

# Az Ómagyar Korpusz bemutatása

Simon Eszter

MTA Nyelvtudományi Intézet

2017. január 13.  
29. Finnugor Szeminárium

# Az előadás vázlata

- 1 A projektek
- 2 A korpusz anyaga
- 3 A feldolgozás lépései
- 4 Korpuszlekérdező felületek

# A projektek

## Magyar **G**eneratív **T**örténeti **Sz**intaxis (MGTSz) 1&2

az OTKA támogatásával  
(OTKA No. 78074 & OTKA No. 112057)  
2009. április – 2013. március & 2015. január – 2018. december  
MTA Nyelvtudományi Intézet  
projektvezető: É. Kiss Katalin

# A projektek célja

## MGTSz 1&2

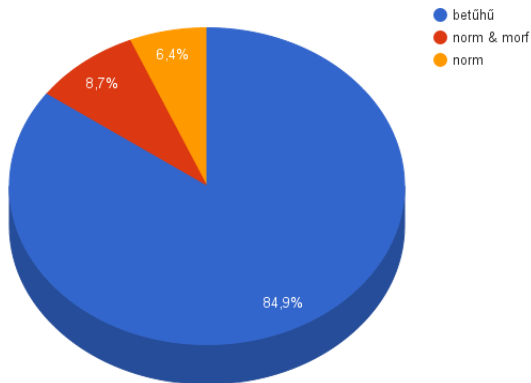
- 1 elméleti cél: szintaktikai változások rekonstrukciója és vizsgálata
- 2 számítógépes munka: egy olyan annotált történeti korpusz létrehozása, amely tartalmazza az összes egybefüggő ómagyar szöveget & a korpusz kibővítése középmagyar korból származó Újszövetség-fordításokkal

# A korpusz anyaga

- 47 ómagyar kódex
- 24 kisebb egybefüggő szövegemlék az ómagyar korból (HB – ME)
- 244 misszilis (Hegedűs–Papp: Középkori leveleink (1541-ig))
  
- Pesti Gábor (1536)
- Sylvester János (1541)
- Heltai Gáspár (1565)
- Károli Gáspár (1590)
- Káldi György (1626)

# A korpusz mérete

2.772.788 token, ebből 419.067 normalizált,  
ebből 242.260 morfológiailag elemzett és egyértelműsített



# Szövegfeldolgozottsági szintek

- 
- (1) kiadott kódex szkennelve  
→ OCR
  - (2) nyers OCR-kimenet  
→ *kézi* javítás
  - (3) betűhű elektronikus forma  
→ *kézi* normalizálás
  - (4) normalizált forma  
→ *automatikus* morfológiai elemzés
  - (5) szótövesített és morfológiailag elemzett forma  
→ *félautomatikus* egyértelműsítés
  - (6) egyértelműsített korpusz
-

# Az annotációs szintek

az egyes szövegszavakhoz tartozó annotációs szintek párhuzamosan alakulnak a szövegfeldolgozottsági szintekkel:

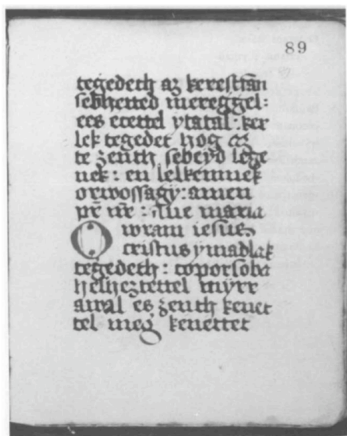
- betűhú forma (3): *ad̂ad*
- normalizált alak (4): *adjad*
- szótő (6) alapján: *ad*
- morfológiai elemzés (6): *V.Sub.S2.Def*



# A korpusz felépítése

oldal	sor	betűhű	norm	lemma	elemzés
001	01	Mÿ	mi	mi	Pro.Nom_gen
001	01	vronknac	Urunknak	Úr	N:P.PxP1.Dat_gen
001	01	iesus	Jézus	Jézus	N:P
001	01	cristusnac	Krisztusnak	Krisztus	N:P.Dat_gen
001	01	gyczeretyre	dicséretére	dicséret	N.PxS3=i.Sub
001	02	es	és	és	C
001	02	gyczewsegere	dicsőségére	dicsőség	N.PxS3.Sub

# 1. Kézzel írott kódexek, nyomtatott kiadások



177  
89r

- tegedeth az keresztian  
sebhetted mereggel :  
ees ecettel ytatal : ker-  
lek tegedet hog az  
-lek tegedet hog az  
5 te zenth sebeyd legé-  
-nek : en lelkemnek  
orwossagy : amen  
pf nr : Aue maria  
O wram iesus  
10 cristus ymadlak  
tegedeth : coporsoba  
helhezttel myrr-  
-awal es zenth kenet-  
-tel meg kénettet

## 2. Szkenelés, OCR

tegedeth az kerestfan  
 sebhethed mereggél :  
 ees écttel ytatal : ker-  
 -lek tegedet hog az  
 5 te zenth sebeyd legé-  
 -nek : en lolkemnek  
 orwossagý : amen  
 pf nf : Aue maria  
 O wram iesus  
 10 cristus ýmadlak  
 tegedeth : coporsoba  
 helhezttel myrr-  
 -awal es zenth kenet-  
 -tel meg kénéttet

177  
 89r

tegedeth az kerestfan  
 sebhethed méreggel :  
 ees ecettel ytatal : ker-  
 -lek tégedet hog az  
 te zenth sebeyd legé-  
 -nek : en lolkemnek  
 orwossagy : ámen  
 pf nf : Aue maria  
 O wram iesus  
 eristus ymadlak  
 tegedeth : coporsoba  
 hellieztettel myrr-  
 -awal es zenth kenet-  
 -tel meg kenéttet

### 3. A betűhű szöveg előállítása

52 latin alapkarakter

42 diakritikus jel

10 szám

15 speciális karakter

34 szövegtagoló és egyéb jel

3 görög betű

Összesen: 156 karakter + ezek kombinációi

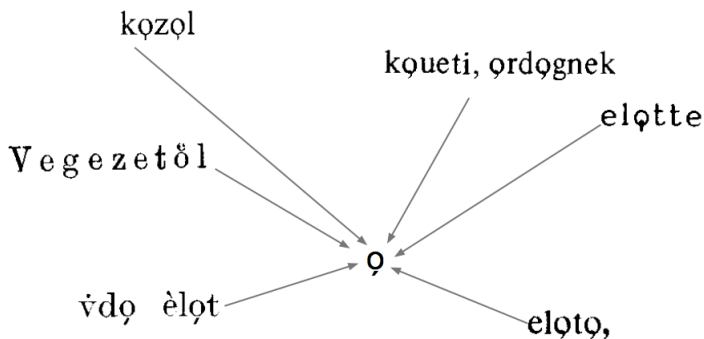
*UTF-8 kódolású sztenderd Unicode karakterek*

# Unicode

- nemzetközi szabvány
- a világ összes nyelvének összes karakterét egy kódolási rendszerbe foglalja
- minden platformon elérhető
- lehetőséget nyújt az alapkarakterek és a diakritikus jelek kombinációjára, pl.  $e + ' + \bar{=} \tilde{e}$

az egész korpuszra kiterjedő szigorúan egységes formátum, DE:  
van olyan régi karakter, amely a Unicode-ban nincsen reprezentálva  
→ helyettesítő karakter:  $L \rightarrow \check{c}$

# Hangjelölés-egységesítés



## 4. Normalizálás

betűhű	normalizált
a varofba	a városba
ahazi	a házi
annèphèz	a néphez
az arpak	az árpák
ānèp	a nép
a" tew	a tú
a' nyaar	a nyár
a · mendenható	a Mindenható

# A normalizálás első alapelve

az összes ma már nem létező szót, toldalékot, morfológiai konstrukciót megtartjuk

## Példa

<i>ýsa</i>	<i>pur</i>	<i>es</i>	<i>chomuv</i>	<i>uogmuc</i>
<i>isa,</i>	<i>por</i>	<i>és</i>	<i>hamu</i>	<i>vagyunk</i>
<i>lata</i>	<i>ø</i>	<i>napat</i>	<i>fèkette</i>	
<i>látá</i>	<i>ő</i>	<i>napát</i>	<i>fekette</i>	



# A normalizálás második alapelve

elhagyjuk az összes fonológiai és helyesírási esetlegességet, vagyis egységes, a mainak megfelelő helyesírásra törekszünk

## Példa

*meden ~ menden ~ mendun ~ mendē ~ mendè ~ miden ~  
minden ~ mynden ~ mýnden ~ mýndē ~ mýden ~ mýnden ~  
mýndew ~ mýnden ~ mýndon ~ mēden ~ mēdèn ~ mēdē ~  
mēdèn ~ mēnden ~ mēden → minden*

# Tokenizálás és mondatra bontás

## Példa

*de sãbãdicz== ==mk mikët a gonostwl*

*de szabadíts meg minket a gonosztól*

*harmal napon halottay bool felthamata*

*harmad napon halottaiból feltámada*

*egmen-@@denic ȝ at't'afiat nē zorongat't'a*

*egymindenik ő atyjafiát nem szorongatja*

# Automatikus előnormalizálás

- tokenizálás: egy sor–egy token, a kötőjeles szavak összevonása, az írásjelek leválasztása
- mondatra bontás: minden írásjel után üres sor
- normalizálás: a már normalizált szövegekből készült lista alapján, amely tartalmazza az egy az egyhez megfeleltethető betűhű–normalizált párokat

*ezzel a tokenek 63%-át lefedjük, és 49%-nyi időt spórolunk meg a teljesen kézi normalizáláshoz képest*

## 5. Morfológiai elemzés és egyértelműsítés

**morfológiai elemzés:** a Humor ó- és középmagyarra szabott változatával

**egyértelműsítés:** kézzel egy webes interfészen keresztül

- a jó elemzés kiválasztható a lehetséges elemzések listájából
- a normalizált szóalak és az elemzés is szerkeszthető, ha szükséges
- az eredmény egy HTML fájl, amit visszakonvertálunk az általunk használt tsv formátumra

# Univerzális Dependencia és Morfológia

- Humor saját formalizmus → sztenderdizáció → Univerzális Dependencia és Morfológia
- most: a morfológiai elemzés konverziója → később: szintaktikai elemzés hozzáadása
- különbségek az ómagyar és a mai magyar között → különbségek a szófaji és az inflexiós címkékben
- különbségek a két annotációs séma között → bizonyos jelenségeket nem lehetett információvesztés nélkül konvertálni

# Metaadatok

- **lókuszelölők**
  - az adott szó helye az eredeti kódexben
  - bibliai könyv-, fejezet- és versszámítás
- **értelmezés:** a normalizált alak mai magyarra való “fordítása”
- **megjegyzés:** TITLE, LANG{nyelv}, ADD, RECO, STRIKE, FAIL, FRAG
- **igekötő:** az elvált igekötőt jelöljük az alapige mellett

# Lekérdezés minden szinten

- a lekérdező lényege, hogy *bármely szinten* meg lehet fogalmazni a lekérdezésünket
- a lekérdezés eredménye független a lekérdezéstől, vagyis például kereshetünk a normalizált szinten, és megjeleníthetjük a betűhűt
- a lekérdezés eredménye: konkordancia vagy gyakorisági lista
- a kontextus ablaka állítható
- a szövegemlék specifikálható

```
http://omagyarokorpusz.nyttud.hu/hu-search.html
```

# A felhasználói felület

## Magyar Generatív Történeti Szintaxis

Kezdőlap » Keresés
hun | eng

**kezdőlap**

**projektum**

**korpusz**

**keresés**

**szövegmélekek**

**kódexek**

**egyéb szövegek**

**kapcsolat**

### Keresés

Adja meg a keresett szó tulajdonságait ...

Betűhi (b)	:	<input type="text"/>	<input type="text"/>	<input type="text"/>			
Normalizált (n)	:	<input type="text"/>	<input type="text"/>	<input type="text"/>			
Szótő (e)	:	<input type="text"/>	<input type="text"/>	<input type="text"/>			
Elemzés (e)	:	<input type="text"/>	<input type="text"/>	<input type="text"/>			
Értelmezés	:	<input type="text"/>	<input type="text"/>	<input type="text"/>			
Igekötő	:	<input type="text"/>	<input type="text"/>	<input type="text"/>			
Megjegyzés	:	<input type="text"/>	<input type="text"/>	<input type="text"/>			
Azonosító	:	<input type="text"/>	<input type="text"/>	<input type="text"/>			

Formátum:		<input type="text" value="konkordancia"/>			
Megjelenítés:		<input type="text" value="minden"/>			
Kontextus <sup>NETA</sup> :		<input type="text" value="csak a találati szakasz"/>			
Nyelvemlék:		<input type="text" value="mind"/>			

Szerk. mód:		<input type="text" value="új lekérdezés"/>			
Kihagyás:		Min: <input type="text" value="0"/>		Max: <input type="text" value="5"/>	<input type="button" value="OK"/>

... vagy szerkesse a lekérdezést közvetlenül. (Guide)

Megjegyzés:

v0.3.8.9 – 2016.11.24. – S. B. | Emdros – sűgő – a Nemzeti Korpuszportál tagja.  
 A korpusz használata esetén kérjük hivatkozni a következő cikkre: Simon Eszter, Sasó Billian: Nyelvtudományi és kulturális örökség, avagy korpuszépítés ómagyar kódexekkel. In: *Általános Nyelvtudományi Tanulmányok*. 2012(XXIV): 243-264.



# A lekérdezés eredménye: konkordancia

[384] Konyvecse - 27r - 1/113219

Mýnden	<b>földreh</b>	ký meneh	az	o	zongesek	,
minden	<b>földre</b>	kimene	az	ő	zöngésük	,
minden	föld	kimegy	az	ő	zöngés	
N:Pro	N.Sub	VPfx.V.Ipf.S3	Det	N:Pro.S3	N.PxP3	

[385] Konyvecse - 27r - 1/113241

zer-@@zed	oketh	fejedelmpl	mýnden	<b>földön</b>	:
szerzéd	őket	fejedelemül	minden	<b>földön</b>	:
szerez	ők	fejedelem	minden	föld	
V.Ipf.S2.Def	N:Pro.P3.Acc	N.Ess	N:Pro	N.Sup	

[386] FestK - 3 - 1/115123

merth	ew	kezeeben	wadnak	<b>fewldnek</b>	mýnden	wégey	:
mert	ő	kezeében	vannak	<b>földnek</b>	minden	végei	,
mert	ő	kéz	van	föld	minden	vég	
C	N:Pro.S3.Nom_gen	N.PxS3.Ine	V.P3	N.Dat_gen	N:Pro	N.PxS3.Pl	

# A lekérdezés eredménye: gyakorisági lista

Query: [w FOCUS w\_4 ~ 'föld']

Number of hits: 596 – Elapsed time: 25s

földön <b>földön</b>	49 db
földnèc <b>földnek</b>	44 db
föld <b>föld</b>	42 db
földèt <b>földet</b>	29 db
földeből <b>földjéből</b>	18 db
földön <b>földön</b> föld N.Sup	18 db

# Párhuzamos Bibliaolvasó

<http://parallelbible.nytud.hu/>

- bibliai fejezet- vagy vers szinten párhuzamosan lehet megjeleníteni különböző bibliafordításokat
- tartalmaz ó-, közép- és mai magyar bibliafordításokat, plusz a King James Bible-t
- terv: bibliafordítások uráli nyelveken és oroszul is

# Köszönöm a figyelmet!

`simon.eszter@nytud.mta.hu`

`http://omagyarkorpusz.nytud.hu/`

`http://parallelbible.nytud.hu/`