

Derivation and Graphs

Duško Vitas

Problems

- The problem of unknown words in corpora processing
- The problem of the description of derivational and the consequences to e-dictionary organization (and dictionaries in general)

A question?

- Should each word be tagged, even if a tag might be wrong, or should words be left unprocessed if correct tags cannot be assigned to them?
- The relation between statistical methods and rule-based methods.

An example - srwac - ej!

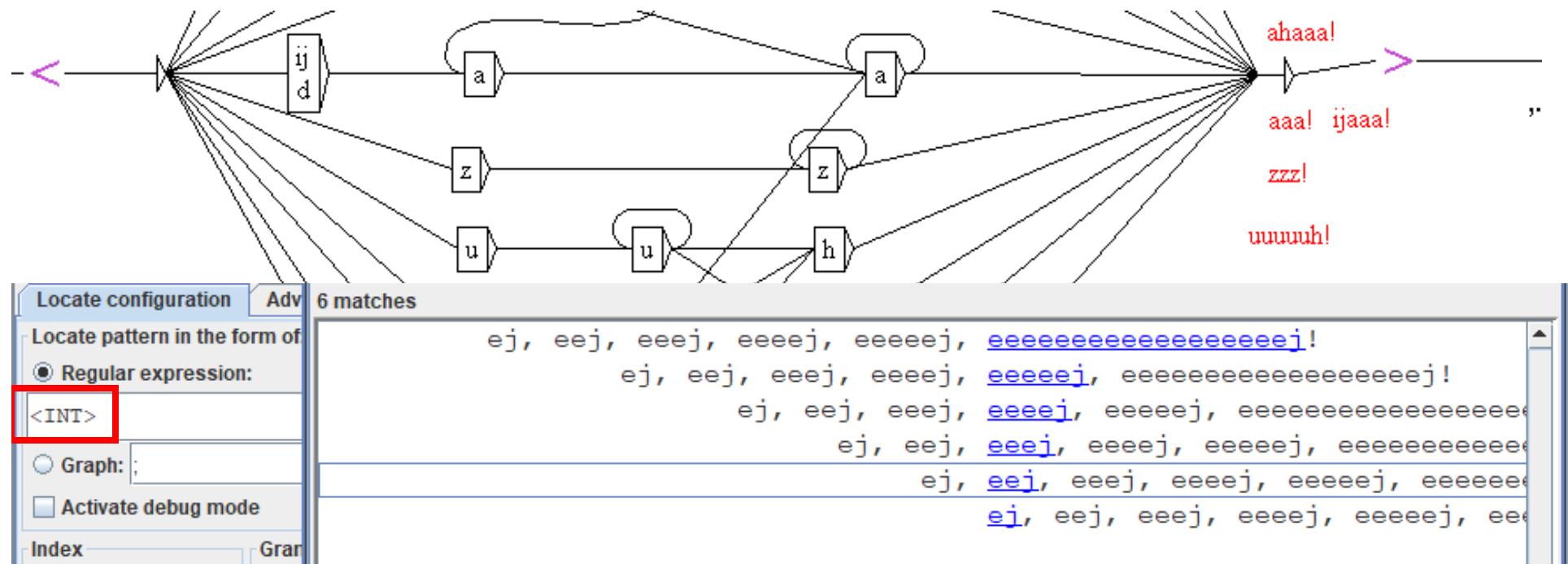
Imicu imala i Agrobanka koja
u poreze na trgovinu **g** Danas
nija Solel Solar Systems, koju
thodne godine. Godišnji rast
martu za 1,5 procenata, dok
noj faks, baš je druželjubiva,
im " efektu Kejt ": šta god da
Riki Fauler, Rori Mekllroj **eeⁿj, n < 20**. **g** Od njih trojice, najbolje je
ijekta IPA 2011 (čita se: **eeⁿj, n < 20**) Amerikanac je bio bolji od Di
Amerikanac je bio bolji od Di
za svetska tržišta nafte, jer
akcija **eeⁿj ~ INT & N & NProp & A &....**
e sasvi

ej /Rgp
ej /Sa
ej /Pd-msan
ej /Ncmsn
ej /Agpfsny
ej /Cc
ej /Pd3fsd
- /AT
Ej /Xf
ej /Pp3fsd

porasla 3,32 odsto uz promet od 6,9
drugi dan samita " Velike osmorke " u
koncern kupio 2009. za 418 mln. dolara
najbrži od septembra 2010. godine.
cena odeće obuće i lične opreme op
odmah mi je dala telefon **g** VOP: Ah
uradila ili nosila vojvotkinja od Kembla
g Od njih trojice, najbolje je
„ **eeⁿj, n < 20** **g** U toku je sadnja sezonske
Pointsa i Rorija Mekllroja. To mu je ova zemlja drugi na svetu potrošač r
n. evra. F
jos. Nara

With dictionary graphs

Input: **ej, eej, eeej, eeeej, eeeeej,
eeeeeeeeeeeeeeeeeej!**



Unknown words

... or unrecognized words by a system of e-dictionaries:

- No matter how comprehensive e-dictionary might be, new unknown words occur with each new text!
- Even the most ordinary newspaper article is usually not covered by e-dictionaries!

An example

The corpus of the weekly magazine **Vreme**
2010-2018:

29.930.758 (422.508 diff) tokens

13.062.037 (422.230) simple forms

...but in ERR:

**129.124 (46.315 diff) simple forms or
0.431% (of the text)**

~ 11% of diff. simple forms

Some examples from err

sveta na portalu www . **wikileaks** . org (u međuvremenu doncite ... ni e vo dobro **sećavanje**", glasi makedonski

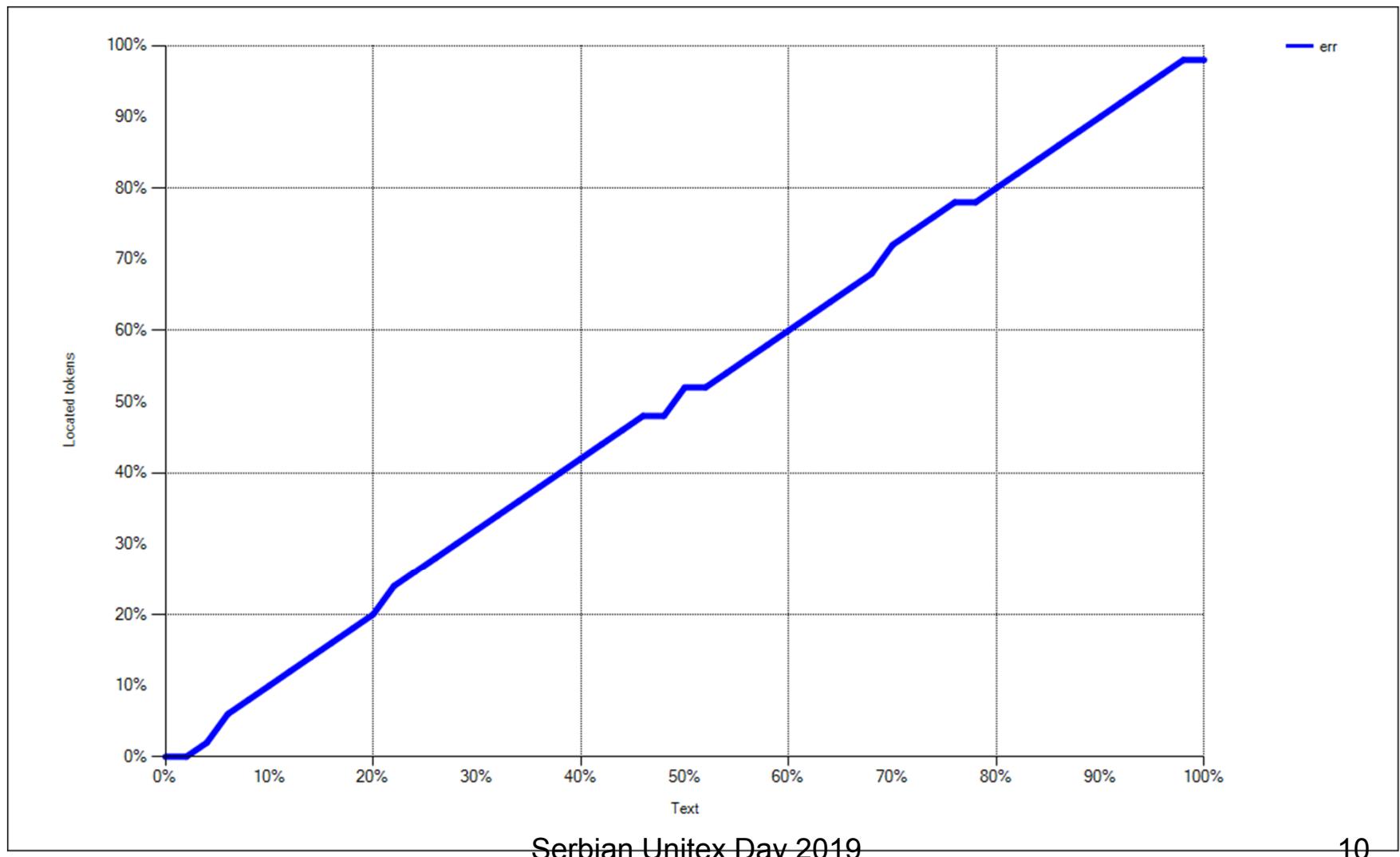
je da Srbija slučaj stavi **ad akta**, pošto je rezultat u

obila (što je, reći će naš **amerikanoFILski** autor, omaž m

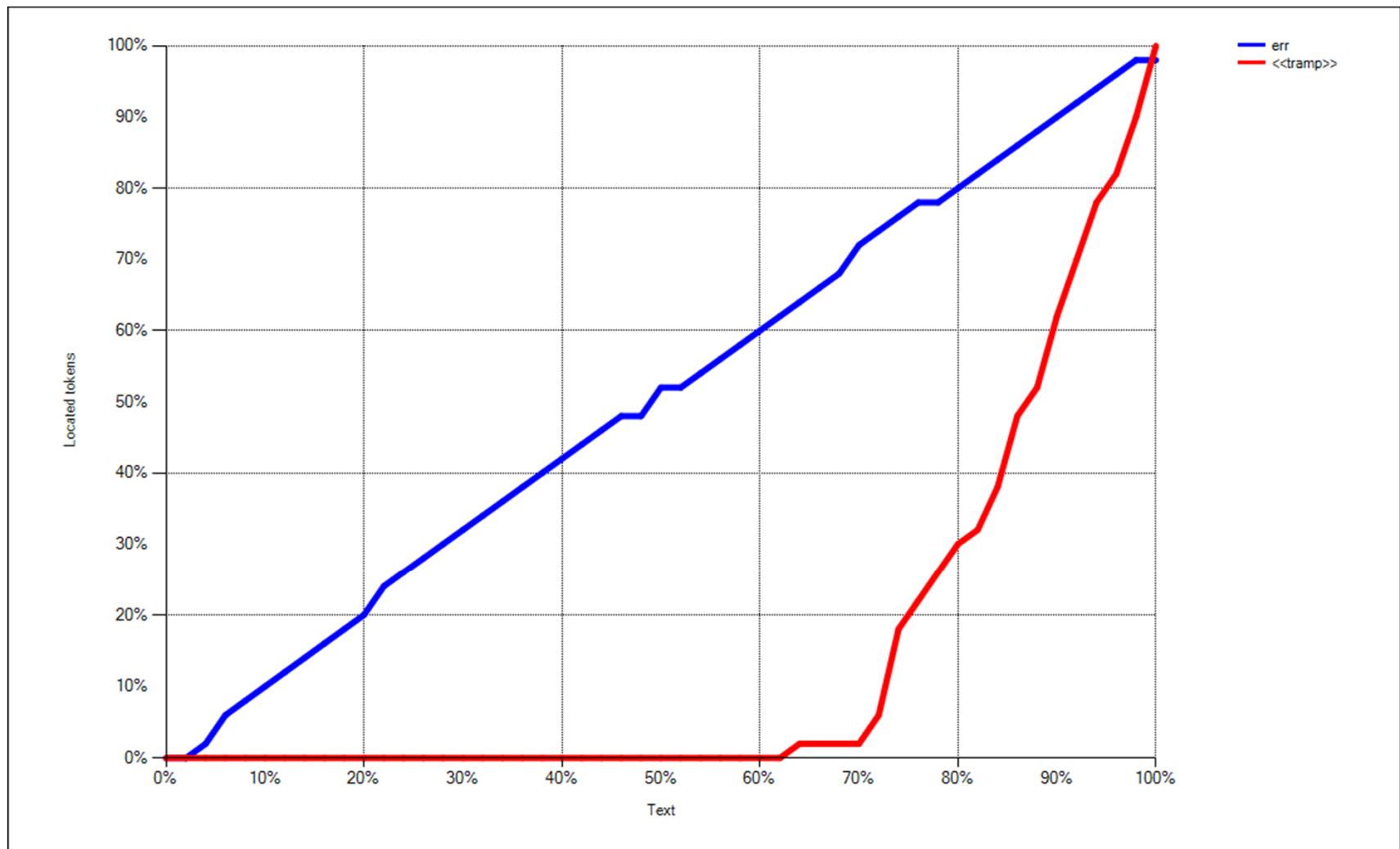
određenih mera protiv **antidejtonskih** političara, na nifestacije, još od doba **antifederalista**, koji su se

avidnu volju i umeće da **atrofirajućim** odnosima

Distribution of err vs. new word



Distribution of err vs. new word



An approximation

- If words tagged by graphs as **Unk** and **ABB** are included, and this set also includes some potentially incorrect words, then numbers are:
**908.326 (877.006 diff) simple forms or
3.034%**
- (correct solutions are included).

Examples of Unk and ABB

N+NProp+Unk:

J. Evaristo, Monti, **Arico** Suarez – Peucelle, **OK**
istrage (takozvana **Chilcot Inquiry**).{S} Javn
sno SIS-a (Secret **Intelligence Service**).{S} **OK?**
iranjem u lepoti" i **Ronaldinjovim** majstorijam **A**
anizmu čoveka.{S} **Biočip** je zapravo mala pl **Not NP**

ABB:

priznao u vestima **B92** da jeste reč o njegov: **B is ABB**
javnosti.{S} I MI5 i **MI6** izgubile su zbog zlou **MI is ABB**
a objašnjavaju da **rad** tajnih službi, uz sve **OK**
a u goste 1994. u **SAD** i 1998. u Francusku, **OK**

<<merkel>> (1300 occ.)

In DLF:

Merkel,.N+NProp+Hum...

Merkelova,Merkelov.A+Pos+N
Prop+Hum+Last+... ???

Merkelici,.N+NProp+Unk

Merkelina,.N+NProp+Unk ?

In ERR:

merkelice

("*merkelice*" ovog sveta)

merkelovci

(*merkelovci* svih partija)

merkelovih

(od svih tih makrona ,
merkelovih , junkera i
rahoja)

merkelovski

Twitter in 10 years



Derivation

Derivational dictionary

Worth, S., A. Kozak, D. Johnson: *Russian Derivational Dictionary*. RAND Corporation, 1970.

МАРКС

suffixes

МАРКС	ЙЗМ
МАРКС	ЙСТ
МАРКС	ЙСТ К А
КВАЗИ	МАРКС ИСТ
НЕ	МАРКС ЙСТ
ПО -	МАРКС ЙСТ СК И
	МАРКС ЙСТ СК ИЙ
АНТИ	МАРКС ЙСТ СК ИИ
ПСЕВДО	МАРКС ЙСТ СК ИЙ
ДО	МАРКС ЙСТ СК ИЙ

Derivational dictionary

Worth, S., A. Kozak, D. Johnson: *Russian Derivational Dictionary*. RAND Corporation, 1970.

prefixes

МАРКС

КВАЗИ
НЕ
ПО -

АНТИ
ПСЕВДО
ДО

МАРКС	ЙЗМ
МАРКС	ЙСТ
МАРКС	ЙСТ К А
МАРКС	ЙСТ
МАРКС	ЙСТ
МАРКС	ЙСТ СК И
МАРКС	ЙСТ СК ИЙ
МАРКС	ЙСТ СК ИЙ
МАРКС	ЙСТ СК ИЙ

Capital letters

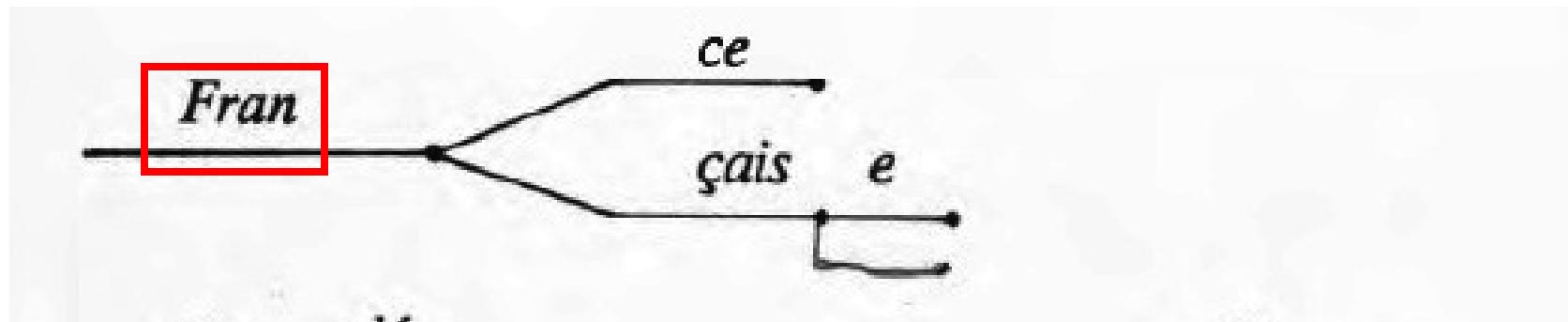
МАРКС

МАРКС ЈЗМ
МАРКС ЈСТ

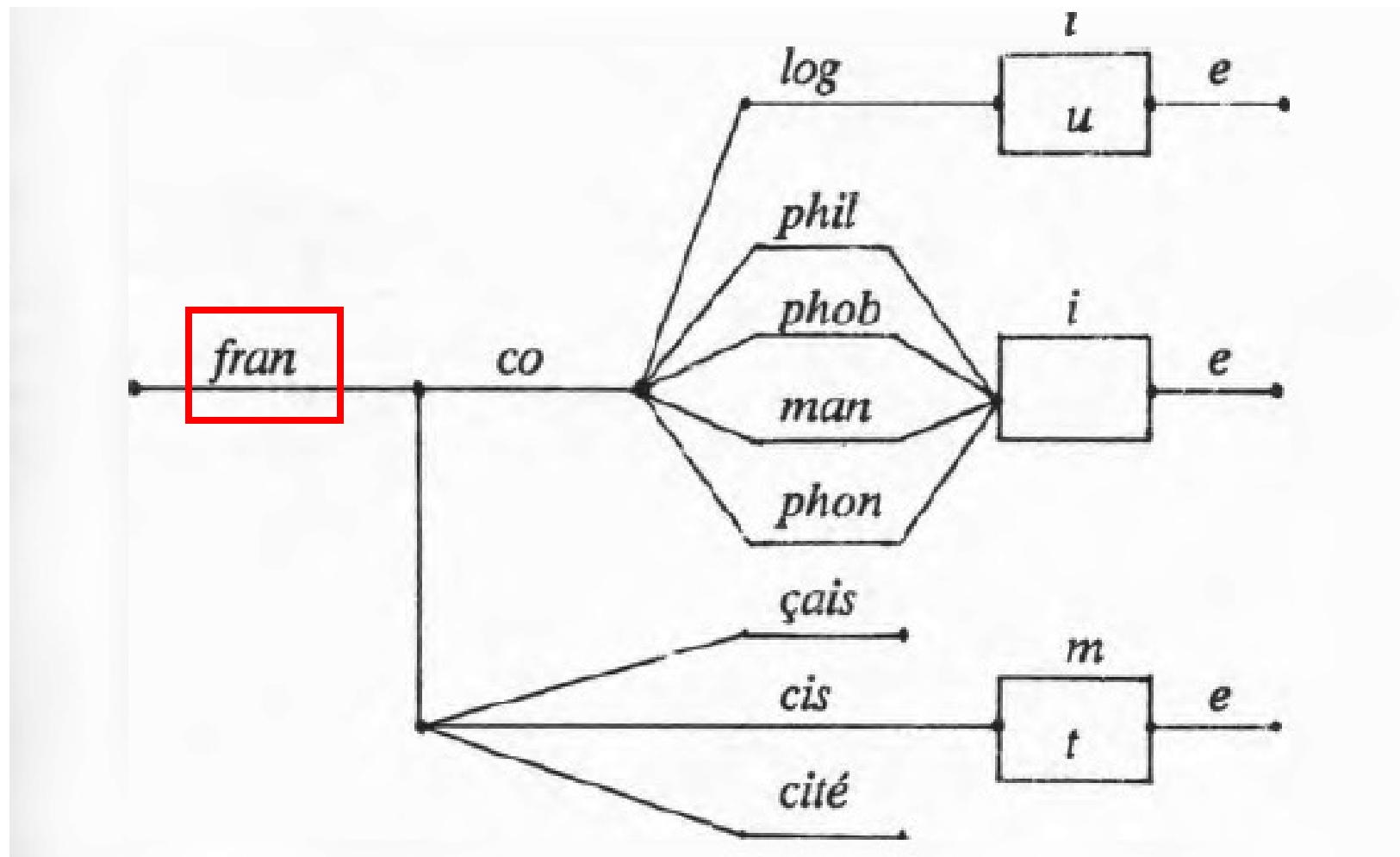
Marks, but marksizam

The representation done by Gross

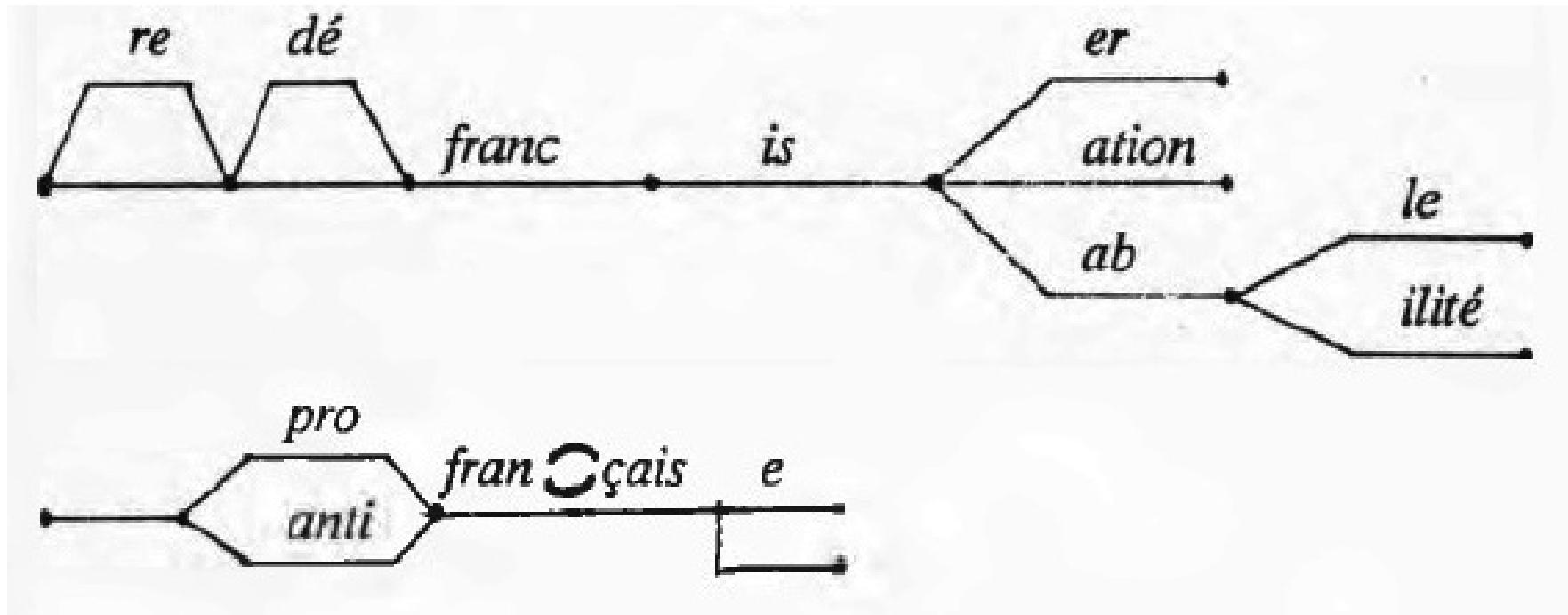
Gross, Maurice: *The Use of Finite Automata in the Lexical Representation of Natural Language*. LNCS 377, pp. 34-50, 1989



The factorization to the right - suffixes



The factorization to the left - prefixes



The case of Serbian

In dictionary RMS6:

gluma

glumac

glumački, -a, -o

glumački

glumiti

glumica

glumičin

glumčev

glumčina

glumčić

n.

n.

adj.
...
related to actors

adv. in a way of actors

v. to act

n. actress (gender motion)

adj. possessive of *glumica* (actresse)

adj. possessive of *glumac* (actor)

n. augmentative of *glumac* (actor)

n. diminutive of *glumac* (actor)

Redundancy (and incompleteness)

In dictionary RMS6:

gluma	n. 1. acting art
glumac	n. actor
glumački, -a, -o	adj. related to actors
glumački	adv. in a way of actors
glumiti	v. to act
glumica	n. actress (gender motion)
glumičin	adj. possessive of <i>glumica</i> (actresse)
glumčev	adj. possessive of <i>glumac</i> (actor)
glumčina	n. augmentative of <i>glumac</i> (actor)
glumčić	n. diminutive of <i>glumac</i> (actor)

Prolex 3.0

Beograd

N

beogradski

A+Rel

Beograđanin

N+Hab:m

beogradjanski

A+Rel

Beograđaninov

A+Hab+Poss

Beograđanka

N+Hab:f

Beograđankin

A+Hab+Poss

Regular derivation

For the animate (human) nouns:

gender motion: *glumac* --> *glumica*

possessive adjectives: *glumac* --> *glumčev*

diminutives: *glumac* --> *glumčić*

augmentatives: *glumac* --> *glumčina*

For adjectives:

negation: *prištojan* --> *neprištojan*

abstract nouns: *mudar* -> *mudrost*

For verbs:

verbal noun: *pevati* --> *pevanje*

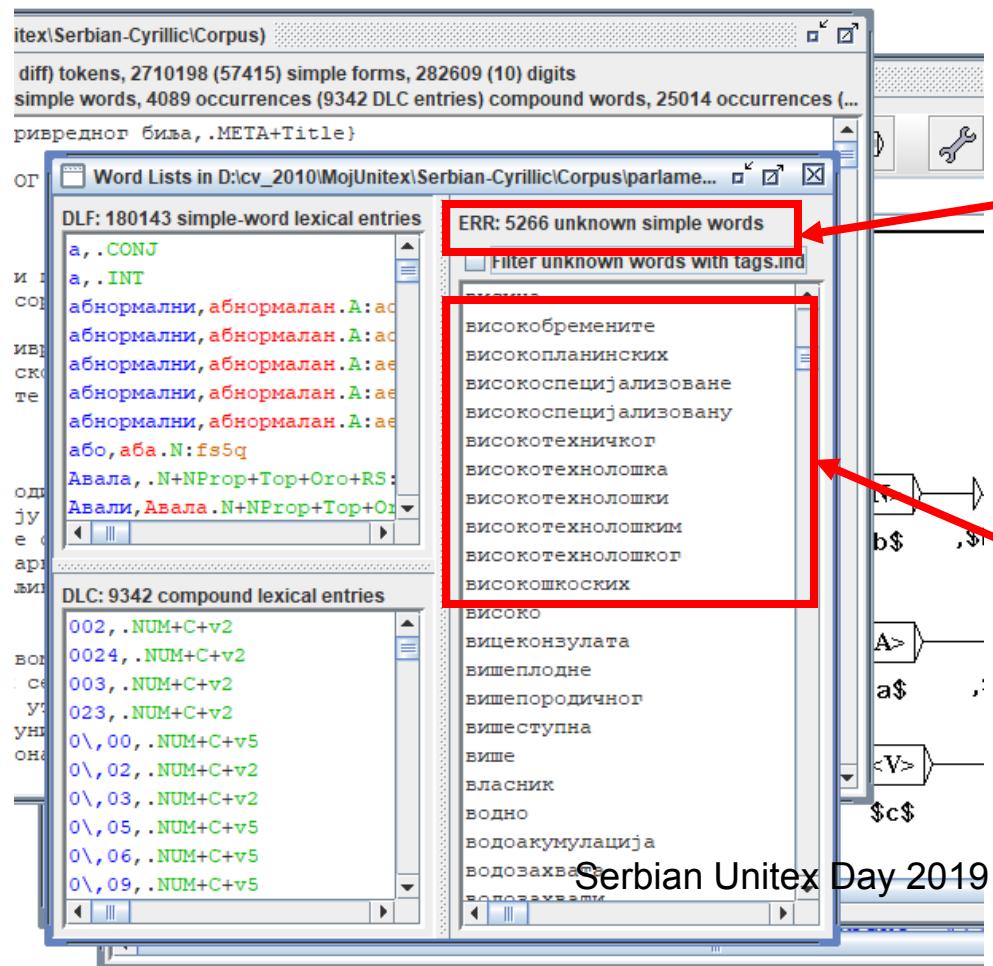
verb forms acting as adjectives. *atrofirati* --> *atrofirajući*

Derivation in Unitex

1. Morphological graphs
2. Super - lemma

Morphological graphs

The corpus of Serbian laws (2013) ~ 3MW

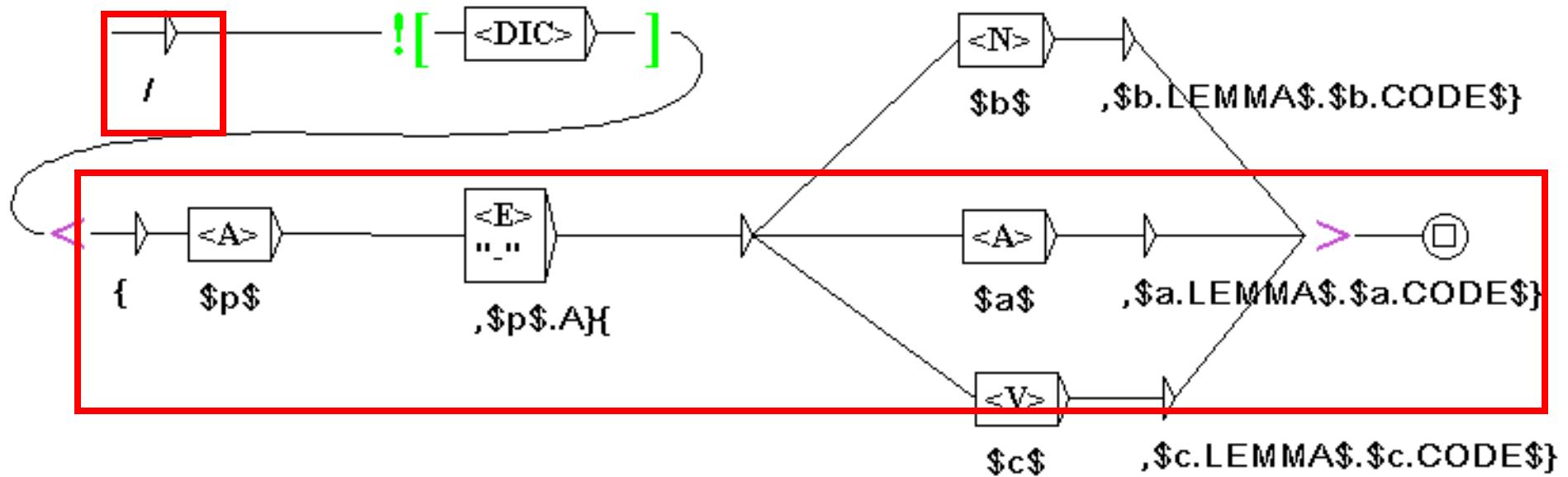


in ERR ~ 5266

visokobremenite
visokoplaninskih
visokospecijalizovane
visokotehničkog
visokotehnološkog

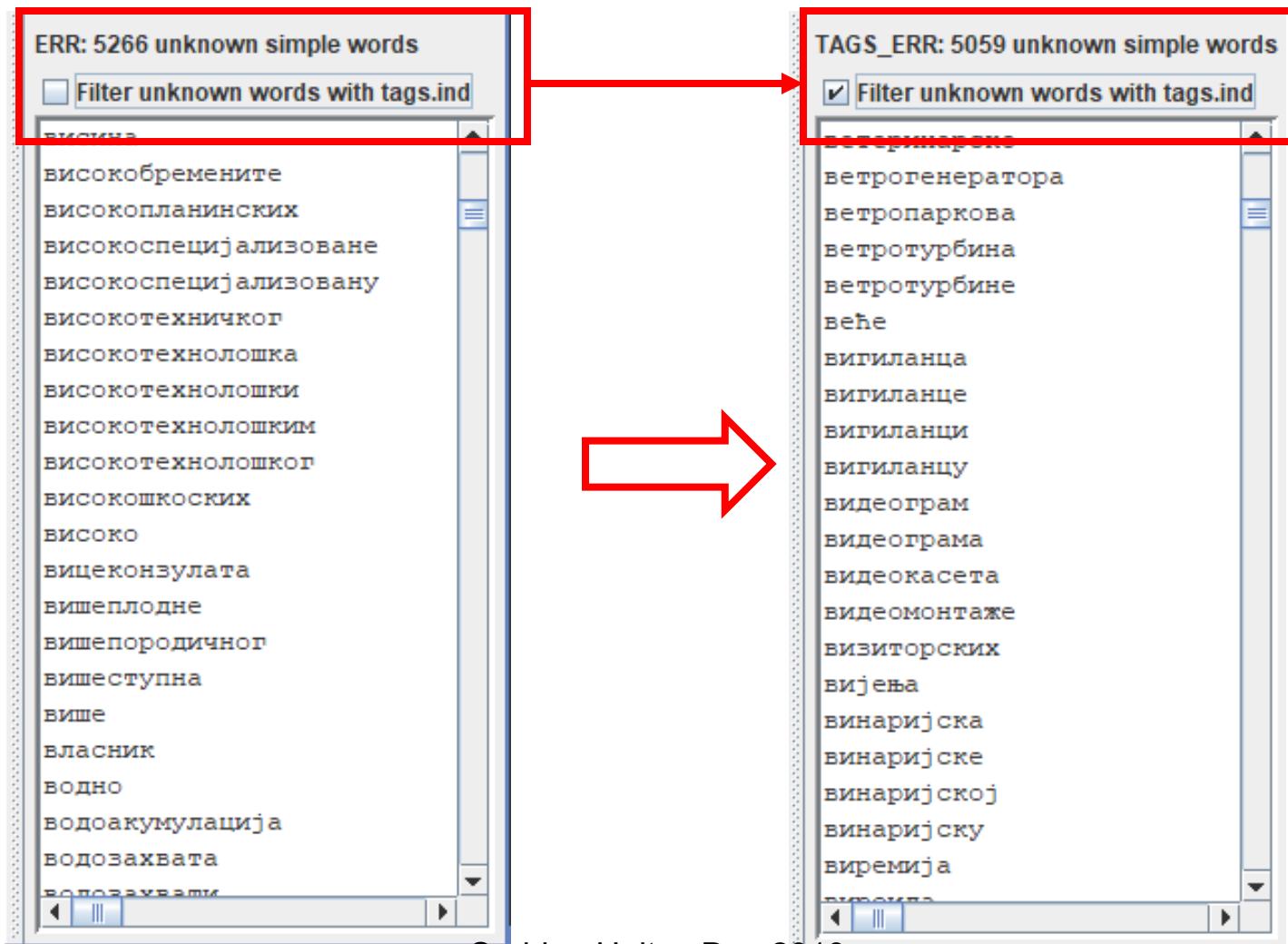
...

~ A<E>A



}Члан 26. Забрањено је превозити: {S1} [/{{високо,.A}{бремените,брименит.A:aefp5g}}](#) женке у перио
}Члан 26. Забрањено је превозити: {S1} [/{{високо,.A}{бремените,брименит.A:aefs2g}}](#) женке у перио
}Члан 26. Забрањено је превозити: {S1} [/{{високо,.A}{бремените,брименит.A:aefw2g}}](#) женке у перио
}Члан 26. Забрањено је превозити: {S1} [/{{високо,.A}{бремените,брименит.A:aefw4g}}](#) женке у перио
}Члан 26. Забрањено је превозити: {S1} [/{{високо,.A}{бремените,брименит.A:aemp4g}}](#) женке у перио
а Члан 17. Заштита екосистема (шумских,
а Члан 17. Заштита екосистема (шумских,
а Члан 17. Заштита екосистема (шумских,
датке који се односе на: {S1} обављање [/{{високо,.A}{планинских,планински.A+PosQ:aefp2g}}](#), воден
датке који се односе на: {S1} обављање [/{{високо,.A}{планинских,планински.A+PosQ:aemp2g}}](#), воден
датке који се односе на: {S1} обављање [/{{високо,.A}{планинских,планински.A+PosQ:aenp2g}}](#), воден
датке који се односе на: {S1} обављање [/{{високо,.A}{специјализоване,специјализован.A+PP+DER=OvatiIrr}}](#)
ика је здравствена установа која обавља [/{{високо,.A}{специјализовану,специјализован.A+PP+DER=OvatiIrr}}](#)

tags.ind



Another way

- Examples of unknown words from the **Vreme** corpus having the prefix **anti-**

čno splasnula rasistička, "antiafrička" hysterija koja je pratila Gvatljivi, baš kao i plan B antiakademskog lobija koji nalaže da ako dine bile su veoma tražene antiapokaliptične pilule i maske koje su ene je kao antinacionalnu, antiautomobilsku i antikerozinsku partiju je, zahvaljujući svojoj antibriskoj politici, nadmašila Liber provinijencija: od onih sa antidejtonskom platformom po kojima se, onalnu, antiautomobilsku i antikerozinsku partiju zahvatio ledeni tij Zemlje. {S}Zelene je kao antinacionalnu, antiautomobilsku i antik

anti.snt (D:\cv_2010\MojUnitex\Serbian-Latin\Corpus\aa-fds)

5 sentence delimiters, 373 (148 diff) tokens, 169 (137) simple forms, 4 (3) digits
157 occurrences (512 DLF entries) simple words, 4 occurrences (6 DLC entries) compound words, 12 occurrences (12 ERR lines) unknown words

I u Libiji je konačno splasnula rasistička, "antiafrička" hysterija koja je pratila Gadafijev pad i povratak njegovih vernih tuareških ratnika u Mali i Niger.

{S} Sve i da je briga iskrena, lobija koji nalaže da ako PA sertifikata PA neograničenom {S} Te godine bile su veoma trsvakom čošku i služile kao za {S} Zelene je kao antinacional talas: njihove teme, koje su istorijom.

{S} Stranka za nezavisnost Uje podršku, u međuvremenu je, za demokrate za dva procenta.

{S} Sad su stvari uglavnom pre antidejtonskom platformom po automatski vraća na Ustav Rep

Word Lists in D:\cv_2010\MojUnitex\Serbian-Latin\Corpus\aa-fds\ant...
DLF: 512 simple-word lexical entries

- ako, .CONJ
- ako, .PAR
- automatski, .A+PosQ:admslg
- automatski, .A+PosQ:aemplg
- automatski, .A+PosQ:aemp5g
- automatski, .A+PosQ:aems4q
- automatski, .A+PosQ:aems5g
- automatski, .ADV+Adj
- B, .ABB+Mes+SI=B
- b, .ABB+Mes+SI=b

DLC: 6 compound lexical entries

- 1992, .NUM+C+v2+NVAL=1992
- dejtonski sporazum, Dejtonski
- dejtonski sporazum, Dejtonski
- u međuvremenu, .ADV+Comp
- u vreme, .PREP+Comp+Tmp+p2
- ujedinjenog kraljevstva, Ujed

ERR: 12 unknown simple words

- Filter unknown words with tags.ind
- antiafrička
- antiakademskog
- antiapokaliptične
- antiautomobilsku
- antibriselskoj
- antidejtonskom
- antikerozinsku
- antinacionalnu
- suspects
- tuareških
- tzv
- usual

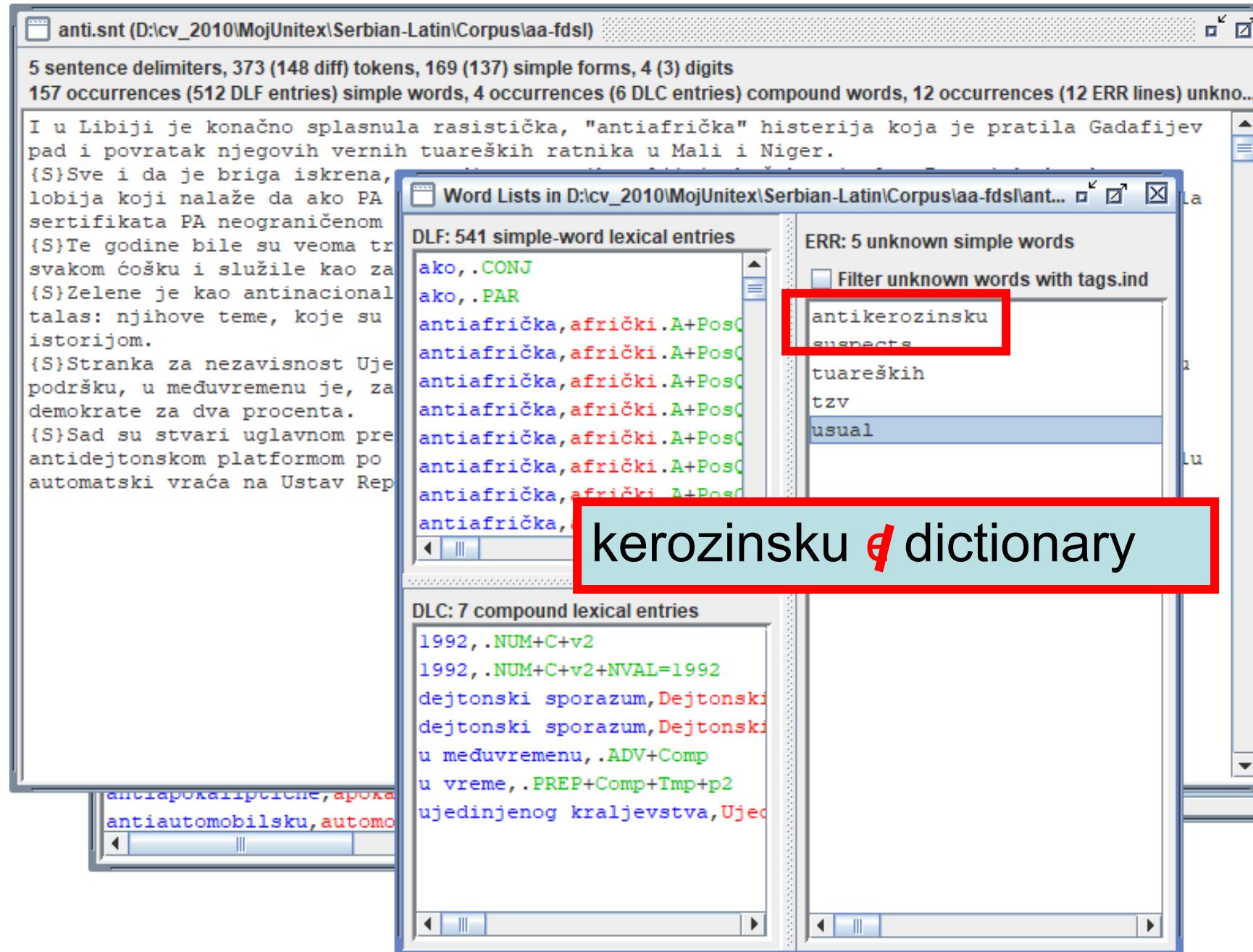
antidejtonskim, dejtonsk

Serbian Unitex Day 2019

$\vdash ![<\text{DIC}>] < \text{anti} < [<\text{A}>] , \$x\$. \text{LEMMA}\$. \$x\$. \text{CODE}\$ \Box$

The screenshot shows a software interface with a search results window. The title bar of the window reads "D:\cv_2010\MojUnitex\Serbian-Latin\Delalanti.txt". The window displays a list of 2138 matches, with the first few entries shown below:

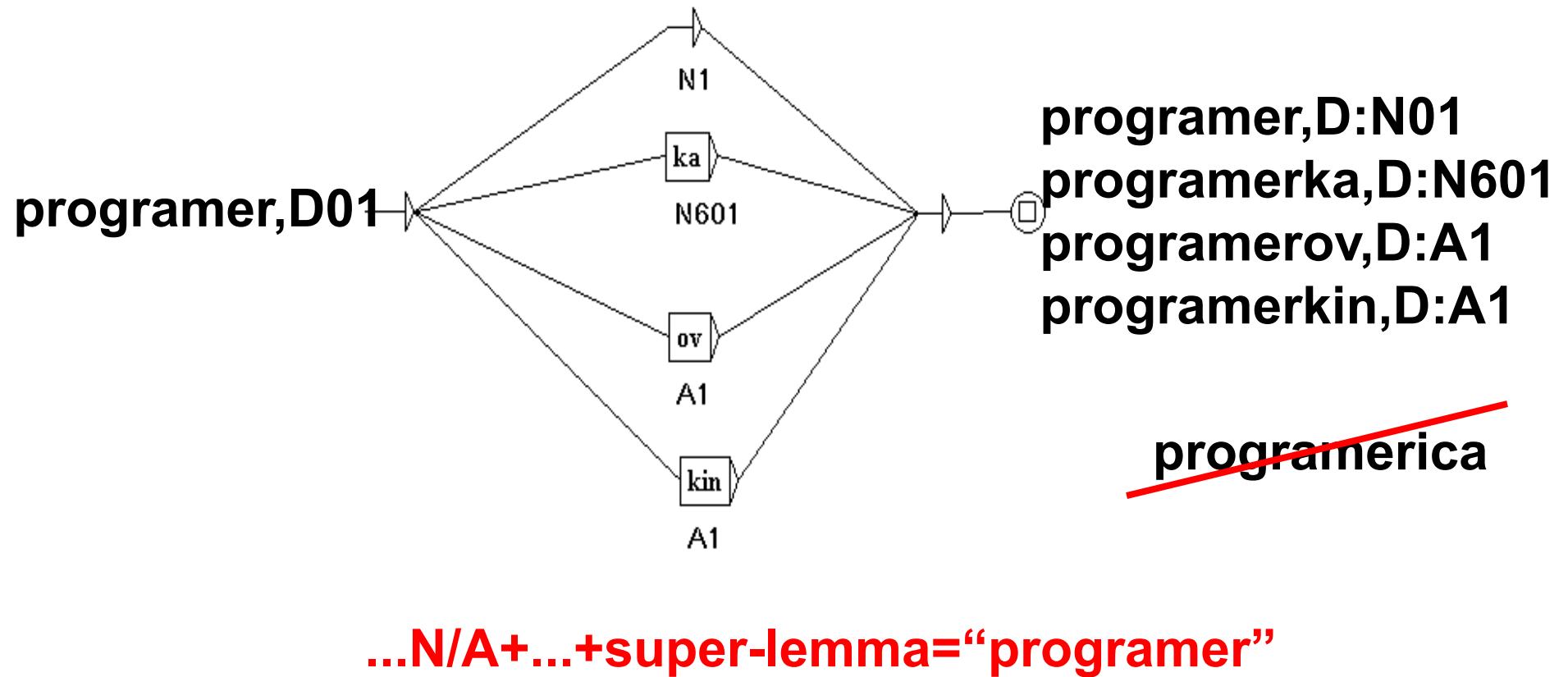
```
//{anti afrička, afrički.A+PosQ+NProp+Top+Kont+SupReg:aefslg  
//{anti afrička, afrički.A+PosQ+NProp+Top+Kont+SupReg:aefs5g  
//{anti afrička, afrički.A+PosQ+NProp+Top+Kont+SupReg:aemw2g  
//{anti afrička, afrički.A+PosQ+NProp+Top+Kont+SupReg:aemw4g  
//{anti afrička, afrički.A+PosQ+NProp+Top+Kont+SupReg:aenw2g  
//{anti afrička, afrički.A+PosQ+NProp+Top+Kont+SupReg:aenw4g  
//{anti afrička, afrički.A+PosQ+NProp+Top+Kont+SupReg:aenplg  
//{anti afrička, afrički.A+PosQ+NProp+Top+Kont+SupReg:aenp4g  
//{anti afrička, afrički.A+PosQ+NProp+Top+Kont+SupReg:aenp5g  
//{anti afrička, afrički.A+PosQ+NProp+Top+Kont+SupReg:akms2g  
//{anti afrička, afrički.A+PosQ+NProp+Top+Kont+SupReg:akms4v  
//{anti afrička, afrički.A+PosQ+NProp+Top+Kont+SupReg:akns2g  
//{anti akademskog, akademski.A+PosQ:adms2g  
//{anti akademskog, akademski.A+PosQ:adms4v  
//{anti akademskog, akademski.A+PosQ:adns2g  
//{anti apokaliptične, apokaliptičan.A:aefw2g  
//{anti apokaliptične, apokaliptičan.A:aefw4g  
//{anti apokaliptične, apokaliptičan.A:aefplg  
//{anti apokaliptične, apokaliptičan.A:aefp4g  
//{anti apokaliptične, apokaliptičan.A:aefp5g  
//{anti apokaliptične, apokaliptičan.A:aefs2g  
//{anti apokaliptične, apokaliptičan.A:aemp4g  
//{anti automobilsku, automobilski.A+PosQ:aefs4g
```



„The super-lemma“

- ... as a mean to gather forms obtained by the regular derivation.
- We should remember that, as an example, the regular derivation produces from the lemma **programer** forms:
 - gender motion: **programer** --> **programerka**
 - possessive adjectives: **programer**--> **programerov**
 - diminutives: **programer** --> **programerčić**
 - augmentatives: **programer** --> **programerčina...**

„The super-lemma“ - graph D01



The super-lemma *programer*

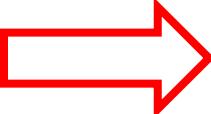
programer,D:N01

DELAS:

programerka,D:N601

programer,N01

programerov,D:A1



programerka,N601+sl="programer"+FG

programerkin,D:A1

programerov,A1+Poss+
sl="programer"

programerkin,A1+Poss
+sl="programer"



DELAF

The classification problems

The processing of unknown words can be divided to three general cases:

1. A really unknown word (neologism, occasionalism, a foreign word,...) – a specific processing
2. A word from DELAS to which a prefix is added (a morphological graph)
3. A word derived by suffixation from a word from DELAS (a generalized inflectional graph)

Possible consequences

- In the example of the noun *glumac* and its derivatives, instead of having ten or more separate entries that refer to *glumac*, only one entry appears with a „code“ that points to possible derivatives.
- Problems? Re-defining of inflectional graphs?
- Applications? Maintenance of dictionaries, but also avoiding redundant lexicographic descriptions.

Thank you!