

Online morfológiai elemzők és szóalak-generátorok kisebb uráli nyelvekhez

Bakró-Nagy Marianne², Endrédi István¹, Fejes László², Novák Attila¹, Oszkó Beatrix², Prószéky Gábor¹, Szeverényi Sándor², Várnai Zsuzsa², Wagner-Nagy Beáta³

¹MorphoLogic

1116 Budapest, Kardhegy utca 5.

{endredy, novak, proszeky}@morphologic.hu

²MTA Nyelvtudományi Intézet

1068 Budapest, Benczúr utca 33.

{bakro, fejes, oszko, szeverenyi, varnai}@nytud.hu

³Universität Hamburg

Institut für Finnougristik/Uralistik

Johnsallee 35, 20148 Hamburg

beata.wagner-nagy@uni-hamburg.de

Kivonat: Cikkünkben egy olyan webhelyet mutatunk be, amelyen több korábbi projekt keretein belül számos kisebb uráli nyelvre készített morfológiai elemzőket, szóalak-generátorokat és korpuszokat tettünk elérhetővé. Az elemzések a webes felületen egy morfológiai és szemantikai egyértelműsítő eszköz formájában jelennek meg és minden nyelvhez virtuális billentyűzet segíti a szövegbevitelt.

1 Bevezetés

Az MTA Nyelvtudományi intézetének Finnugor és Nyelvtörténeti Osztálya és a MorphoLogic közötti együttműködés számos veszélyeztetett kisebb uráli nyelv morfológiájának számítógépes feldolgozására 2001-ben kezdődött a NKFP-5/135/01 számú *Komplex uráli nyelvészeti adatbázis* című projektum keretében. Ezt három OTKA projekt követte, amelyeknek keretében a finnugor nyelvcsalád permi ágához tartozó komi és udmurt nyelvre¹, az északi szamojéd nganaszan nyelvre² és a két obi-ugor nyelv, a manysi és a hanti három nyelvjárására³ készített számítógépes morfológiákat sikerült olyan szintre fejleszteni, hogy azokat a tudományos közösség számára publikálhatónak éreztük.

A morfológiákat, valamint a kipróbálásukhoz használható szövegeket webes felületen keresztül tettük elérhetővé, amelyet a MorphoLogic üzemeltet, és a <http://www.morphologic.hu/urali/index.php> címen érhető el.

¹ OTKA T 048309 *Permi nyelvészeti adatbázisok*

² OTKA K 60807 *A nganaszan nyelv számítógépes morfológiai elemzése*

³ OTKA NF 71707 *Obi-ugor morfológiai elemzők és korpuszok*

2 A morfológiák

A finnugor nyelvek (komi, udmurt, manysi, hanti) elemzésére a MorphoLogic *Humor* elemzőjét használjuk. A nganaszan morfológiát a Xerox *xfst* eszközének felhasználásával készítettük el.

A komi és az udmurt beszélőinek száma a többi itt bemutatott nagyon erősen veszélyeztetett nyelvvel ellentétben viszonylag jelentős, a permi nyelvek számottevő irodalommal és könyvkiadással rendelkeznek, ezért ezekre a nyelvekre olyan elemzőket készítettünk, amelyek a sztenderd cirill helyesírással írott szövegek elemzésére képesek. A kisebb, elsősorban csak beszélt nyelvként élő nyelvekre olyan elemzőket készítettünk, amelyek latin betűs fonologikus átírást használnak. Az utóbbiak esetében komoly problémát jelentett a lejegyzések következtelensége, illetve az, hogy a különböző szövegkiadásokban jelentősen különböző átírásokat használtak. A manysi esetében ugyanazon nyelvjárás (az északi) három korpuszának (ChVog⁴, WT⁵, VNGY⁶) feldolgozásához három különböző elemzőt kellett készítenünk.

A weboldal létrehozását elsősorban az a cél motiválta, hogy ezekkel a nyelvekkel kapcsolatos nyelvi adatokat hozzáférhetővé tegyük a tudományos kutatóközösség minél szélesebb köre számára. Ezért lehetőség szerint glosszákkal együtt jelenítjük meg az elemzéseket, hogy azok a nyelvet nem beszélő kutatók számára is értelmezhetőek legyenek. A közeljövőben befejeződő obi-ugor OTKA projektben már kifejezetten cél volt az elemzők tótárának angol glosszákkal való ellátása is. Az egyértelműsítő felület ily módon egyben szemantikai egyértelműsítésre is használható.

A weboldalon jelenleg az alábbi morfológiák érhetőek el:

nyelv	glosszázás	tótár	tótár	toldaléktár mérete ⁷
		lemma	jelentés	
nganaszan	magyar	4200	4775	310
komi (zürjén)	nincs (orosz glosszák hozzáadását tervezzük)	36000		193
udmurt	magyar	13500	18500	286
északi manysi (WT)	magyar, német, angol	3820	4200	376

Az év végéig az alábbi morfológiák kerülnek még fel az oldalra:

északi (ChVog)	manysi	magyar, német, angol	1250	1530	271
északi (VNGY)	manysi	magyar, német, angol	15600	16500	297

⁴ Kálmán Béla: *Chrestomathia Vogulica*. Tankönyvkiadó, Budapest. (1989)

⁵ Kálmán Béla: *Wogulische Texte mit einem Glossar*. Akadémiai Kiadó, Budapest. (1976)

⁶ Munkácsi Bernát: *Vogul népköltési gyűjtemény*. 1–4. Budapest. (1892–1921)

⁷ morféma, ill. lexikalizált morfémakombináció

szinjai hanti	magyar és angol	2300	2500	138
kazimi hanti	magyar és angol	1750	1950	151

3 A webes felület

Az elemzők esetében a kiválasztott szöveget a megfelelő ablakba másolva a felhasználó megkapja a szövegben szereplő szavak lehetséges morfológiai elemzéseit és az elemzésekben szereplő tömorfémák jelentését. A webes felületen valamennyi nyelven hozzáférhetőek olyan példaszövegek, amelyekkel az elemző kipróbálható. Virtuális billentyűzet segítségével a felhasználó maga is gépelhet be szövegeket.

Uráli morfológiai elemzők és szóalak-generátorok

© 2010, MTA Nyelvtudományi Intézet, MorphoLogic

The screenshot shows the MorphoLogic web interface. At the top, there are two radio buttons: 'Elemzés' (selected) and 'Generálás'. Below them is a text input area containing several lines of text in a non-Latin script, likely Hanti. The text includes words like 'χosa', 'ōls', 'man', 'wāŋi', 'ōls', 'akw-mat-ērtn', 'χottaŋ', 'minne', 'nomtn', 'joχtuwās', 'āmp-niēlam', 'tūp-sup', 'wārs', 'ponal-ŋēr', 'χāp-sup', 'wārs', 'naluw-nariytaste', 'χāpe', 'tūpe', 'wis', 'ta', 'towī', 'ta', 'mimi', 'ti-mos', 'ēryi', 'ponal-ŋēr', 'χāp-supt'em', 'šāw-šaw-šāw', 'āmp-niēlam', 'tūp-supt'em', 'pōl-pol-pōl...'. Below the text input is a keyboard layout for 'Mansi Latin' and a control bar with buttons for 'Elemzés' and 'Generálás'. The 'Elemzés' button is selected.

Az elemzéseket megjelenítő webes felület egyben kézi egyértelműsítő eszközként is szolgál: a többértelmű szavak elemzéseit pop-up ablakban jelennek meg, ha az egeret egy többértelmű szó fölé moztatjuk, ezek közül egérrel választhatunk.

mān	pāwluw	ŋāpat,	saran-pāwəlŋ
mān[N Pro]=mān+[NOM]	pāwəl[N]=pāw+l+w[PxPl1]+[NAG]	ŋapa[N dial_Sy]=ŋapa+[LOC]	saran-pāwəl[N]=saran-pāwəl+n[LAT]
en: we+[NOM]	en: village+[PxPl1]+[NAG]	en: adjacency+[LOC]	en: Saranpaul (large Mansi-Komi village)+[LAT]
de: wir+[NOM]	de: Dorflein+[PxPl1]+[NAG]	de: Nahe+[LOC]	de: Saranpaul, ein großes wogulisch-syrjanisches
hu: mi+[NOM]	hu: falu, falucska+[PxPl1]+[NAG]	hu: közel+[LOC]	hu: Szaranpaul (nagy mansi-komi falu)+[LAT]

[...] wātat,	janiy	ŋapa[N dial_Sy]=ŋapat+[LOC]	ōli,
wāta_pasan-[N dial_Ob]=wāta+[LOC]	janiy[A]=janiy+[NAG]	en: adjacency+[LOC]	ōli[V]=ōli+i[VxPrsSg3]
en: edge (of table)+[LOC]	en: big, large+[NAG]	de: Nahe+[LOC]	en: to be, to exist+[VxPrsSg3]
de: (Tisch)kante+[LOC]	de: groß+[NAG]	hu: közel+[LOC]	[NAG] de: sein+[VxPrsSg3]
hu: (asztal) széle+[LOC]	hu: nagy+[NAG]	ŋapa[N dial_Sy]=ŋapat+[Pi]+[NAG]	hu: van+[VxPrsSg3]

[...] jāŋkolmay	lāwawe.	ŋapa[N dial_Sy]=ŋapat+[Pi]+[NAG]
jāŋkolma[N]=jāŋkolma+γ[TRE]	lāwə[V]=lāw+a+	hu: közel+[Pi]+[NAG]
en: swamp where berries grow+[TRE]	en: to say+[Pa]	ŋapa[N dial_Sy]=ŋapat+[PxSg3]+[NAG]
de: Moor, wo Sumpfbeeren wachsen+[TRE]	de: sagen+[Pa]	en: adjacency+[PxSg3]+[NAG]
hu: bogys mocsár+[TRE]	hu: mond+[Pa]	de: Nahe+[PxSg3]+[NAG]
		hu: közel+[PxSg3]+[NAG]

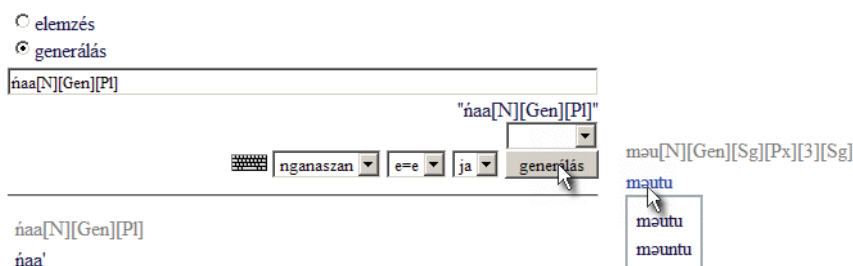
Az elkészült elemzések, illetve azok egyértelműsített változata elmenthető, az elmentett változatot a böngészőbe betöltve, az esetlegesen félbehagyott egyértelműsítő munka később folytatható.



Egyértelműsítés közben javíthatóak az elemzendő szövegben előforduló elgépelések is, ezután a szöveg újraelemeztethető. Ilyenkor nem vesznek el a korábban hozott egyértelműsítési döntések.



A webes felületen keresztül nemcsak morfológiai elemzők, hanem szóalak-generátorok is elérhetők az egyes nyelvekhez. Ha egy adott morfémásorozat több formában is megjelenhet, akkor a generátor kimenete az elemző többértelmű kimenetének megjelenítéséhez hasonlóan jelenik meg a webes felületen, a lehetséges szóalakváltozatok itt is az egérmutatót a generált szóalak fölé mozgatva megjelenő pop-up ablakban láthatóak.



Terveink között szerepel, hogy a weblapot egyértelműsített korpuszokkal és korpuszkereső szolgáltatással egészítsük ki.

Bibliográfia

1. Beesley, K.R., Karttunen, L.: Finite State Morphology. CSLI Publications. Stanford University, Stanford (2003)
2. Novák A.: Milyen a jó Humor? In: Magyar Számítógépes Nyelvészeti Konferencia (MSZNY 2003). Szegedi Tudományegyetem, Szeged (2003) 138–145
3. Novák, A.: Language resources for Uralic minority languages. In: Proceedings of the SALT MIL Work-shop at LREC-2008: Collaboration: interoperability between people in the creation of language resources for less-resourced languages. Marrakech (2008) 27–32
4. Prószycki, G., Novák, A.: Computational Morphologies for Small Uralic Languages. In: Arppe, A., Carlson, L., Lindén, K., Piitulainen, J., Suominen, M., Vainio, M., Westerland, H., Yli-Jyrä, A. (szerk.): *Inquiries into Words, Constraints and Contexts Festschrift in the Honour of Kimmo Koskenniemi on his 60th Birthday*. Gummerus Printing, Saarijärvi/CSLI Publications, Stanford (2005) 116–125