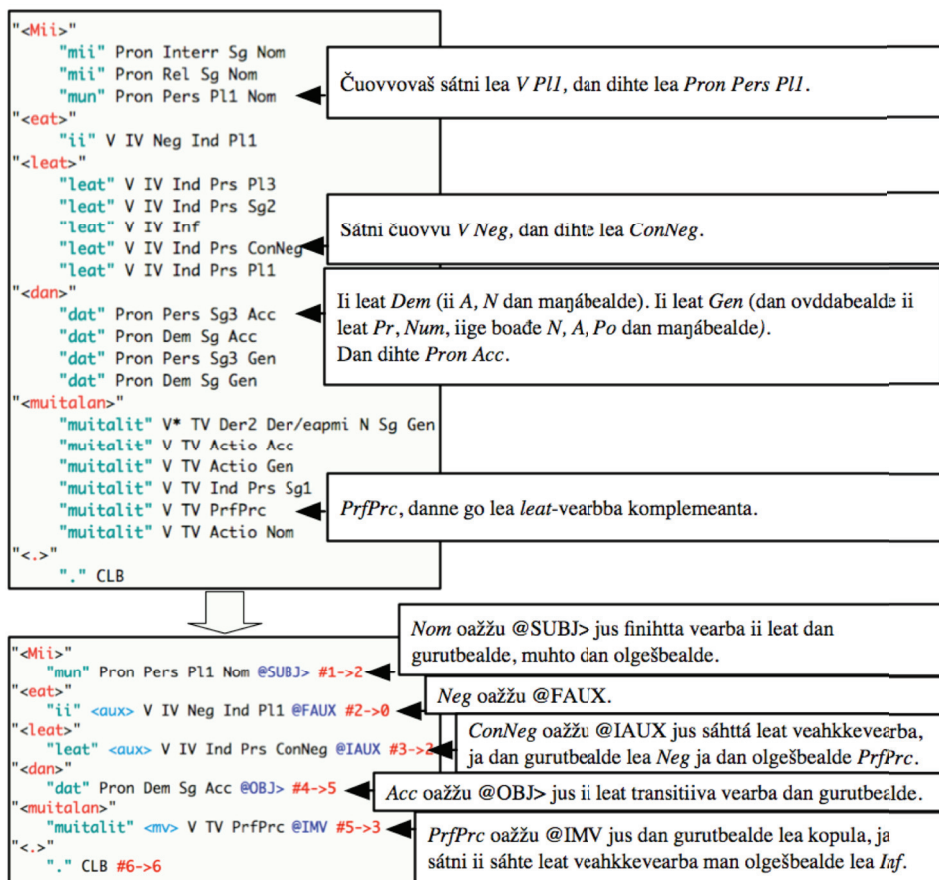
3.1. Homonymiijaid disambigueren

Oktilis teavsttas lea ollu homonymiija. Davvisámegiellat sátnehámiin leat gaskamearálaččat 2,7 vejolaš morfologalaš analysa, julevsámegiellat sátnehámiin fas 2,2 vejolaš analysa. Davvisámegiellat leat eanet homonymat earenoamážit danne go julevsámegiella lea konservatiivvalaččat go davvisámegiella. (Trosterud – Wiechetek 2007.)

Go olmmoš gullá sámegiella, de hárve lea eahpečielggas maid hubmi oaivvilda, danne go sánit gullet man nu kontekstii, ja olbmossat leat sihke syntávssalaš ja semánttalaš dieđut. Maiddá automáhtalaš analysa galggašii leat guovttečilggolašvuodaid haga.

Homonymiija sáhtá disambigueret sihke statistihkalaččat ja grammatihkalaččat. Kapihttalis 2 letne ákkastallan manne leat válljen bargat grammatihkalaččat eatge statistihkalaččat. Grammatihkalaččat sáhtá bargat guovtti ládje, badjin–vulos, dahje vuollin–bajás. Badjin–vulos-analysáhtor geahččaladdá hypotesaid das makkár syntávssalaš struktuvra cealkagis lea, ja disambiguere dan mielde. Muhto lunddolaš giella ii álo čuovo syntávssalaš njuolggadusaid; hubmi sáhtá rievdatit oaivila gaskan cealkaga, son sáhtá lasihit guhkes čilgehusaid, dahje cealkka šaddá njulgestaga ilá kompleaksa. Dalle badjin–vulos-analysáhtor ii gávna makkárge analysa. Maiddá buorit badjin–vulos-analysáhtorat hárve máhttet analyseret eanet go 60 % cealkagiin oktilis teavsttas.

Giellatekno analysáhtor geavaha grammatihkalaš vuollin–bajás-analysáhtora mii nagoda addit analysa maiddá cealkafragmeanttaide ja kompleaksa cealkagiidda. Gávnojit mángga lágan analysáhtorat, ja Giellatekno analysáhtor vuodđuduvvá ráddjehusgrammatihkkii (*constraint grammar*) (Karlsson – Voutilainen – Anttila 1995; Tapanainen 1996). Grammatihkas leat ráddjehusat (*constraints*) das makkár konteavsttain iešguhtet morfologalaš analysa sáhtá leat. Ráddjehusgrammatihkkii gullet njuolggadusat, ja njuolggadusat kompilerejuvvojit vislcg3-prográmmain (Visl-group 2008). Govus 6 čájeha mo njuolggadusat doibmet, omd. go *mii*-sátnehápmi lea cealkaga álggus ja *eat*-vearbba ovddabealde, de lea rivttes analysa



Govus 6. Ráddjehusnjuolggadusaiguin morfologalaš analysa disambiguerejuvvo ja dasa lassin lasihuvvojit syntáksagilkorat, main lea @, ja dependeansagilkorat: #iežas sajádat → oaivvi sajádat.

Muhtun dáhpáhusain lea dárbbášlaš maiddái čujuhit sániid semánttalaš sisdollui:

- (3) *Mun boran eatni gievkkanis.*
- (4) *Mun boran láibbi gievkkanis.*

Cealkagis (3) lea *eatni* genetiiva danne go čujuha olbmui, ja cealkagis (4) *láibbi* lea akkusatiiva seammá syntávssalaš konteavsttas. Vai analysáhtor sáhtá válljet rivttes kásusa, de sáhttit lasihit gilkorra (<hum>) buot leksemaide mat čujuhit olbmuide, dahje čohkket diekkár leksemaid seahtaide. Ráddjehusgrammatihka njuolggadus čujuha juogo gilkorii dahje sehttii, ja sáhtá dainna lágiin váldit vuhtii sáni semánttalaš sisdoalu.

Davvisámegiela lea akkusatiivvas ja genetiivvas homonymiija olles substantiivaparadigmas, ja opposišuvdna lea dušše lohkosániin. Morfologalaš analysáhtor addá dattetge sihke akkusatiiva- ja genetiivagilkoriid maiddá substantiivvaide. Dán homonymiija lea váttis disambigueret danne go akkusatiiva ja genetiiva sáhttet leat seammá syntávssalaš konteavsttas, nugo bajábeal ovdamearkkas. Mánngga oktavuodas lea dattetge ávkkálaš earuhit akkusatiivva ja genetiivva, omd. jus dihtor galgá máhttit jorgalit gielas nubbái. Maiddá dalle go automáhtalaččat lasiha syntávssalaš funkšuvdnagilkoriid sániide, lea ávkkálaš earuhit akkusatiivva ja genetiivva.

Davvisámegiela morfologalaš analysáhtor gárvánii vuosttažin, ja dan dihte leatge bargan eanemusat davvisámi disambiguáhtoriin. Dan F-score⁵ lea 0,99 sátneluohkkádisambigueremis ja 0,94 olles morfologalaš disambigueremis. Julevsámegiela F-score leat fas 0,95 ja 0,88 (skábmamánu 2009). Eará sámegielaide eat leat vuos ráhkadan disambiguáhtora.

3.2. Syntávssalaš gilkorat

Mii geavahit maiddá ráddjehusgrammatihka merket syntávssa ja dependeansa gilkoriid sániide. Sámegiela lea viehka friddja sátnortnet, ja mii leat válljen addit olles 50 funkšuvdnagilkora. Omd. lasihuvvojit iešguđetlágan subjeaktagilkorat dan mielde leago finihtta vearba olgeš vai gurut bealde, ja leago infinihtta vearbba subjeaktan. Disambiguáhtor ja syntávssalaš analysáhtor leat muhtun muddui ovttahttojuvvon danne go muhtun syntávssalaš gilkorat leat veahkkin disambigueremiin. Davvisámegiela syntávssalaš analysáhtoriin oažžu F-score 0,93, julevsámegiela analysáhtoriin fas 0,86 (skábmamánu 2009).

Álggus mii ráhkadeimmet sierra syntávssalaš analysáhtora juohke sámegillii, muhto leat ovttahttigoahtán buot sámegielaid ovttá analysáhtorii. Syntávssalaš erohusat gielaid gaskkas váldojuvvojit vuhtii go lea vejolaš čujuhit morfologalaš analysáhtora buktosa giellagilkorii (<sme> <smj> <sma>). Omd. lullisámegiela kopula ii dárbbas leat mielde eksistentiála iige habitiiva cealkagis, ja dalle cealkagis ii leat ollenge finihtta vearba. Davvisámegiela

⁵ F-score lea vuohki mihtidit dárkilvuoda, ja presišuvnna ja deaivanmeari gaskamearálaš logu. Presišuvnna oažžu go juohká «rivttes analysa» -logu buot bohtosiiguin, ja deaivanmeari go juohká «rivttes analysa» -logu buot analysaiguin maid dihtor livččii galgan gávdnat. Buoremus vejolaš F-score lea 1, heajumus lea 0.

cealkka finihhta vearbba haga analyserejuvvošii fragmeantan dahje elliptalaš cealkkan. Nubbi ovdamearka lea habitiiva kásus mii lea iešguđetlágan gielas gillii – lokatiiva (davvisámegielas), inessiiva (julevsámegielas) ja genetiiva (lullisámegielas). (Antonsen – Trosterud – Wiechetek 2010.)

Min syntávssalaš gilkorat vuodđuduvvojit kompromissii sámi grammatihkalaš árbevieru ja ráddjehusgrammatihka konvenšuvnnaid gaskkas. Válderohus lea ahte ráddjehusgrammatihkka lea lineára vuogádat, ja gilkorat addojuvvojit sátnehámiide, eaige gihpuide:

(5) *Mu.*@>N *čeahci.*@SUBJ> *duddjo.*@+FMAINV *čáppa.*@>N *niibbi.*@<OBJ.

Cealkagis (5) dušše *čeahci* oažžu @SUBJ> gilkora (@-mearka earuha syntávssalaš gilkora morfologalaš gilkoriin). *Mu* oažžu @>N, mii muitala ahte sátnehápmi modifisere substantiivva mii lea olgeš bealde. Ja de ferte lohkki dulkot @>N ja @SUBJ> kombinašuvnna gihppun.

Syntávssalaš analysáhtora buktosis sáhtttá ráddjehusgrammatihkain dahkat dependeansamuora dependeansagilkoriiguin, nugo govvošis 6. Ruohtas lea merkejuvvon 0. Cealkaga sánit nummarastojuvvojit dađistaga ja njuolla čujuha sáni oaivái, omd. 2→1, mii máksá ahte cealkaga nubbi sátni lea cealkaga vuosttaš sáni dependeanta.

Giellatekno/Divvun lea ráhkadeamen analysáhtoriid sámegeielaide. Tabealla 8 čájeha man muttus bargu lei skábmamánuš 2009.

Tabealla 8. Sámegeielaide analysáhtorat. Dan sajis go ráhkadit sierra syntávssalaš analysáhtora lullisámegiela várás, de leat ovttahttiigoahtán buot sámegeielaide ovtta analysáhtorii, mii ii leat vuos gárvvis.

Giella	Morfologalaš analysáhtor	Syntávssalaš analysáhtor	Dependeansa-analysáhtor
Davvisámegiella	101.000 leksema (57.000 leksema)	3450 njuolggadusa – 50 syntávssalaš gilkora	270 njuolggadusa
Julevsámegiella	47.000 leksema	874 njuolggadusa	
Lullisámegiella	30.000 leksema – válmmaš 2010	–	

4. Geavaheaddjiprográmmat

Grammatihkalaš analysáhtor geavahuvvo juo geavaheaddjiid prográmmain.

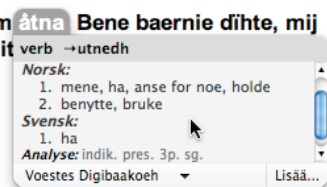
Čállindárkkistanprográmma Divvun (Gaup – Moshagen – Omma – Palismaa – Pieski – Trosterud 2006) vuodđuduvvá morfologalaš automáhtaide mat genererejit buot normatiiva sátnehámiid nugo čilgejuvvo kapihttalis 2.3. Prográmma válbmanii davvisámegillii ja julevsámegillii jagi 2007 loahpageahčen, ja jagi 2010 gárvána maiddá lullisámegiel veršuvdna. Divvun-prográmmas lea funkšuvdna mii geavaheaddjái evttoha riehta čállon sátnehámiid dihto algoritmma mielde.

Sámegielaid siskkáldas morfofonologalaš proseassat dahket oahpahallái vátisin diehtit guđe leksemii sátnehápmi gullá. Dušše 7,9 % sániin davvisámegiel oktilis teavsttas leat seammaláganat go lemmahápmi (Antonsen – Gerstenberger – Moshagen – Trosterud 2009). Dan dihte Giellatekno/Divvun lea morfologalaš generáhtoriin lasihan sihke leksemaid sojahanhámiid ja grammatihkalaš sániid *Vuosttaš Digisánit-* ja *Áarjelsaemien Digibaakoeh*-nammasaš digitálalaš sátnegirjiide. Sátnegirji sihke mitala guđe grammatihkalaš sánis lea sáhka (gč. govvosa 7), ja addá dieđuid sojahanparadigma birra, nu gohčoduvvon čoavddahámiid, vai geavaheaddji oaidná morfofonologalaš proseassaid mat gullet leksemii. Dákkár sátnegirji lea álki ráhkadit go gielas lea morfologalaš analysáhtor, ja go gávdno guovttegielalaš sátnelistu maid sáhtta geavahit.

Øhpehtimmie ihkuve aajkan

Jarle Jonassen áarjelsaemien báatsoeburrie jñh fuelhkiem **átna** Bene baernie dñhte, mij bovtsiluvnie. Jñjtje lea politihkerinie abpe tijjen. Akte polit raasth særta, jallh dejtie nááhkehte jis daerpies.

Almmuhuvvon: 24.08.2009



Govus 7. Voestes Digibaakoeh-prográmma dovdá sojahuvvon sániid. Prográmma sáhtta viežžat <http://giellatekno.uit.no/words/dicts/index-sma.sme.html>

Giellatekno lea ráhkadan giellaoahppanprográmmaid, OAHPA ja VISL (Antonsen – Baal – Trosterud 2009, Antonsen – Huhmarniemi – Trosterud 2009). Morfologalaš ja syntávssalaš analysa dahká vejolažžan ii dušše mitalit leago oahppi vástádus riehta vai ii, muhto maiddá analyseret

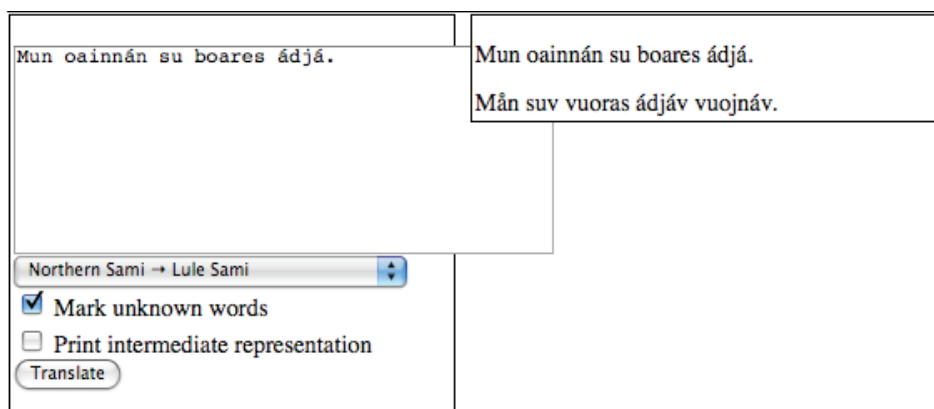
boasttuvuođaid. Guovtti oasseprográmmas (*Vasta* ja *Sahka*) oahppi beassá ieš formuleret vástádusa dihtora jearaldagaide. Dihtora jearaldat ja oahpahalli vástádus addojuvvojit oktan analysáhtorii, mii analysere daid morfologalaččat ja syntávssalaččat. Analysáhtor lasiha guovttelágan gilkkoriid, nugo govvosis 8:

- Semánttalaš gilkkora mii muitala maid oahppi lea vástidan vai dihtora čuoovvovaš jearaldat dahje kommentára heive oahppi vástádussii (&dia-hivsset danne go oahppi lea vástidan ahte TV galgá hivssegii).
- Feailagilkkora mii muitala ahte oahppi lea geavahan boasttu sojaheami. Dáinna son oažžu bagadeaddji máhcahusa sierra lásas (&grm-missing-III danne go oahppi vástádusas ii leat illatiiva).

```
"<Gude>"
  "guhte" Pron Interr Sg Gen &grm-missing-III
"<latnjii>"
  "latnja" N Sg Ill
"<moai>"
  "mun" Pron Pers Du1 Nom
"<bidje>"
  "bidjat" V TV Ind Prs Du1
"<mu>"
  "mun" Pron Pers Sg1 Gen
"<TV>"
  "TV" N ACR Sg Acc
"<Aqst>"
  "^sahka" QDL gosa_bidjat_TV &dia-hivsset
"<Moai>"
  "mun" Pron Pers Du1 Nom
"<bidje>"
  "bidjat" V TV Ind Prs Du1
"<TV>"
  "TV" N ACR Sg Gen
"<hivssegis>"
  "hivsset" N Sg Loc
"<.>"
  "." CLB
```

Govus 8. Sahka-prográmma lasiha jearaldat/vástádus-bárrii sihke navigeren- ja grammatihkkagilkkora. Gilkkoriid vuodul vuogádat addá máhcahusa oahppái. Prográmma lea interneahtras: <http://oahpa.uit.no>

Mášenjorgaleamis geavahuvvo eanaš statistihkalaš lahknanvuohki, dego Google Translate-programmas. Dakkár teknologiija ii heive sámegillii, go dárbbáša stuorát teakstačoakkáldagaid go mat davvisámegillii gávnojit. Sámegillii lea baicca lunddolaš geavahit grammatihkalaš analysáhtoriid. Mášenjorgalanprogramma Apertium lea ráhkaduvvon jorgalit fuolkegielaid gaskkas (nugo katalána- ja spánskkagiela gaskkas). Apertium-programma doaibmá rabas gáldokodaprinsihpa mielde (<http://wiki.apertium.org>), ja dihtorlingvisttat olles máilmmis oassálastet bargui lasihit odđa giellapáraid dasa. Dál leat 35 doaibmi giellabára, ja 27 jorgalanbára leat álggahuvvon. Daidda gullá prošeakta jorgalit davvisámegielaš julevsámegillii (Tyers – Wiechetek – Trosterud 2009) (gč. govvosa 9).



Govus 9. Apertium-programmain lea vejolaš jorgalit davvisámegielaš julevsámegillii. Programma lea internehtas: <http://xixona.dlsi.ua.es/testing/>

5. Servodatlaš váikkuhusat

5.1. Giellateknologiija servodagas

Dábáleamos giellateknologalaš heiveheapmi lea dokumeantaohcan. Dokumeantaid siskkáldas vuoruheapmi lea mearriduvvon liñkastruktuvrra bokte: dat dokumeanta vuoitá masa čujuhit mánga liñkka. Ovdalgo vuorua, de ferte gávdat áššáiguoskevaš dokumeantaid. Sámegeiela oktavuodas morfologalaš variašuvdna lea dehálaš faktor: jos ohcá sáni *giella*, de giellateknologiija haga ii gávna sániid *gielas*, *gillii*, *giellan*. Máddagastima bokte mañimus sátnehámi goitge gávdná (-n váldojuvvo

eret). Davvisámegiela máddagastima bokte substantiivvain leat gaskamearálaččat 3,1 máddaga (gč. kap. 2.1), nu ahte 2/3 vejolaš relevánta sátnehámiin leat oaidnemeahtumat. Morfologalaš analysa livččii dan dihte dehálaš dokumeantaohcamii.

Buriid giellateknologalaš resurssaid haga ii oktage giella sáhte boahtteáiggis:

- doaibmat hálddahaslaš giellan. (Teaksta ferte dárkkistuvvot, mánggalágan skovit ráhkaduvvot, olbmot dárbbasit čoahkkái-geasuid. Buot dát galgá dahkkot automáhtalaččat, giellateknologalaš reaiduiguin.)
- doaibmat guovttegielalaš hálddahasas. (Máנגgagielalaš veršuvvnat seamma teavsttas ráhkaduvvojit automáhtalaččat mášenjorgaleami bokte. Teavsttain dárkkistuvvo automáhtalaččat ahte tearbmageavaheapmi lea konsekveanta.)
- vurket dokumeanttaid digitála arkiiivvaide. (Gáldomášiinnat geavahit giellateknologiija sihke dokumeanttaid ja dieđuid klassifiseremii ja maiddái vurkejuvvon teavsttaid gávdnamii. Livččii maiddái vuogas jus sáhtášii oheat dihto sisdoalu beroškeahtá das guđe gillii teaksta lea čállojuvvon.)

Go oaidná vejolašvuodaid maid dihtor attášii, de olmmoš ákkastalašii geavahit dárogiela ja eará eanetlogu giela dakkár sajis masa ii gávdno giellateknologiija. Manne galgá vurket dokumeantta dakkár čállojuvvon gillii mii dahká dan veadjemeahttumin fas gávdnat?

Maŋimuš logi jagi leat dáhpáhuvvan stuora ovdáneamit suorggis mas leat omd. semánttalaš fierpmádagat ja sániid mearkkašumiid (*word-sense*) disambigueren⁶. Dát bargu ii leat vuos ihtán dábálaš geavaheaddji dihtorii, muhto dat boahdá, ja jus mii eat daga maidege, de dat boahdá dušše stuora gielaide. Dađistaga go gávppálaš giellateknologalaš prográmmat buorránit, ja go almmolaš hálddahas vuodđuduvvá giellateknologiijai, de stuorru erohus gielaid gaskkas main leat ja main eai leat dakkár resurssat.

Dáidda geavatlaš čuovvumušaide lassin, de mii lingvisttat maiddái sáhttit

⁶ Omd. sánis fierbmi leat mánga mearkkašumi, ja go disambiguere, de gávdná guđe mearkkašumis lea sáhka dihto konteavsttas.

ávkkástallat resurssaiguin: Bures ovdánan giellateknologiija dagašii vejolažžan guorahallat giela ollu beaktileappot go giellateknologiija haga. Ja de lea vel duođaid dehálaš sivva: Eanaš gielaide máilmmis ii leat gávppálaččat gánnáhahtti dahkat vuodđobarggu mii dárbbasuvvo giellateknologalaš geavaheaddjiprográmmaide. Bargu báhcá akademalaš dutkanásahusaid vuoruheami duohkai. Muhto lingvistii leage mávssolaš bargu čállit ollislaš referánsagrammatihka nu ollu gielaide go vejolaš.

5.2. Olaheamit ja ođđa hástalusat

Sámegielaide stuorimus servodatlaš váikkuhus dán rádjái lea leamaš dihtora vuogádaga ja boallobeavddi lokaliseren. Unicode-standárdda, Sámi dihtorlávdegotti ja Ruota standardiserenlávdegotti ánsun juohke dihtoris leat dál sámegiel boallobeavdi ja bustávát.

Čuovvovaš dehálaš boadus lei Divvun-prográmma, mii ee. lea dahkan vejolažžan buvttadit sámegiel aviissaid jođáneappot (Guhkes lávki ovddasguvlui 2007), ja maid Ávvira váldodoaimmaheaddji atná buoremus veahkkeneavvun maid sámegielat aviisa lea ožžon ođđa áiggis (Mii dárbbasit «čuorbbi» 2008).

Boahtteáiggis májggagielašvuohta šaddá dehálaččat aht' dehálaččat. Norgga stáhtahálddaha ođđa giellapolitihka mielde Norgga stáhta lea internehtas njealjegiela (guokte dárogiela, engelasgiella, davvisámegiella)⁷. Dakkár politihka lea guhkit áigeperspektiivvas vejolaš čađahit dušše mášenjorgalemiin.

Sámi ásašusain ja servviin lea muhtumin ággan čállit dárogillii dat ahte eai buot lohkkat hálddaš sámegiela. Ja nu sámegiella dávjá báhcá jorgalusgiellan, dan muddui go lea áigi ja resurssat dasa. Mášenjorgalanprográmma sámegielas dárogillii dagašii vejolažžan čállit njuolggá sámegillii, go prográmma veahkehivččii olbmuid geat eai ipmir sámegiela.

Stuora hástalus lea jorgalit doarvái teavsttaid julevsámegillii ja lullisámegillii. Sámelága giellanjuolggadusaid ollašuttimii eai gávdno doarvái čállit

⁷ Geahča omd. <http://rádđehus.no/>. Neahttasiiddus leat njeallje giellamolssaeavttu: Bokmål, Nynorsk, Sámegiella, English.

eaige jorgaleaddjit julev- ja lullisámegillii. Sáhtta leat váttis geargat buot skuvlagirjiid jorgalahttit ovdalگو čuovvovaš oahppoplánaođastus boahá. Buorre veahkin dan bargui livččii mášenjorgaleapmi davvisámegiela bokte.

6. Konklusuvdna

Buot digitálalaš sámegiela teakstagedáhallaan vuodđuduvvá dasa ahte lea vejolaš čállit ja lohkat sámegiela digitálalaš hámis. Mañimuš logi jagi lea dát leamaš vejolaš buot sámegielaide, muhto eai buot almmolaš registarat leat atnigoahtán ođđa teknologiija, nu ahte sámegielat dieđuid ii velge leat vejolaš registreret omd. Brønnøysundregistariidda iige Álbmotregistarii.

Mañimuš jagiid leat boahtán syntávssalaš ja morfologalaš analysaprográmmat davvi-, julev- ja lullisámegillii. Nugo čájehuvvon 2. kapihttalis, de sámegielaide ii sáhte analyseret seammá lahkonañvugiiguin go gielaid main ii leat nu rikkis morfologiija. Sámegielaide grammatihkalaš lahkonañvuohki lea dárbbalaš. Analysaprográmmat barget guovtti proseassas, nubbi lasiha gehčosiid sáni máddagii ja nubbi dahká morfofonologalaš proseassaid. Boađus lea prográmma mii sáhtta analyseret ja genereret juohke sámegiel sátnehámi. Syntávssalaš analyserenprográmma vállje rivttes morfologalaš analysa, addá syntávssalaš funkšuvnna ja dependeansaanalysa dan vuodul makkár relatiiva sajadat sátnehámis lea cealkagis.

Analysáhtorat veahkehit min ipmirdit sámi grammatihka, muhto dain leat maiddái dát dehálaš doaimmat: Daiguin sáhtta analyseret stuora teakstačoakkáldagaid vai gielladutkit sáhttet generaliseret. Ja daid vuodul lea vejolaš ráhkadit prográmmaid mat leat dehálaččat buot sámegiela-giidda: Divvunprográmmaid, interaktiivalaš sátnegirjiid, giellaoahppan-prográmmaid ja dieđuidgedáhallañprográmmaid.

Dáin bohtosiin leat servodatlaš čuovvumušat. Jus sámegiella galgá leat mielde ođđa máñggagielalaš servodagas, ja jus davvisámegiella ja eará sámegielat galget sáhttit doaimmat hálddahuñgiellañ ja oahpahuñgiellañ, de sámegielain fertedit leat seammá buorit veahkkeneavvut go eará gielain. Lea ain ollu bargu dan muddui ollet, muhto vuodđobargu lea juo dahkkon – dihtor máhtta sámegiela grammatihka.

Materiála

- DNT = (2005) Det nye testamentet. – <http://www.bibelen.no>
NT = (1611) New Testament, King James Bible. – <http://www.kingjamesbibleonline.org>
NAČ = (1994) Finnmárkku eatnamiid ja čázádagaid geavaheapmi historjjálaš geahččamiin. Norgga almmolaš čielggadeamit 1994:21 S. – http://www.regjeringen.no/se/dep/krd/Dokumeanttat/NA-at/1993/nac_199421.html?id=139744
OT = (1998) Ođđa testameanta. – <http://www.bibelen.no>

Girjjálašvuhta

- ANTONSEN, LENE – BAAL, BERIT ÁNNE BALS – HUHMARNIEMI, SAARA – TROSTERUD, TROND 2009: Dihtor ja giela válljenvejolašvuodát – gielalaš ja pedagogalaš čuolmmat. – Johanna Ijäs – Nils Øivind Helander (doaim.), *Sáhkavuoruin sáhkan. Sáme giela ja sámi girjjálašvuoda muhtin áige guovdilis dutkanfáttát* s. 87–102. Dieđut 1/2009. Guovdageaidnu: Sámi allaskuvla.
- ANTONSEN, LENE – GERSTENBERGER, CIPRIAN-VIRGIL – MOSHAGEN, SJUR NØRSTEBØ – TROSTERUD, TROND 2009: Ei intelligent elektronisk ordbok for samisk. – *LexicoNordica* Volum 16 s. 271–283. Oslo: Nordisk forening for leksikografi.
- ANTONSEN, LENE – HUHMARNIEMI, SAARA – TROSTERUD, TROND 2009: Interactive pedagogical programs based on constraint grammar. *Proceedings of the 17th Nordic Conference of Computational Linguistics*. Nealt Proceedings Series 4. Tartu: Tartu University Library. – <http://hdl.handle.net/10062/9546> (16.04.10).
- ANTONSEN, LENE – TROSTERUD, TROND – WIECHETEK, LINDA (2010): Reusing Grammatical Resources for New Languages. *Proceedings of the International conference on Language Resources and Evaluation LREC 2010*. Stroudsburg: The Association for Computational Linguistics. – http://www.lrec-conf.org/proceedings/lrec2010/pdf/254_Paper.pdf (10.06.10).
- BEESEY, KENNETH R. – KARTTUNEN, LAURI 2003: *Finite State Morphology*. Stanford, California: CSLI publications in Computational Linguistics.
- GAUP, BØRRE – MOSHAGEN, SJUR N. – OMMA, THOMAS – PALISMAA, MAREN –

- PIESKI, TOMI – TROSTERUD, TROND 2006: From Xerox to Aspell: A First Prototype of a North Sámi Speller Based on TWOL Technology. – Anssi Yli-Jyrä – Lauri Karttunen – J. Karhumäki (doaim.), *Finite-State Methods and Natural Language Processing*. Lecture Notes in Computer Science 4002, s. 306–307. Berlin – Heidelberg: Springer-Verlag. – <http://www.springerlink.com/content/an651qt0g45k55u1/> (16.04.10).
- Guhkes lávki ovddasguvlui [váldočála]. 2007. – *Min Áigi*, 2.6.2007, Nr 42, s. 2.
- KARLSSON, FRED – VOUTILAINEN, ARTO – HEIKKILÄ, JUHA – ANTTILA, ARTO 1995: *Constraint grammar. A language-independent system for parsing unrestricted text*. Berlin – New York: Mouton de Gruyter.
- Mii dárbbasit «čuorbbi» [váldočála]. 2008. – *Ávvir*, 15.11.2008, Nr 144, s. 2.
- MOSHAGEN, SJUR – OMMAN, THOMAS – PIESKI, TOMI 2008: *Goallosteapmi Divvun-reaidduin*. Tromsø: Universitetet i Tromsø. – http://giellatekno.uit.no/background/Goallosteapmi_Divvun.pdf (30.11.2009).
- MOSHAGEN, SJUR – SAMMALLAHTI, PEKKA – TROSTERUD, TROND 2004: Twol at work. – Antti Arppe – Lauri Carlson – Krister Lindén – Jussi Piitulainen – Mickael Suominen – Martti Vainio – Hanna Westerlund – Anssi Yli-Jyrä (doaim.), *Inquiries into Words, Constraints and Contexts* s. 94–105. Stanford, California: CSLI.
- TAPANAINEN, PASI 1996: *The Constraint Grammar Parser CG-2. Publications of the Department of General Linguistics, 27*. Helsinki: University of Helsinki.
- TROSTERUD, TROND – WIECHETEK, LINDA 2007: Disambiguering av homonymi i nord- og lulesamisk. – Jussi Ylikoski – Ante Aikio (doaim.), *Sámit, sánit, sátnehámit. Riepmočála Pekka Sammallahtii miessemánu 21. beaivve 2007* s. 401–421. Suomalais-Ugrilaisen Seuran Toimituksia 253. Helsinki: Suomalais-Ugrilainen Seura. – http://www.sgr.fi/sust/sust253/sust253_trosterudjawiechetek.pdf (16.04.10).
- TYERS, FRANCIS M. – WIECHETEK, LINDA – TROSTERUD, TROND 2009: Developing Prototypes for Machine Translation between Two Sámi Languages. *Proceedings of the 13th Annual Conference of the European Association of Machine Translation, EAMT09*. Allschwil: European Association for Machine Translation. – <http://www.mt-archive.info/EAMT-2009-Tyers-1.pdf> (16.04.10).
- VISL-GROUP 2008: *Constraint Grammar*. Odense: University of Southern Denmark. – http://beta.visl.sdu.dk/constraint_grammar.html (30.11.09).

Why the computer should know its Sami grammar

Language technology constitutes the foundation for the necessary infrastructure needed for any language to function in a modern literary society. The Sami languages differ from the languages for which most such technology is developed in two important ways: The body of text available (either Sami or bilingual Sami – majority language) constitutes but a fraction of what is available for Western European state languages, and the Sami languages have morphological structures far more complex than the ones for most of the Western European state languages.

The article argues that the answer to this challenge is to build a grammar-based language technology for the Sami languages, and presents ongoing work fulfilling this goal. It is shown how morphophonological processes and inflectional and derivational morphology may be modelled as finite-state transducers, and combined with a syntactic component consisting of context-sensitive constraint grammar rules, to constitute a robust grammatical analyser capable of both analysing running text, and generating any word form. The speech communities of the Sami languages are not large enough to uphold a language technology industry, but the grammar-based language model is interesting for theoretical linguists as well.

Practical applications derived from the basic grammatical analysers include spell-checkers, interactive computer-assisted language learning programs, and machine translation.

Lene Antonsen – Trond Trosterud

Romssa universitehta / University of Tromsø

lene.antonsen@uit.no

trond.trosterud@uit.no