

A készülő MGTSz adatbázis felépítése

Blaho Sylvia, Sass Bálint & Simon Eszter

MTA Nyelvtudományi Intézet

2010. február 4.

Az előadás vázlata

- 1 A projekt bemutatása
 - A szöveg feldolgozásának szintjei
 - A korpusz felépítése
 - Egységes beviteli formátum
 - Kézi kódolás
 - A kódolási szabályzat
 - A normalizálás alapelvei
- 2 A szövegek egyszerűsített átirata
- 3 „Régi magyar konkordancia”



Outline

- 1 A projekt bemutatása
 - A szöveg feldolgozásának szintjei
 - A korpusz felépítése
 - Egységes beviteli formátum
 - Kézi kódolás
 - A kódolási szabályzat
 - A normalizálás alapelvei
- 2 A szövegek egyszerűsített átirata
- 3 „Régi magyar konkordancia”



A projekt

A projekt:

Magyar Generatív Történeti Szintaxis (MGTSz)

OTKA projekt

É. Kiss Katalin vezetésével

2009.04.01.–2013.03.31.



A projekt célja

elektronikus nyelvtörténeti adatbázis:

a teljes ómagyar és válogatott középmagyar anyag

- i. összegyűjtjük és egységesítjük a már meglévő elektronikus nyelvtörténeti anyagokat
- ii. a számítógép által olvasható és feldolgozható formára hozzuk az elektronikusan nem elérhetőeket
- iii. a betűhű változat mellett előállítunk egy egyszerűsített változatot is
- iv. normalizáljuk a szövegeket
- v. a korpuszt morfológiailag elemezzük és egyértelműsítjük
- vi. a korpusz egy részét szintaktikailag is elemezzük.



A szöveg feldolgozásának szintjei

Az eredeti kódextől a morfológiailag elemzett elektronikusan tárolt szövegig a következő szintek vannak:

- (0) fac simile → feldolgozás, kiadás
- (1) betűhű átirat → szkennelés, OCR
- (2) OCR-ezett szkennelés → javítás, kódolás
- (3a) betűhű szöveges elektronikus forma (txt) → egyszerűsítési szabályok



A szöveg feldolgozásának szintjei

- (3b) egyszerűsített változat → normalizálás
- (4) normalizált alak → automatikus morfológiai elemzés
- (5) elemzett forma → automatikus morfológiai egyértelműsítés
- (6) egyértelműsített forma



A korpusz felépítése

A korpuszban minden egyes szövegszó mellett szerepelni fognak a következő adatok:

- betűhú forma (3a): *adÿad*
- egyszerűsített alak (3b): *adyad*
- normalizált alak (4): *adjad*
- szótő (6) alapján: *ad*
- morfológiai elemzés (6) alapján: *ad[V.Sub.S2.Def]*



Lekérdezés minden szinten

a lekérdező lényege, hogy *bármely* szinten meg lehet fogalmazni a lekérdezésünket

Példa

„Milyen szavak szerepelnek egy igealak és egy igekötő között?” :
(6)

gyakorisági lista a korpusz egy részéből: a szótöveken (6) alapján
az „m̄”-et tartalmazó szavak: (3a)



Kézi kódolás

A kódolók munkájának lényege a (3a) és a (4) alak előállítás.

egységes beviteli formátum: a *kódolási szabályzat* alapján



A kódolási szabályzat

- lókuszelölők
- Prószéky-kódolás
- mondatrabontás
- a szkriptor javításainak kódolása
- a normalizálás alapelvei



A normalizálás alapelvei 1.

A ma nem létező összes szót, toldalékot, morfológiai konstrukciót meg kell tartani, nem szabad, hogy ilyen információ elvesszen.

(3a)	(4)	értelmezés
villamik	villamik	villámlik/villanik
isa	isa	bizony
iesek	jeszek	jövök



A normalizálás alapelvei 2.

El kell hagyni az összes fonológiai és helyesírási esetlegességet, egységes, amennyire lehet, a mainak megfelelő helyesírásra kell törekedni.

(3a)

(4)

me8den

minden

menden

minden

minden

minden

algyu

ágyú

agyu

ágyú

srumlast

ostromlást



Outline

- 1 A projekt bemutatása
 - A szöveg feldolgozásának szintjei
 - A korpusz felépítése
 - Egységes beviteli formátum
 - Kézi kódolás
 - A kódolási szabályzat
 - A normalizálás alapelvei
- 2 A szövegek egyszerűsített átirata
- 3 „Régi magyar konkordancia”



„Betűhű” változat

- a betűhű változat elkészítésekor nem a szövegek kézzel írott változatát, hanem az általunk használt átírat szerkesztőjének konvencióit követtük.
- ezért előfordul, hogy a különböző forrás alapján feldolgozott szövegek más és más egyszerűsítéseket tartalmaznak az eredetihez képest.

Példa

- *Jókai kódex: 3 vs. ̄3 vs. 3̄ → a szerkesztő mindhármát 3-ként írja át*



„Betűhű” változat

- a betűhű változat elkészítésekor nem a szövegek kézzel írott változatát, hanem az általunk használt átírat szerkesztőjének konvencióit követtük.
- ezért előfordul, hogy a különböző forrás alapján feldolgozott szövegek más és más egyszerűsítéseket tartalmaznak az eredetihez képest.

Példa

- *Jókai kódex: ʒ vs. ʒ vs. ʒ → a szerkesztő mindhármat ʒ-ként írja át*
- *s vs. ſ, z vs. ʒ, y vs. ý vs. ÿ*
- *palatalizált mássalhangzók: t̃ vs. t̃*



„Betűhű” változat

- a betűhű változat elkészítésekor nem a szövegek kézzel írott változatát, hanem az általunk használt átírat szerkesztőjének konvencióit követtük.
- ezért előfordul, hogy a különböző forrás alapján feldolgozott szövegek más és más egyszerűsítéseket tartalmaznak az eredetihez képest.

Példa

- *Jókai kódex: ʒ vs. ʒ vs. ʒ → a szerkesztő mindhármát ʒ-ként írja át*
- *s vs. ſ, z vs. ʒ, y vs. ý vs. ÿ*
- *palatalizált mássalhangzók: t̃ vs. t̃*



Egyszerűsített átirat

- célja, hogy
 - az ómagyar szövegek olvasásában nem járatos felhasználó számára megkönnyítse az olvasást,
 - minimalizálja a speciális karakterek használatából adódó problémák előfordulásának esélyét
- elkészítésekor a következő szempontokat tartottuk szem előtt:
 - nyelvészeti relevancia
 - gyakrabban használt karakterek
 - egységesítés



Nyelvészeti relevancia

Az egyszerűsített változatból kimaradnak olyan speciális karakterek és diakritikumok, amelyek paleográfiai, stb. jelentőséggel bírnak, nyelvészetivel azonban nem.

Példa

\int → s

3 → z

ÿ → y

ÿ → y



Gyakrabban használt karakterek

A következő szempont az volt, hogy az egyszerűsített átíratban minél kevesebb, a magyar betűkészletben nem meglévő karakter és mellékjel szerepeljen.

Példa

í → ty

ő → ö



Gyakrabban használt karakterek

Ahol mégis a magyarban nem használt betűre volt szükség, előnyben részesítettük az ismertebb, európai nyelvekben gyakran használt karaktereket.

Példa

ē → ě



Gyakrabban használt karakterek

Ahol mégis a magyarban nem használt betűre volt szükség, előnyben részesítettük az ismertebb, európai nyelvekben gyakran használt karaktereket.

Példa

ē → ě

Ez egyrészt az olvasás és a régi magyar példák számítógépes bevitelének megkönnyítését szolgálja (pl. cikkekben, prezentációk segédanyagaiban), másrészt csökkenti a különböző programok fonthasználatából adódó hibák esélyét.



Egységesítés

Azokat a nyelvészetileg releváns karaktereket és diakritikumokat, amelyek különböző nyelvelmékekben eltérő módon jelölik ugyanazt a betűt, illetve tulajdonságot, igyekeztünk egységesíteni.

Példa

\acute{t} → ty

\ddot{t} → ty

ty → ty



Inkonzisztencia

- egyes kódexek között:

Példa

$\dot{g} \rightarrow [g]$ vs. $g \rightarrow [g]$
 $g \rightarrow [t]$ $\dot{g} \rightarrow [t]$

- ugyanazon a kódexen belül:

Példa

Müncheni kódex:

	<i>tő</i>	<i>g</i>	<i>ḡ</i>
<i>gyermek</i>	22	0	
<i>hogy</i>	670	0	
<i>míg</i>	5	3	
<i>ország</i>	137	51	



Outline

- 1 A projekt bemutatása
 - A szöveg feldolgozásának szintjei
 - A korpusz felépítése
 - Egységes beviteli formátum
 - Kézi kódolás
 - A kódolási szabályzat
 - A normalizálás alapelvei
- 2 A szövegek egyszerűsített átirata
- 3 „Régi magyar konkordancia”



Korpusz

- Formátum: egységes, szabványos, hordozható
 - *XML*: adatbázisszerkezet
egymásba ágyazott objektumok + attribtumok
~ mondat és szó + a megfelelő szintek adatai szavanként
 - *UTF-8*: különleges karakterek kódolására
- Anyag
 - *Cél*: teljes ómagyar kori anyag
 - *Jelenleg*:

	(3a)	(4)	(6)
kisebb nyelvemlékek	+	+	
Bécsi kódex	+		
Jókai kódex	+		+



Lekérdezőfelület

Régi magyar konkordancia

Adjon meg egy lekérdezést [\(Guide\)](#)

... vagy adja meg a keresett szó alábbi tulajdonságait!

Megjegyzés:

v0.23 – 2010.02.03. – [S. E.](#) | [Emlékező](#)

Betűhő alak (3a) (teljes):

Normalizált alak (4) (teljes):

Szótő (5) (teljes):

Elemzés (6) (teljes):

Értelmezés (teljes):

Igekötő (teljes):

Megjegyzés (teljes):

Prezentáció:

Eredmény – konkordancia

Régi magyar konkordancia

Adjon meg egy lekérdezést [\(Guide\)](#)

... vagy adja meg a keresett szó alábbi tulajdonságait!

[W FOCUS w_4 ~ 'A4\(\{jonh*}]

Megjegyzés:

Mehet

Törölés

v0.23 – 2010.02.03. – [S. B.](#) | [Erdős](#)

Betűhű alak (3a)	(teljes):	<input type="text"/>
Normalizált alak (4)	eleje:	<input type="text" value="jonh"/>
Szótő (6)	(teljes):	<input type="text"/>
Elemzés (6)	(teljes):	<input type="text"/>
Értelmezés	(teljes):	<input type="text"/>
Igekötő	(teljes):	<input type="text"/>
Megjegyzés	(teljes):	<input type="text"/>

OK

Prezentáció:

2010-02-04 02:56:18

Lekérdezés: [W FOCUS w_4 ~ 'A4\(\{jonh*}]

Találati szavak száma: 6 – Futási idő: 3s

[1] MS - 103a5 - B E43D

eo és minden minden ereinek erősnek dlian olyan lezen leszen ionha jonha, (szíve) mit mint paurek pávának.

[2] OMS - 5 - B E471

en én iunhum jonhom (szívem) ^{DIFFANA} buq búval farad / fárad,

[3] OMS - 10 - B E47B

en én io-hum jonhom (szívem) ^{DIFFANA} delothya alélátja. (alélása) ^{MORFO{noun}}

[4] OMS - 18 - B E204

en én iunhumnok jonhomnak (szívemnek) ^{DIFFANA} bel bel (beiső) tua búja, qui ki (ami) lumha soha nym nem kyul ^{FAIL} tyul: hül.

[5] SzV - 85 - B E285

Nagy nagy bws bús yonhai jonhhal (szívet) ^{DIFFANA} meg vy33a teenek meg-visszatérének, meg-vissza

Eredmény – gyakorisági lista

Régi magyar konkordancia

Adjon meg egy lekérdezést [\(Guide\)](#)

... vagy adja meg a keresett szó alábbi tulajdonságait!

Megjegyzés:

v0.23 – 2010.02.03. – [S. B.](#) | [Erdős](#) Betűhő alak (3a) (teljes): Normalizált alak (4) (teljes): Szótő (6) (teljes): Elemzés (6) (teljes): Értelmezés (teljes): Igeköté (teljes): Megjegyzés (teljes):

Prezentáció:

2010-02-04 02:56:42

Lekérdezés: [W FOCUS w_4 ~ 'A4\(\(nem\)\)4*']

Találati szavak száma: 36 – Futási idő: 3s

nem	nem	22 db
né	nem	3 db
num	nem	3 db
Nem	nem	3 db
nym	nem	1 db
Num	nem	1 db
(n)em	nem	1 db
RECO MORFO{tagadószó}		
(n)em	nem	1 db
Nem==	nem	1 db
RECO		
(n)im==	nem	1 db

Példák

(1/4)

- névutók lekérdezése (Hegedűs Vera)

[W FOCUS w_6e ~ 'Pp']

→ 340 db, ebből ragozott 100 db (mind E/3)

- ragozott főnévi igenevek lekérdezése (Tóth Ildikó)

[W FOCUS w_6e ~ 'Inf.Px']

→ 132 db, eloszlás:

	1	2	3
egyes szám	23	23	70
többes szám	5	6	5



Példák

(2/4)

„Hol nincs ott a névelő, pedig várnánk?” (Egedi Barbara)

Olyat keresünk, ami *nincs* ott.

→ Adjunk meg konkrét ilyen helyzeteket.

Lekérdezés: definit ige után tárgyesetű főnév

[W FOCUS w_6e ~ 'V.*Def']

[W FOCUS w_6e ~ 'N.*Acc']

Egy találat

„Es azért ewkewztewk zent ferencz
czudalatost **gyczerÿuala teremtwtt**”



Példák

(3/4)

„Míg a mai magyarban a tagadott igekötős ige (egy-két szerkezetet kivéve) fordított szórendű, a korai ómagyar korban az igekötő az esetek nagyobb részében megelőzi a tagadott igét.” (É. Kiss Katalin)

```
[W FOCUS w_6e ~ 'Mod']
```

```
[W FOCUS w_6e ~ 'V\.'
```

```
[W FOCUS w_6e ~ 'Vpfx']
```

```
[W FOCUS w_6e ~ 'Vpfx']
```

```
[W FOCUS w_6e ~ 'Mod']
```

```
[W FOCUS w_6e ~ 'V\.'
```

Egy találat

„Ver touaba **ký** nem futott”

Példák

(4/4)

„Míg a mai magyarban a tagadószó hordozza a tagadást, és a se-névmások csupán a tagadószóval egyeztetett alakok, a korai ómagyar korban a se-névmásoknak is lehetett tagadó erejük.” (É. Kiss Katalin)

Lekérdezés: 'senki/semmi' után tagadószótól különböző szó

```
[W FOCUS w_6s ~ 'ˆ6s\(\(se[nm][km]i\)\)\$']
```

```
[W FOCUS NOT(w_6e ~ 'ˆ6e\(\(Mod\)\)\$')]
```

Egy találat

„mendenestewlfoguan maganac **semýtt meg** tarttuan”



Elérhetőség

A lekérdezőfelület szabadon elérhető:

<http://corpus.nytud.hu/rmk>

Köszönjük a figyelmet!

