



Nyelvelemzés sajátkezűleg a magyar INTEX rendszer

Váradi Tamás
varadi@nytud.hu

Vázlat

- A történet eddig
- Az INTEX rendszer
- A magyar modul
- Az INTEX korlátai
- premier előtt: NooJ
- konklúziók, további teendők

Rövid történeti háttér

- LADL – Maurice Gross
 - lexikai nyelvtan, véges állapotú technológia
- 1993 INTEX **Max Silberztein**
- elterjedtség
 - több mint 20 kutatóintézetben alkalmazzák
 - francia, angol, szerb, portugál, spanyol
 - éves konferenciák
- Magyar?
 - Premier – itt és most!

Integrált nyelvelemző környezet

- Három eszköz egyben
 - korpuszelemző
 - nyelvtanfejlesztő
 - oktató
- Egységes (véges állapotú) technológia
- Robosztus, jól kidolgozott lexikon
- morfoszintaktikai, szemantikai jegyek
- Grafikus felület

INTEX a korpuszkezelő eszköz

- **.txt** szövegfile → azonnali lekérdezés
- keresés
 - reguláris kifejezések
 - gráfok (nyelvtanok)
- kimenet: konkordancia

A szöveg betöltése

The screenshot shows the 'Intex - Current Language Is Hungarian' application window. The main window title is 'Text: C:\Intex\Hungarian\Corpus\news.txt'. The text area contains several paragraphs of Hungarian news text. A 'Warning' dialog box is overlaid on the text, asking 'Do you want to pre-process the text?' with 'Igen' (Yes) and 'Nem' (No) buttons. A callout bubble points to the dialog box with the text 'Kérünk előfeldolgozást?' (We request pre-processing?).

Text units are lines/paragraphs.

Folytatódik az IRB auditálása
Hétfőn folytatódik az Investicná a Rozvojoová Banka (IRB) auditja, az OTP Bank szakemberei és tanácsadói ismét felkeresik a bank adatszobáját - közölte Wolf László, az OTP Bank vezérigazgató-helyettese az MTI érdeklődésére.

Befejeződött az IRB auditálása
Befejeződött az IRB auditálása - tájékoztatták a pénzügyetnél kedden az MTI-t.

Az Ernst and Young lehet az MNB könyvvizsgálója
Az Ernst and Young Könyvvizsgáló Kft-t jelölték az Állami Számvevőszék (ÁSZ) elnöke - tájékoztatták kedden az MTI-t.

Sport - a Manchester United maradt a leggazdagabb csapat
Immár negyedik éve a brit Manchester United a leggazdagabb csapat a világon - közölte pénteken a Deloitte társvezetője az MTI-vel.

A UAL is elpártolt az Andersentől
2002.05.03. 14:06

Warning
Do you want to pre-process the text?
Igen Nem

Kérünk előfeldolgozást?

egyszerű szöveglekérdezés

The screenshot displays the Intex software interface with the following components:

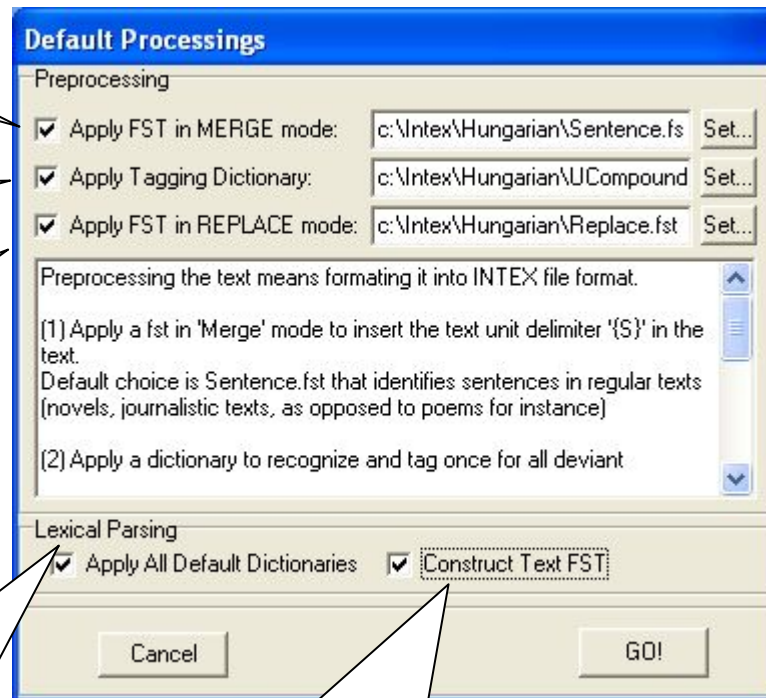
- Window Title:** Intex - Current Language Is Hungarian
- Menu Bar:** Text, DELA, FSGraph, Lexicon-Grammar, Edit, Windows, Info
- Text File:** C:\Intex\Hungarian\Corpus\news.snt
- Concordance File:** C:\Intex\Hungarian\Corpus\news_snt\concord.rtf
- Status Bar:** Search complete: 200 sequences were indexed
- Main Text Area:** A concordance list for the word "szerint". Each entry consists of a snippet of text from a news article, followed by the word "szerint" (underlined), and then the full sentence containing the word. The text is in Hungarian and discusses various news items from 2002.
- Dialog Box: Display indexed sequences...**
 - Extract:** Matching Text Units (selected), Unmatching Units
 - Show Matching Sequences in Context:** (checked)
 - Lengths of Contexts:** Left Col: 40 chars., Right Col: 55 chars.
 - Sort According To:** Center, Right Col.
 - Buttons:** Build concordance

A szöveg előfeldolgozása

Mondatokra bontás
{S} címke beillesztésével

Többszavas
kifejezések bejelölése

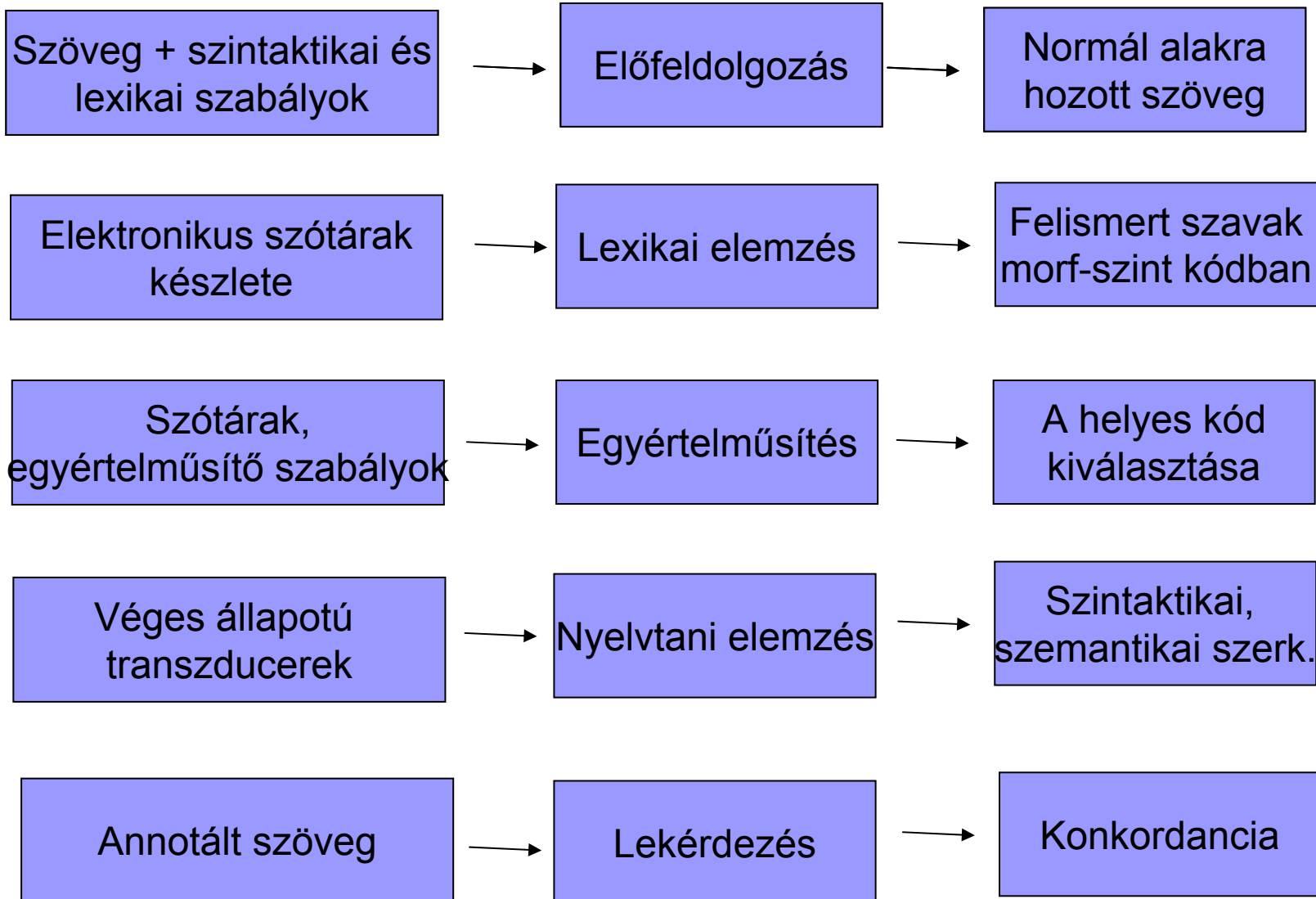
Nem „szabványos” alakok
standard alakra konvertálása
(aszem-> azt hiszem)



Az alapszótárak lefuttatása

Mondatok véges állapotú
automatákba szervezése

A korpusz teljes feldolgozása



Lexikai erőforrások

■ szótárak

- véges számú, lexikalizált elemek
- morfológiai, szintaktikai, szemantikai jegyek

■ lexikai szabályok (lexikai transzducerek)

- hagyományos szótárak által ignorált elemek
 - nyitott lexikai osztályok (pl. képzett elemek)
 - számok, dátumok, nevek (nyílt tokenosztályok)

Az INTEX szótárak

- központi szerep
- robusztus, gyors,
- nyitott, bővíthető
- 4 féle szótár
 - DELAF egyelemű szóalakok
 - DELACF többszavas kifejezések(szóalakok)
 - DELAE idiomatikus kifejezések
 - DELAS morfológiai generáló szótár

A DELAF szótár

Szóalak

Szótári alak
(lemma)

adásvételtől,adásvétel.N+abstract:m
addig,az.PRO:g
addigi,addigi.A:0
addiginál,addigi.A:1
adhat,adhat.V:e3
adhatják,adhat.V:Tt3

Szófaj

morfológiai
kód

Szintaktikai,
szemantikai jegy

A DELACF szótár

Egye
ezze
talá
Belü

```
Bayer Corp,.N+proper:0  
Bayerische Hypo,.N+proper:0  
Beck György,.N+proper:0  
Befektetési Rt,.N+proper:0  
Bellis Rt,.N+proper:0  
Bencze József,.N+proper:0  
Beogradska Banka,.N+proper:0  
Bernard Ebbers,.N+proper:0  
Berva Rt,.N+proper:0  
Best Buy,.N+proper:0  
Best Buy,.N+proper:0  
Beton&t Rt,.N+proper:0  
Betonút Rt,.N+proper:0  
Betonútépítő Rt,.N+proper:0  
Bicinium Kft,.N+proper:0  
Bild Zeitung,.N+proper:0  
Bill Gates,.N+proper:0  
Binder Rt,.N+proper:0  
Biodiszkont Rt,.N+proper:0  
Biologicals Kft,.N+proper:0  
Body Shop,.N+proper:0  
Boeing Co,.N+proper:0
```

A lexikai kódolás

- Az alapszótár toldalékolt szóalakok listája
- Több százezer szavas szótárak
- Tetszőleges információ kódolható szint./szem. jegy formájában
- Morfológiai jegy kódja csak egy betű lehet
- Célszerű igazodni a kialakult gyakorlathoz
- A szótárak láncba szervezhetők és ki-be kapcsolhatók

Text: C:\Intex\Hungarian\Corpus\news.snt

Display tags

17948 delimited units, 35

27430 simple words, 158

C:\Intex\Hungarian\Corpus\news.snt

{S}Folytatódik a

{S}Hétfőn folyta

szakemberei és t

OTP Bank vezérig

{S}Befejeződött

{S}Befejeződött

{S}Az Ernst and

{S}Az Ernst and

könyvvizsgálóján

hétfőn közlemény

{S}Sport - a Man

{S}Immár negyedi

csapata bevétel

{S}MTI is elné

rtverek esetében

héten aláírta az

plazmatévék gyár

l a társaság rep

esellschaft pénz

kban a Deutsche

ereze a tavalyin

lfoglalva. A Fi

Text Vocabulary stored in C:\Intex\Hungarian\Corpus\news_snt

Edit DLF: 27520 simple-word lexical entries

Applied lexical resources:

-- For Simple Words:
newsml_szotar.bin

Update

Edit DLM: 0 analyzed tokens

Edit ERR: 158 unknown simple words

ajánlára apadlóra azFTSE azMTI áfá
án ától ával ben Bizotság bólból
csüörtöki csütöörtökön dolárral
dollár Dupljára együttműködés
ejgzett elleni előtt először
eredményről Ericssonnak es

Edit Frozen expressions:

Edit Compounds:

Tokens...

FST-Text

Lekérdezés képlettel

<van> = lemma,
illeszkedik a van
összes ragozott
alakjára

+ = a logikai „vagy”
művelet jele

a zárójelben
szereplő elemek
közül az egyik

Szóköz = konkatenáció

Locate pattern in the form of:

- Regular expression:
- Graph:
- Tags & Recognized Simple Words
- Recognized Compounds
- Recognized Frozen Expressions

Index

- Shortest matches
- Longest matches
- All matches

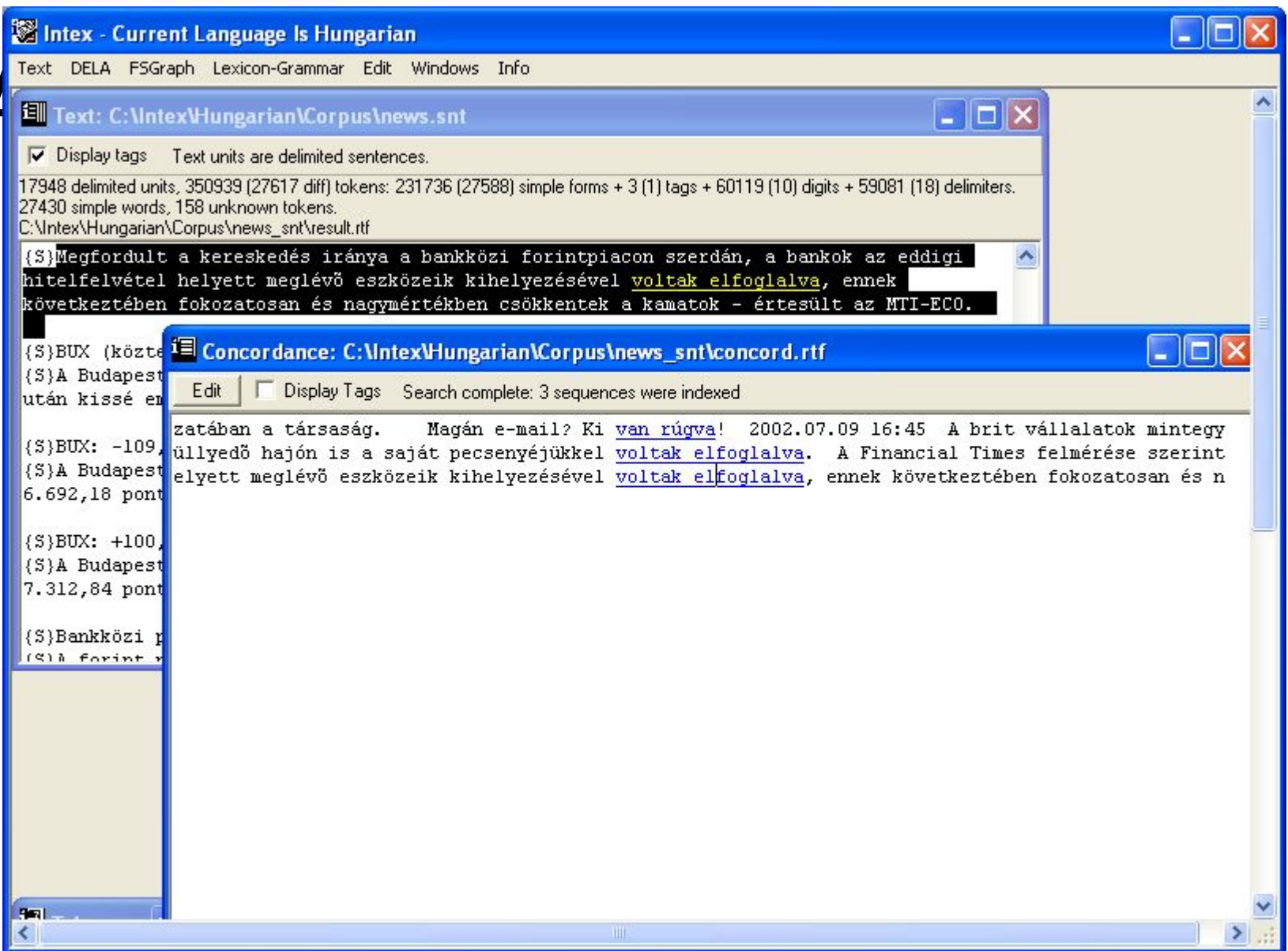
FST outputs

- Are not taken into account
- Merge with input text
- Replace recognized sequences

Search limitation

- Stop after matches
- Index all utterances in text
- 1 example per match

Use Text Index



Intex - Current Language Is Hungarian

Text DELA FSGraph Lexicon-Grammar Edit Windows Info

Text: C:\Intex\Hungarian\Corpus\news.snt

Display tags Text units are delimited sentences.

17948 delimited units, 350939 (27617 diff) tokens: 231736 (27588) simple forms + 3 (1) tags + 60119 (10) digits + 59081 (18) delimiters.
27430 simple words, 158 unknown tokens.

C:\Intex\Hungarian\Corpus\news_snt\result.rtf

{S}Megfordult a kereskedés iránya a bankközi forintpiacon szerdán, a bankok az eddigi hitelfelvétel helyett meglévő eszközeik kihelyezésével voltak elfoglalva, ennek következtében fokozatosan és nagymértékben csökkentek a kamatok - értesült az MTI-ECO.

Concordance: C:\Intex\Hungarian\Corpus\news_snt\concord.rtf

Edit Display Tags Search complete: 3 sequences were indexed

{S}BUX (közte...
{S}A Budapest...
után kissé en...
zatában a társaság. Magán e-mail? Ki van rúgva! 2002.07.09 16:45 A brit vállalatok mintegy
{S}BUX: -109...
{S}A Budapest...
6.692,18 pont...
{S}BUX: +100...
{S}A Budapest...
7.312,84 pont...
{S}Bankközi p...
{S}A forint v...

Lekérdezés gráffal

Kiinduló pont
a vizsgált szöveg
tetszőleges pontja

itt vagyunk, ha
a szövegben
szerepel a **van**
egy ragozott
alakja

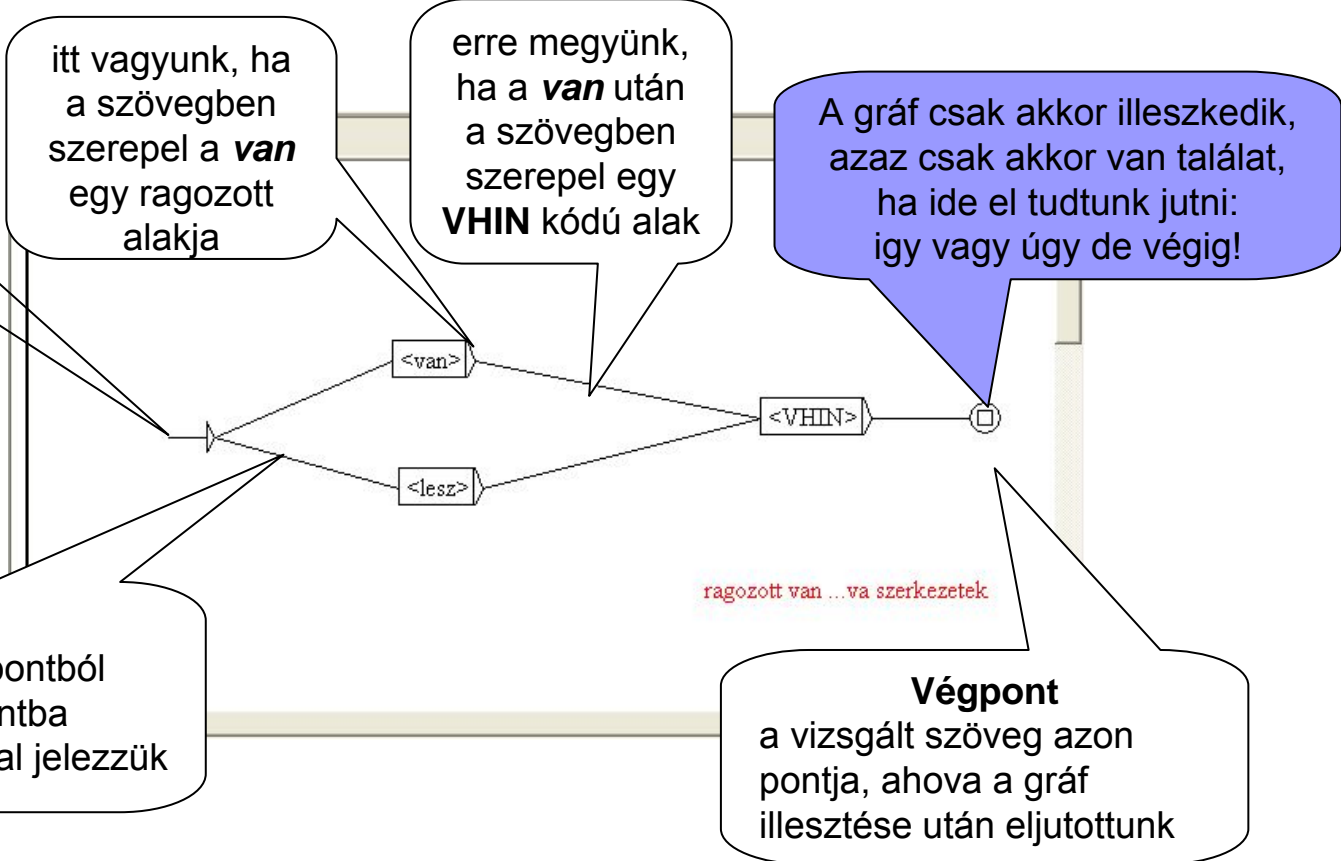
erre megyünk,
ha a **van** után
a szövegben
szerepel egy
VHIN kódú alak

A gráf csak akkor illeszkedik,
azaz csak akkor van találat,
ha ide el tudtunk jutni:
így vagy úgy de végig!

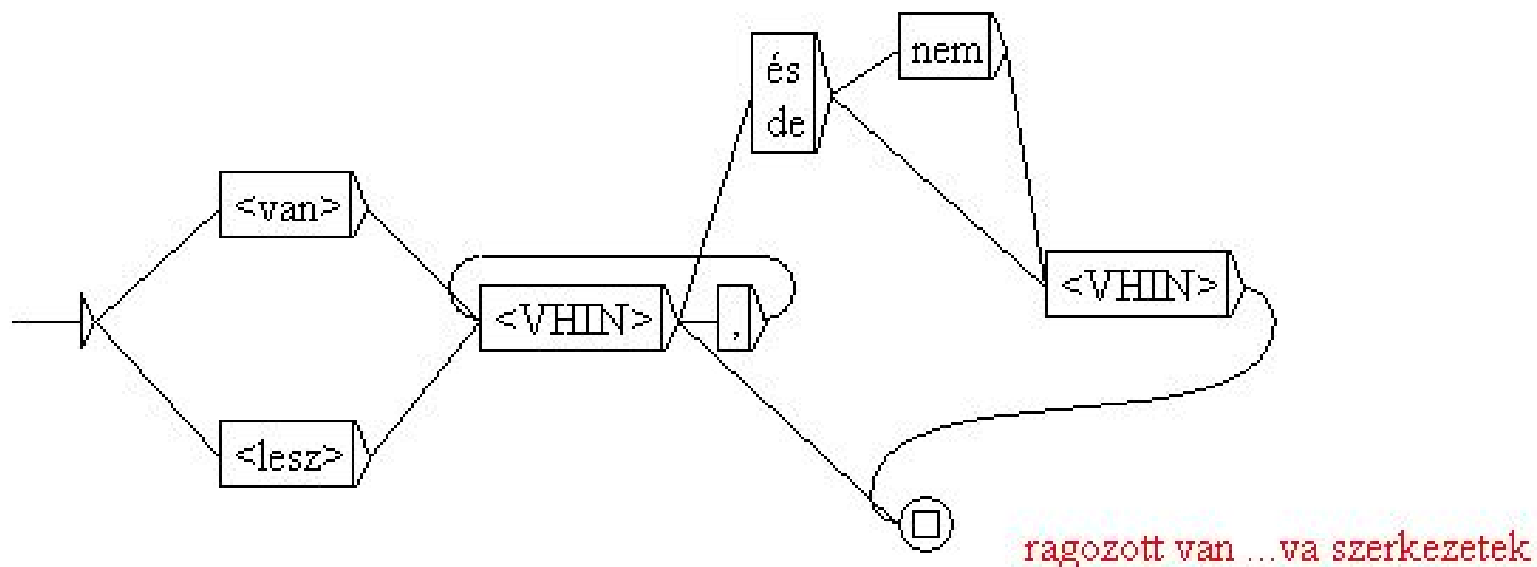
a **vagy** relációt
ugyanabból a pontból
ugyanabba a pontba
vezető több ággal jelezzük

ragozott van ...va szerkezetek

Végpont
a vizsgált szöveg azon
pontja, ahova a gráf
illesztése után eljutottunk

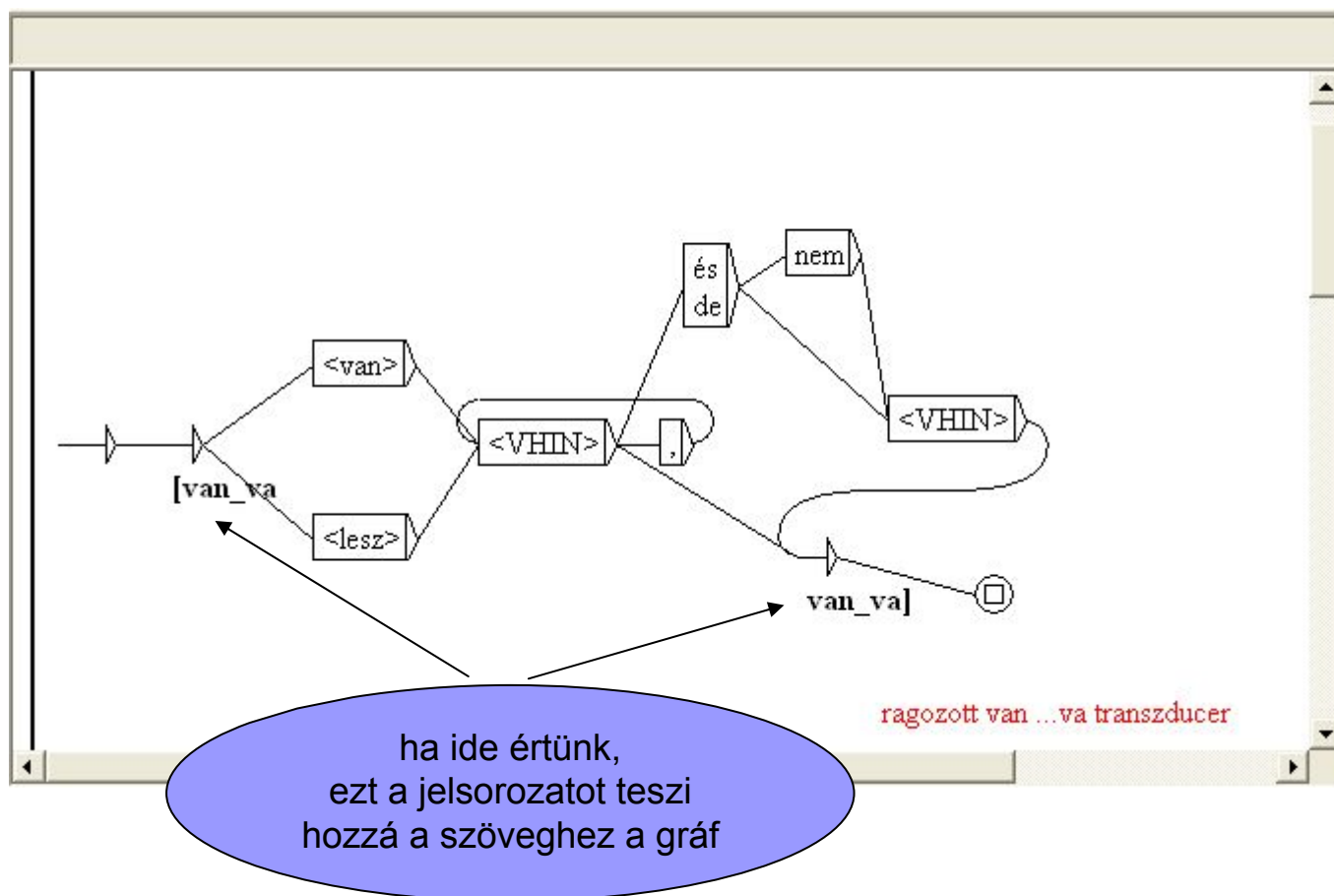


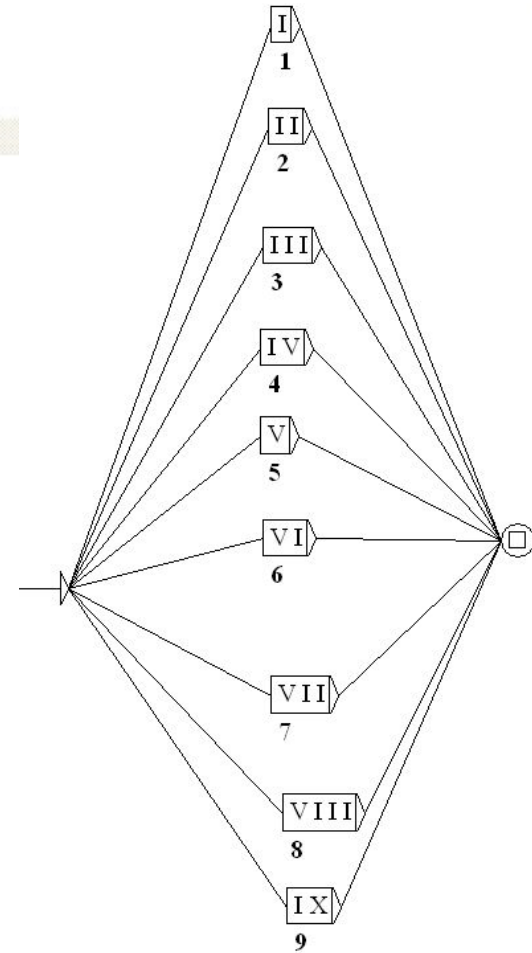
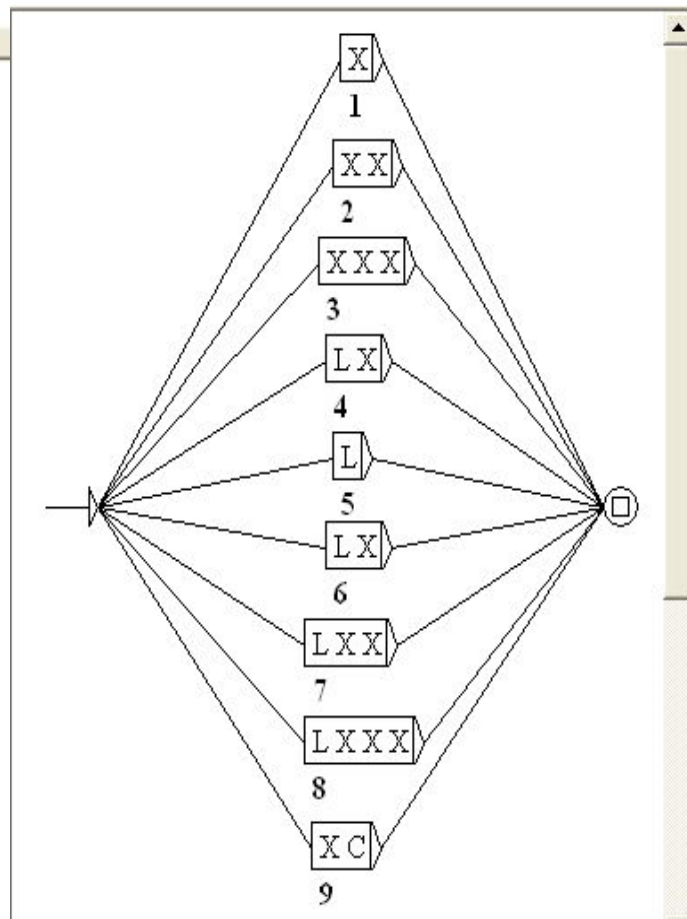
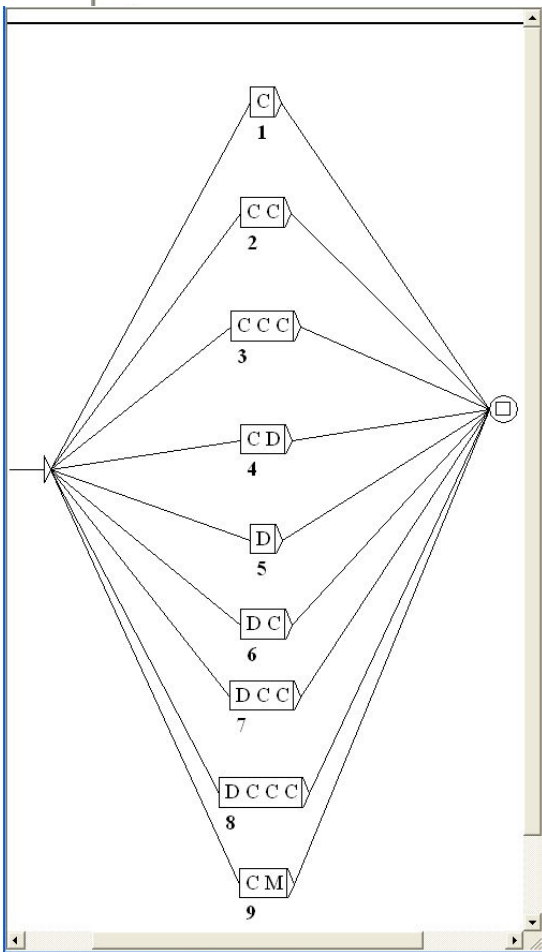
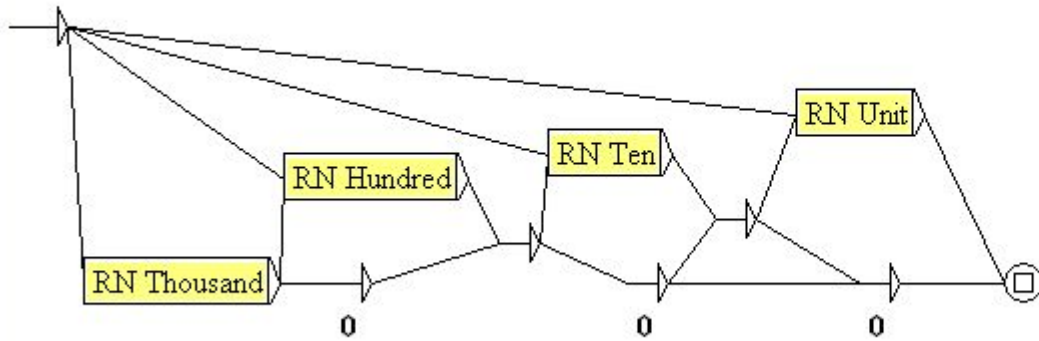
véges állapotú felismerő gráf (fsa)



véges állapotú transzducer (fst)

nemcsak felismer, hanem ki is ad jelsorozatot





Lexikai szabályok

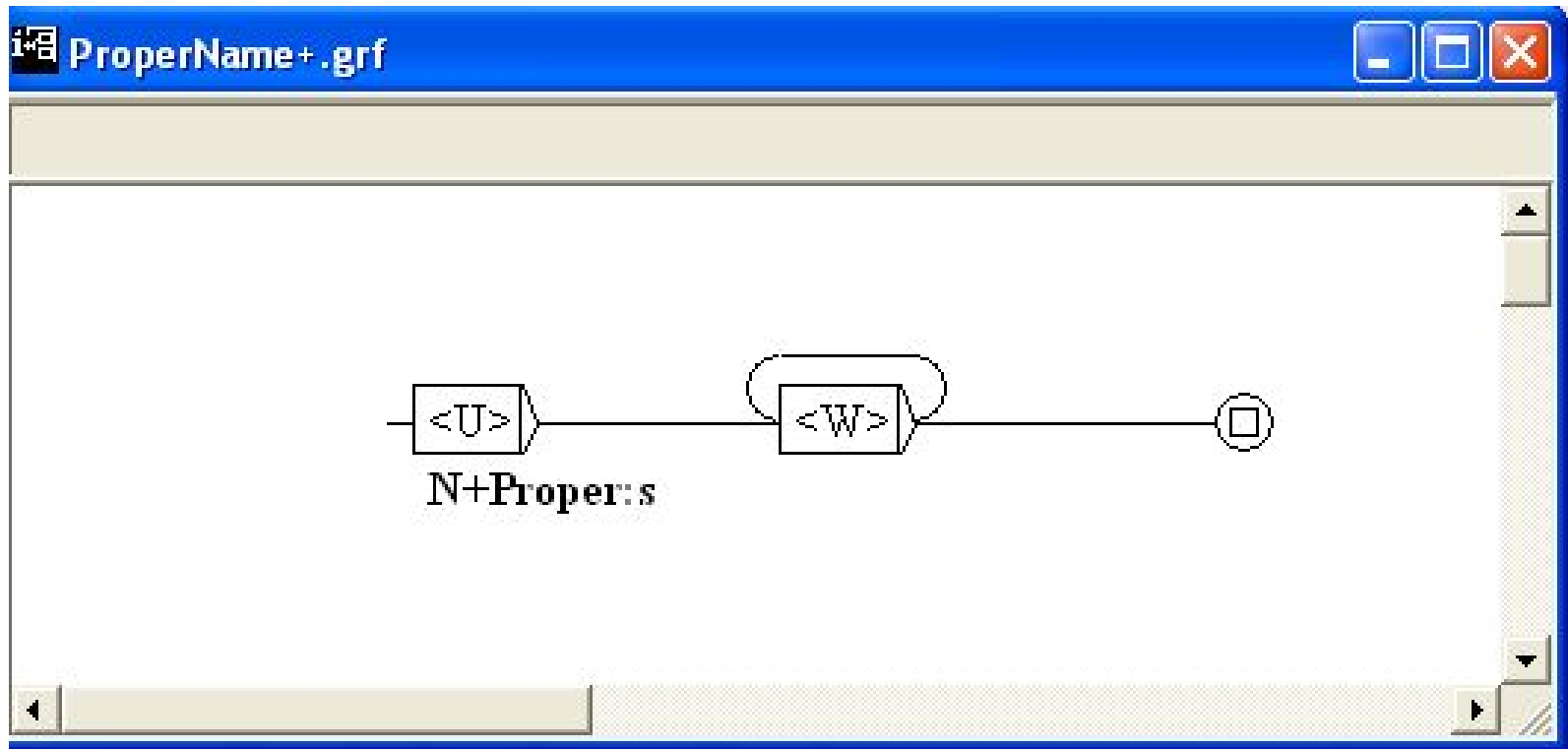
■ Lexikai transzducerek

- képzők, szóösszetételek elemzésére
- tulajdonnevek felismerésére

■ Lexikai megszorításokkal ellenőrizhetjük az elemek érvényességét

■ Kimenetük új szótári elemek, melyeket hozzátehetünk az alapszótárhoz

végző tulajdonnév szabály

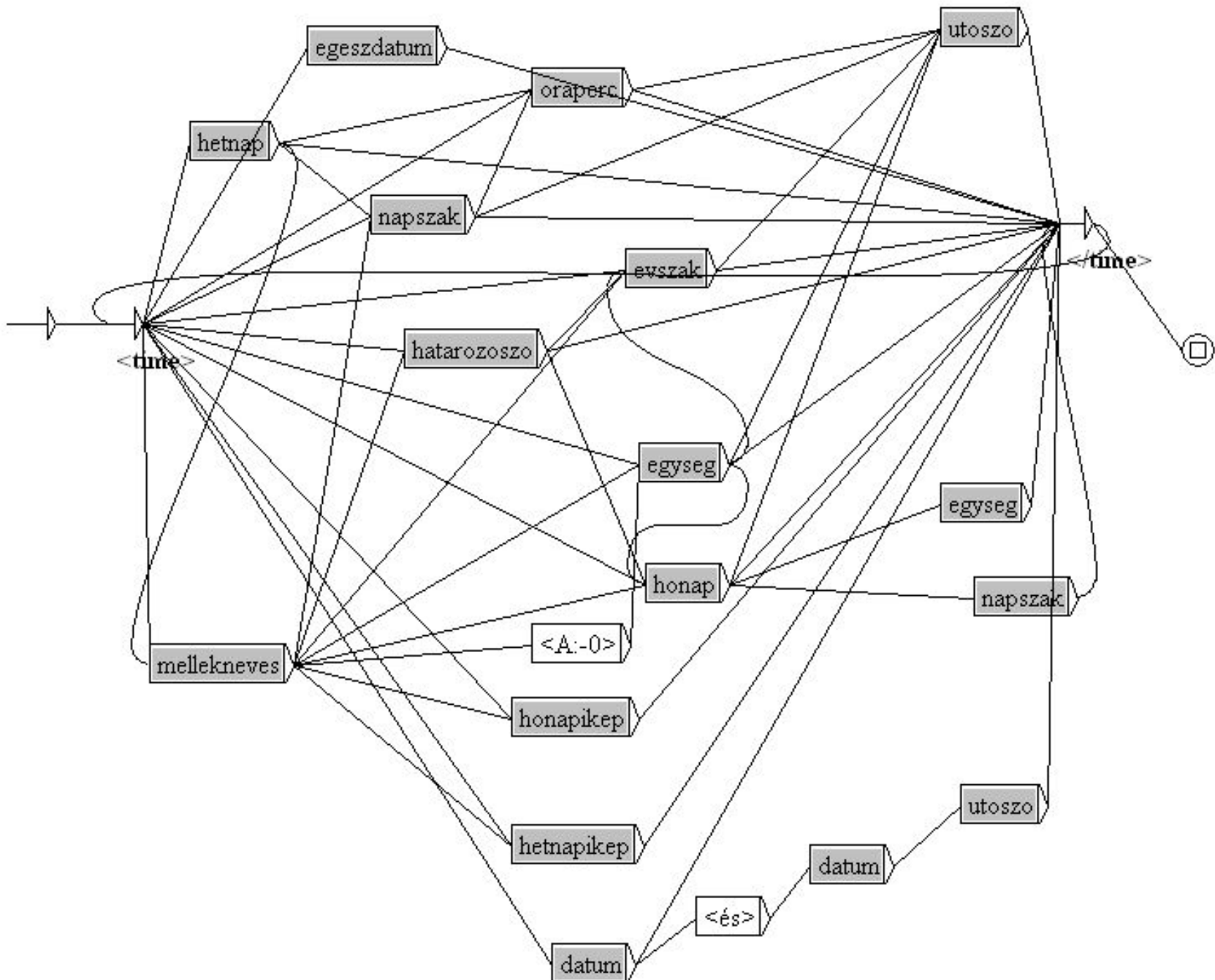


lexikai megszorítások

Tetszőleges
betűsorozatot
beolvasunk és a \$Tő
nevű változóba
teszünk

itt történik az ellenőrzés:
az Tő alakját e-vel
megtoldott alaknak
léteznie kell igeként a
szótárban

Ha minden rendben
megírhatjuk a szótári
bejegyzést:
storable,store.V+able



Magyar változat

- Kulcskérdés: a lexikai modul elkészítése
- szótár egyben morfológiai elemzést is ad
- a magyar morfológia listázással nem megoldható
- Kompromisszum:
 - óriási korpuszok elemzett szókincse
 - HUMOR morfológiai elemzővel elemezve
 - Magyar Nemzeti Szövegtár 150 m szó
 - Magyar webkorpusz (MOKK) 2500m szó

Az INTEX korlátai

- Gyors, hatékony, de sok beépített korlát
 - egyszerre egy nyelvvel lehet dolgozni
 - zárt rendszer – könnyű bevinni szövegeket, de nehéz kivinni őket
 - nincs xml kimenet
 - egyszerre egyetlen szöveggel lehet dolgozni
 - a szöveg annotációja lépcsőzetes elemzésnél nehézkes

JÖN! NooJ!

- Teljesen átdolgozott, többnyelvű, rugalmas, áttekinthetőbb rendszer
- INTEX utódja – elveiben azonos, de továbbfejlesztett változat
- Első nyilvános változat: december 15-én!
- A magyar változat hónapok óta előkészületben

Magyar NooJ

- Működő morfológiai elemző több változatban
 - nyers erő
 - morfológiai gráfok jegyek alkalmazásával
 - online elemző szótár paradigmatablák használatával
- Az Értelmező Kézisztár szókincse teljeskörű ragozással

NooJ - [news.dic]

File Lab Project Info Edit Windows

Dictionary contains 59213 entries

```
# NooJ V1
# Dictionary
#
# Input Language is: hu
#
# Alphabetical order is not required.
# Number of fields must be constant.
#
# Use any number of inflectional description files (.nof or .nog)
# Special Command: #use description.nof
# Special Command: #use description.nog
# DESCRIPTION FILES MUST BE STORED IN SAME DIRECTORY AS DICTIONARY
#
# Special Features: +FLX +UNAME
#
# Special Characters: '\ ' ' " ' + ' , ' # ' ' '
#

abaposztó, N+base+nostemvar+vowel+back+j+accnolink+nonlower
abba, PRO+base+nostemvar+vowel+back+j+accnolink+nonlower
abban, PRO+base+nostemvar
abbé, N+base+nostemvar+vowel+back+j+accnolink+nonlower
abbeli, A+base+nostemvar+vowel+front+unround+j+accnolink+nonlower
abból, PRO+base+nostemvar
abbreviatúra, N+obl+shorten+vowel+back+j+accnolink+nonlower
abbreviatúra, N+base+shorten+vowel+back+j+accnolink+nonlower
ábécé, N+base+nostemvar+vowel+front+unround+j+accnolink+nonlower
ábécérend, N+base+nostemvar+cons+front+unround+j+acclink+nonlower
ábécés, A+base+nostemvar+cons+front+unround+nonj+accnolink+nonlower
ábécés, N+base+nostemvar+cons+front+unround+nonj+accnolink+nonlower
ábécéskönyv, N+base+nostemvar+cons+front+round+nonj+acclink+nonlower
aberráció, N+base+nostemvar+vowel+back+j+accnolink+nonlower
ablak, N+base+nostemvar+cons+back+j+acclink+nonlower
ablakbélés, N+base+nostemvar+cons+front+unround+nonj+accnolink+nonlower
ablakbiztosítás, N+base+nostemvar+cons+back+nonj+accnolink+nonlower
ablakdeszká, N+obl+shorten+vowel+back+j+accnolink+nonlower
ablakdeszka, N+base+shorten+vowel+back+j+accnolink+nonlower
```

NooJ - [Vocabulary for Text: 1104.not [Modified]]

File Lab Project Info Edit Windows

12 Lexical entries:

Select All

Export lexemes

Freq	Entry	Lemma	Category	1	3	ACC	e	i	INS	PL	POS	PS	SUB	t
1	alannyal	alany	N	-	-	-	-	-	+	-	-	-	-	-
1	alanyt	alany	N	-	-	+	-	-	-	-	-	-	-	-
1	gyökeivel	gyök	N	-	+	-	+	+	+	-	-	+	-	-
1	gyökkel	gyök	N	-	-	-	-	-	+	-	-	-	-	-
1	gyökömet	gyök	N	+	-	+	+	-	-	-	-	+	-	-
1	gyököt	gyök	N	-	-	+	-	-	-	-	-	-	-	-
1	hajóimét	hajó	N	+	-	+	+	+	-	-	+	+	-	-
1	hajóival	hajó	N	-	+	-	+	+	+	-	-	+	-	-
1	hajóké	hajó	N	-	-	-	-	-	-	+	+	-	-	-
1	hajóval	hajó	N	-	-	-	-	-	+	-	-	-	-	-
1	házunkra	ház	N	+	-	-	-	-	-	-	-	+	+	+
1	házzal	ház	N	-	-	-	-	-	+	-	-	-	-	-

0 Unknown Tokens:

Select All

Export Unknowns

Freq	Word Form

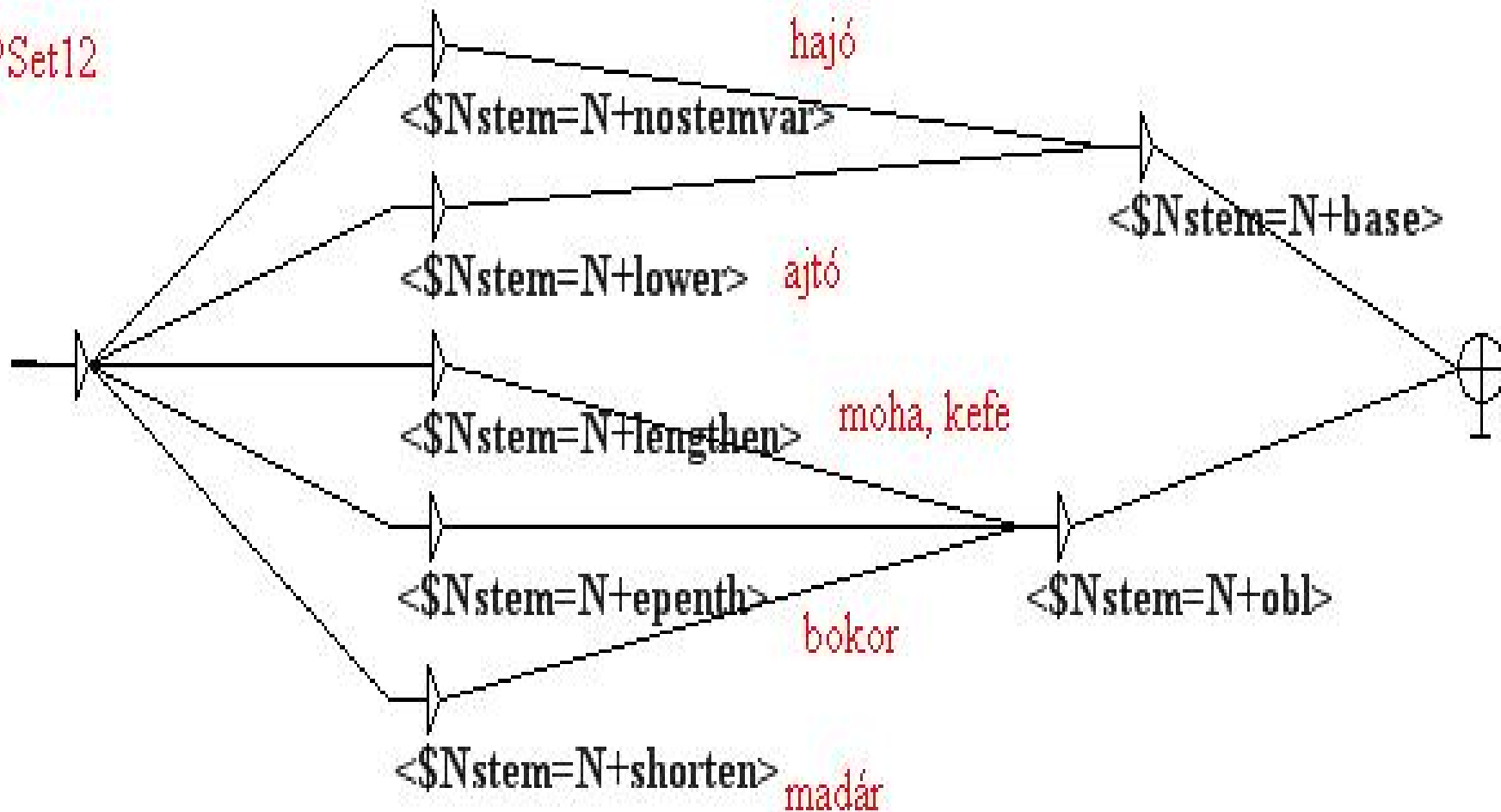
Reset

Filter

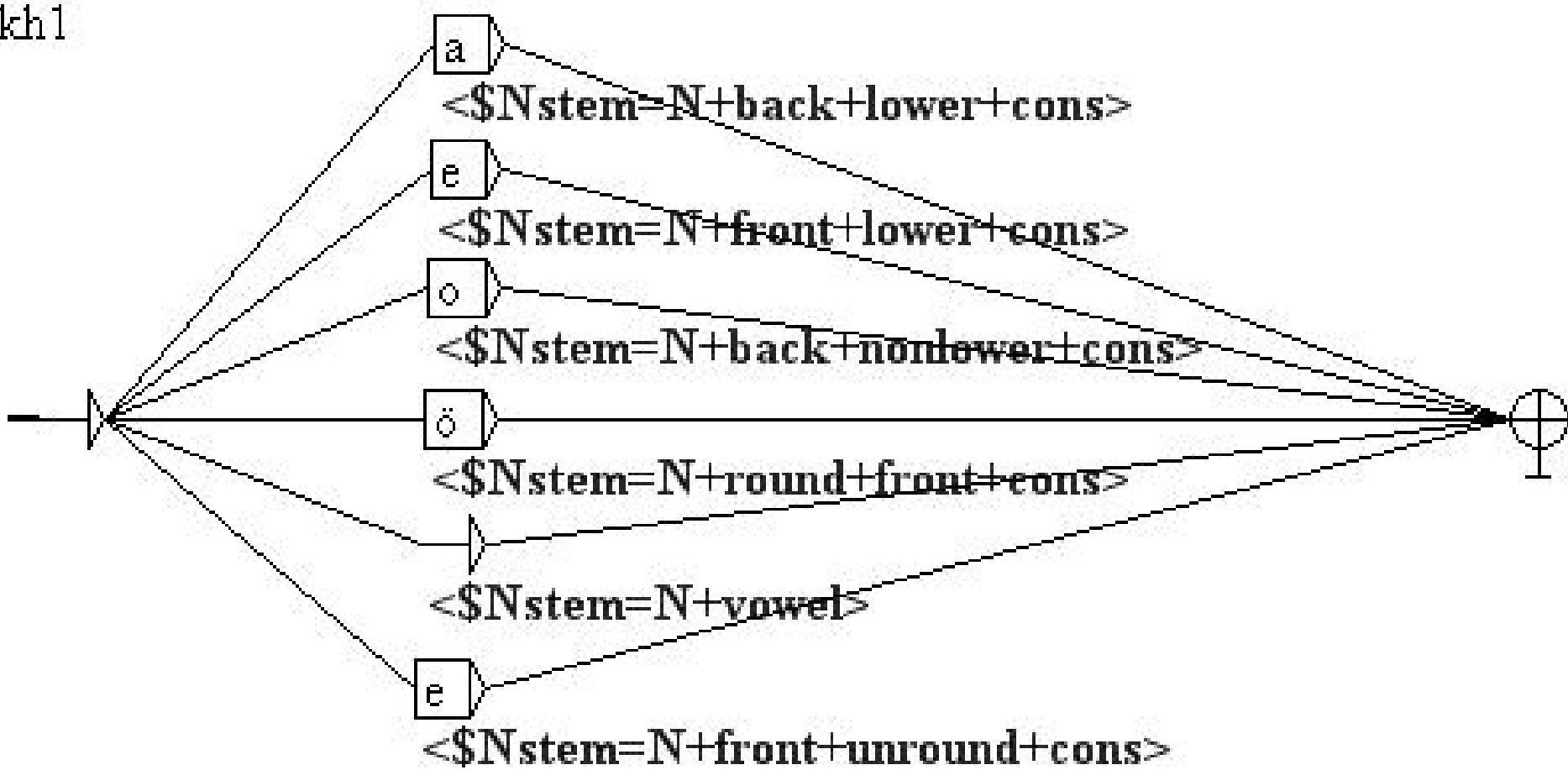


A tövvariánsjegyek ellenőrzése

var1 PSet12



kh1



Összegzés

- INTEX/NooJ általános célú integrált nyelvelemző rendszer
- Nemcsak számítógépes nyelvészeknek ...
- <http://corpus.nytud.hu/INTEX>
- Közreadjuk a magyar nyelvi eszközeinket a szakmai nagyközönség számára
- Várunk érdeklődőket, önkéntes segítőket, a kutatási eredmények megosztását

Végül, de nem utolsósorban...

- Külön köszönet kollégáimnak, akiknek a munkájáról szólt ez a beszámoló:

- Gábor Kata
- Nagy Viktor
- Vajda Péter
- Sass Bálint

- Oravecz Csaba
- Dancsecs Erzsébet
- Mészáros Ágnes
- Héja Enikő