

# A kibővített *Magyar történeti szövegtár* új keresőfelülete

Sass Bálint

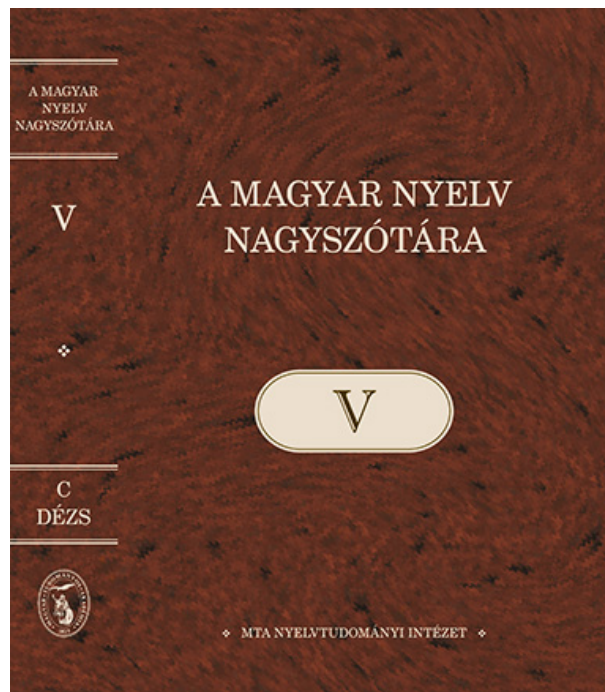
MTA Nyelvtudományi Intézet

sass.balint@nytud.mta.hu

A nyelvtörténeti kutatások újabb eredményei IX.

2016. április 27., Szeged

# Nszt és Mtsz



<http://nszt.nytud.hu/nszt.html>

*Magyar történeti szövegtár*

- 30 millió szövegszó
- 1772-2010-ig ~ 240 év
- 2015-ben készült el a bővítés

1.

## **Az Mtsz használata**

# Az Mtsz felülete

egyszerű keresés: *de viszont*

*Ami látszik:*

- nagybetű/kisbetű nem számít – sőt: f
- strukturális információk (oldal, bekezdés, (vers)sor): **zölddel**
- találatok időrendben

*Ami nem látszik:*

- évszám katt: részletes bibliográfiai adatok
- találat katt: nagyobb kontextus
- *funkciók*: alkorpuszok, mentés, szűrés, gyaklisták, kollokációk
- CQL = Corpus Query Language – formális lekérdezőnyelv

elérhetőség: <http://clara.nytud.hu/mtsz>

# CQL – reguláris kifejezések

Bizonyos tulajdonságú karaktersorozatok megadására.

- . tetszőleges karakter
- \* a megelőző karakterből bármennyi (0 vagy több)
- [ab] 'a' vagy 'b' karakter

*Példák:*

- (1) tejf.l
- (2) .\*
- (3) .\*bb
- (4) nélk[üúű]l

# CQL

## Corpus Query Language

[..] egy tokenre vonatkozó megkötések

$x="y"$   $x$  attribútum értéke legyen  $y$  – Mtsz: csak *word* attribútum van

$x!="y"$   $x$  attribútum értéke *ne* legyen  $y$

& és kapcsolat megkötések között

*Példák:*

1. [] []

2. [word="majd"]

3. "majd"

4. [word!="a.\*"]

# Mtsz példalekérdés

*Feladat. Keressünk ilyen: tárgyestű szó + múltidejű E/3 ige!*

# Mtsz példalekérdés

*Feladat. Keressünk ilyen: tárgyestű szó + múltidejű E/3 ige!*

" . \* t " " . \* t t "



# Mtsz példalekérdés

*Feladat. Keressünk olyet: tárgyesetű szó + múltidejű E/3 ige!*

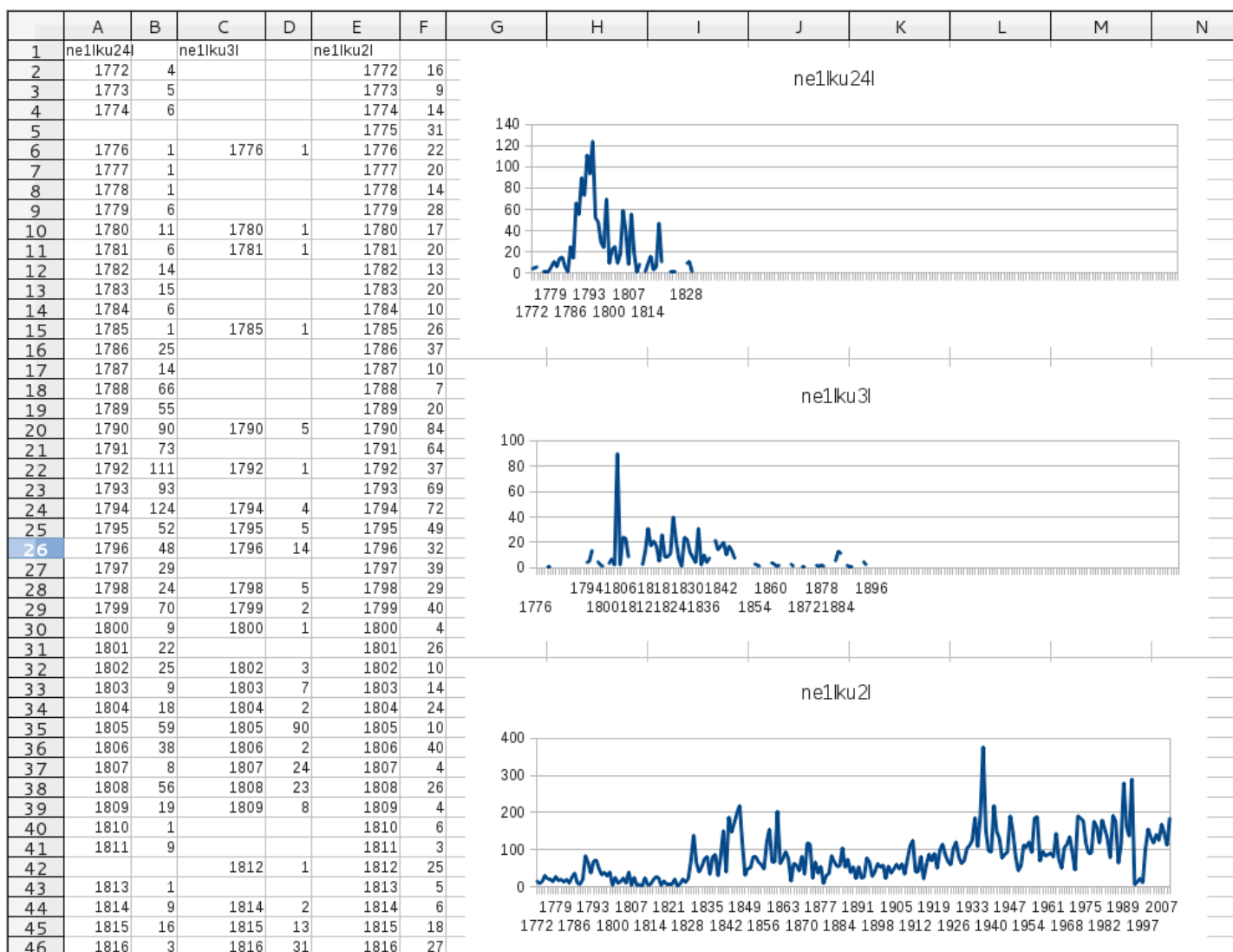
" . \* t " " . \* t t "

*most itt – ???*

# Mtsz példalekérdezés – a vizsgálat lépései

1. CQL: ".\*t" ".\*tt"
  2. Gyakoriságok / szóalakok
  3.  $p \rightarrow$  erőt vett
  4. Milyen szó jön utána?  $\rightarrow$  Gyakoriságok: 1R
  5.  $p \rightarrow$  rajta
  6. Rendezés / jobb: hogy *mi* vesz erőt rajta
- $\rightarrow$  félelem, féltékenység, habozás, kacagás, kishitűség, kíváncsiság ...

# Diakrón vizsgálat: *nélkül* helyesírása



2.

## Az Mtsz építése

# Az Mtsz építése

Mi mindent kell csinálni ahhoz,  
hogy sima szövegből ilyen korpuszlekérdezőfelület legyen?

1. Unicode (UTF-8) karakterkódolás    é ✓    ő ✓    í ✓    ö ✓
2. XML – saját / TEI

*korpuszkezelő rendszer:* NoSketchEngine (NoSkE)

`https://nlp.fi.muni.cz/trac/noske`

3. NoSkE XML ← XSLT

# NoSkE XML

```
<doc mtsz_id="7021030" author="Baróti Szabó Dávid" wdate="1808" ...>
  <oldal oldalszam="18">
    <par>
      Ám
      de
      vízont
      hallá
      <g/>
      '
      hogy
      ...
      <br/>
      ... további tokenek
    </par>
    ... további bekezdések
  </oldal>
  ... további oldalak
</doc>
```

# NoSkE XML

```
<doc mtsz_id="7021030" author="Baróti Szabó Dávid" wdate="1808" ...>
  <oldal oldalszam="18">
    <par>
      Ám
      de
      vifzont
      hallá
      <g/>
      '
      hogy
      ...
      <br/>
      ... további tokenek
    </par>
    ... további bekezdések
  </oldal>
  ... további oldalak
</doc>
```

# Az Mtsz építése

Mi mindent kell csinálni ahhoz,  
hogy sima szövegből ilyen korpuszlekérdezőfelület legyen?

1. Unicode (UTF-8) karakterkódolás      é ✓ ő ✓ í ✓ ö ✓
2. XML – saját / TEI

*korpuszkezelő rendszer:* NoSketchEngine (NoSkE)

`https://nlp.fi.muni.cz/trac/noske`

3. NoSkE XML ← XSLT
4. tokenizálás (...és egyéb elemző lépések)
5. „indexálás” = a korpusz betöltése a korpuszkezelőbe



# Nemzeti Korpuszportál (NKP)

Együtt, egy helyen minden meglévő...

- magyar nyelvű, online lekérdezhető korpusz
- korpuszlekérdező funkció

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

*Cél:* a korpuszok népszerűsítése a szakma és a nagyközönség felé

*Távlati cél:* automatizálás

# A kibővített *Magyar történeti szövegtár* új keresőfelülete

Sass Bálint

MTA Nyelvtudományi Intézet

sass.balint@nytud.mta.hu

A nyelvtörténeti kutatások újabb eredményei IX.

2016. április 27., Szeged