**18TH EURALEX INTERNATIONAL CONGRESS**

Lexicography in global contexts
17-21 July 2018
Ljubljana, Slovenia

# The Dictionary of the Serbian Academy: from the Text to the Lexical Database

https://e-knjige.ff.uni-lj.si/znanstvena-zalozba/catalog/view/118/211/2975-1

Ranka Stanković[1]  Rada Stijović[2]  Duško Vitas[1]  Cvetana Krstev[1]  Olga Sabo[2]

[1] Human Language Technology group at the University of Belgrade
[2] Institute for Serbian Language, Serbian Academy of Sciences and Arts Serbia

**EURALEX**

# Motivation

❖ **Dictionary of Serbian Academy or DSA**

✓ **19 volumes printed, 20th in print, ~15 more to come**

✓ **covers written resources of the standard Serbo-Croatian (from 19th century to the present day) of all Shtokavian dialects**

✓ **complex microstructure**

❖ **Task:**

✓ **automatic transformation of the digitized text of DSA into various standard structured formats and into a lexical database**

✓ **use the lexical base of the Dictionary for research purposes and for the production of different derived lexicographic products**

# Model of Lexical Database

> **Analysis of formatting conventions used for typesetting dictionary entries and identification of triggers:**
> 1. **headword group,**
> 2. **grammatical data,**
> 3. **etymology,**
> 4. **sense,**
> 5. **multiword expressions,**
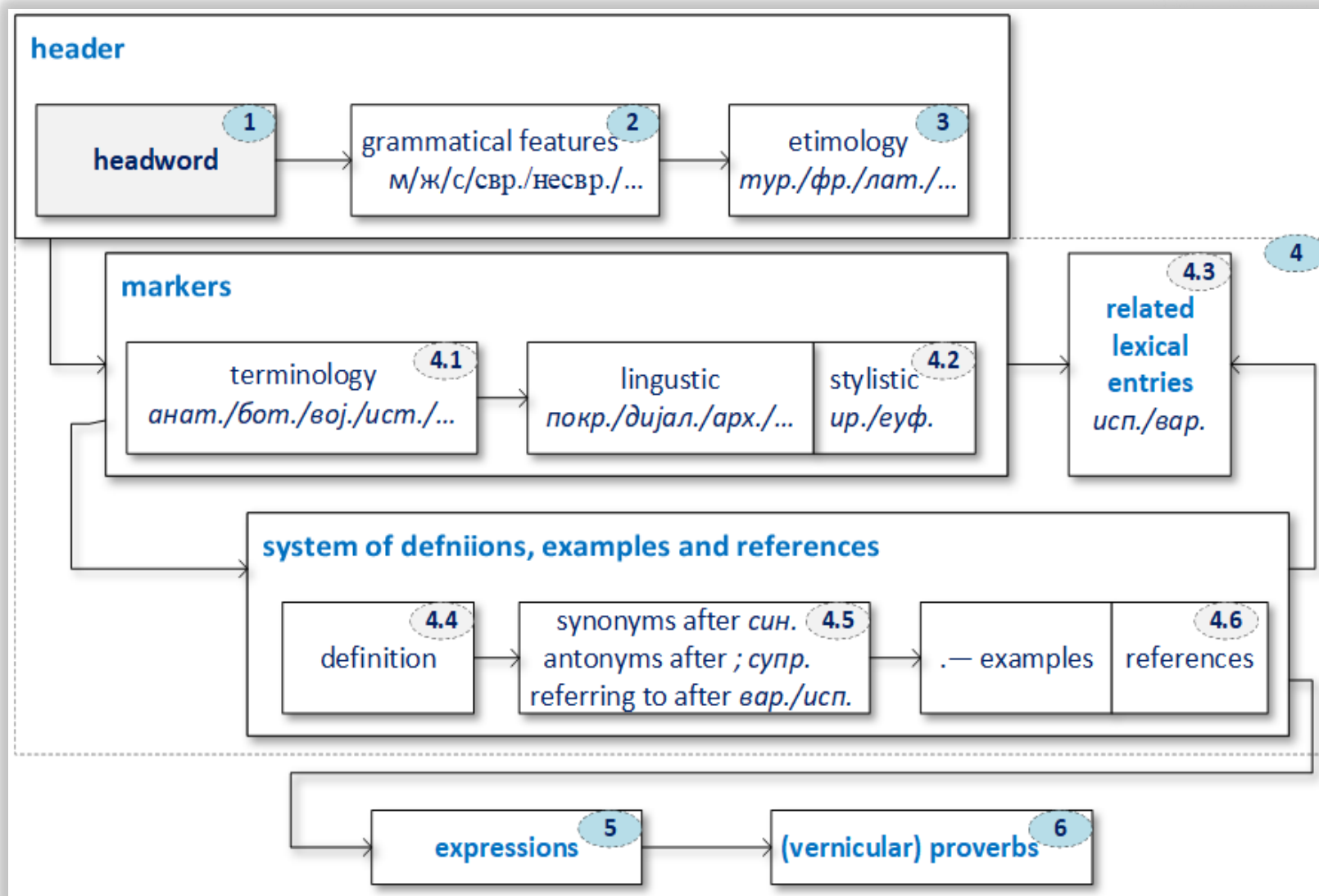> 6. **proverbs.**

# Model of Lexical Database

➢ **Dictionary markers**
- ✓ **semantic, phonetic, morphological, syntactic, normative, functional, stylistic, domain of use,..**
- ✓ **371 abbreviations (markers) mapped as data category values**
- ✓ **30 data categories, further grouped in data-category set**

# Typographic conventions and triggers as applied to nouns

| Element | | Example | typography | trigger begin | Trigger end |
|---|---|---|---|---|---|
| Headword group | | | | | |
| | lemma | па̏леоце̄н | \<nl\> bold | | comma or trigger begin |
| | gramm. data | -а | | hyphen | |
| | lemma | палео̀це̄н | \<nl\> bold | и | comma or trigger begin |
| | gramm. data | -ена | | hyphen | |
| gramm. data | | м | | item in a list | |
| Etymology | | palaiós kainós | | Open parenthesis + item in a list, e.g. „(грч.“ | closing parenthesis |
| Sense | | | | 1, 2, 3 or а, б, в or I, II or trigger begin | |
| | terminological markers | геол. | | item in a list | trigger begin |
| | linguistic/markers | / | | item in a list | trigger begin |
| | related words | / | | Some punctuation marks; item in a list | trigger begin |
| | definition | прва, најстарија епоха палеогена. | italic | | trigger begin |
| | synonyms, antonyms, related | / | | item in a list | |
| Example | example text | Формације геолошке се даље дијеле... | | dash | |
| | Bibliographic references | Д-П1, 17 | | Open parentheses | Closing parenthesis |
| MWE | | / | | Изр. | |
| Proverbs. | | / | | НПосл. | |

# Microstructure of dictionary articles

**header**

| | | |
|---|---|---|
| **headword** (1) | → | grammatical features м/ж/с/свр./несвр./... (2) | → | etimology тур./фр./лат./... (3) |

**markers** (4)

terminology анат./бот./воj./ист./... (4.1) → lingustic покр./диjал./арх./... stylistic ир./еуф. (4.2) → related lexical entries исп./вар. (4.3)

**system of defniions, examples and references**

definition (4.4) → synonyms after *син.* antonyms after *; супр.* referring to after *вар./исп.* (4.5) → .— examples references (4.6)

expressions (5) → (vernicular) proverbs (6)

# Transformation: dictionary article text form -> lexical database

EURALEX

**Typographic conventions and triggers as applied to nouns**

па̏леоце̄н, -а и палео̀це̄н, -ена м (грч. palaiós kainós) геол. *прва, најстарија епоха палеогена.* — Формације [геолошке] се даље дијеле на ... епохе. Тако се ... терцијар [састоји] од пет: палеоцена, еоцена, олигоцена, миоцена и плиоцена (Д–П 1, 17). (Калм. Р. 1, 81; Р. МС).

{RegEx}

C# .NET

*Segmentation, Pattern recognition, Alignment of markers*

| па̏леоце̄н | -а | палео̀це̄н | -ена | м |
|---|---|---|---|---|
| грч. | palaiós kainós | геол. | | |

*прва, најстарија епоха палеогена.*

Формације [геолошке] се даље дијеле на ... епохе. Тако се ... терцијар [састоји] од пет: палеоцена, еоцена, олигоцена, миоцена и плиоцена

| Д–П 1, 17 | Калм. Р. 1, 81; Р. МС |
|---|---|

```xml
<entry n="3971">
  <form type="lemma">
    <orth> па̏леоце̄н </orth>,</form>
  <form type="inflected">-а</form>
  <form type="lemma">
    <orth> палео̀це̄н </orth>,</form>
  <form type="inflected">-а</form>
  <gramGrp>
    <gen>м</gen>
  </gramGrp>
  <sense>
    <etym>(<lang>грч.</lang> palaiós kainós) </etym>
    <usg type="dom"> геол.</usg>
    <def> прва, најстарија епоха палеогена.</def>
    <cit> Формације [геолошке] се даље дијеле на
    ... епохе. Тако се ... терцијар [састоји] од пет:
    палеоцена, еоцена, олигоцена, миоцена и
    плиоцена <bibl>(Д–П 1, 17).</bibl>(Калм. Р. 1,
    81; Р. МС) </cit>
  </sense>
</entry>
```
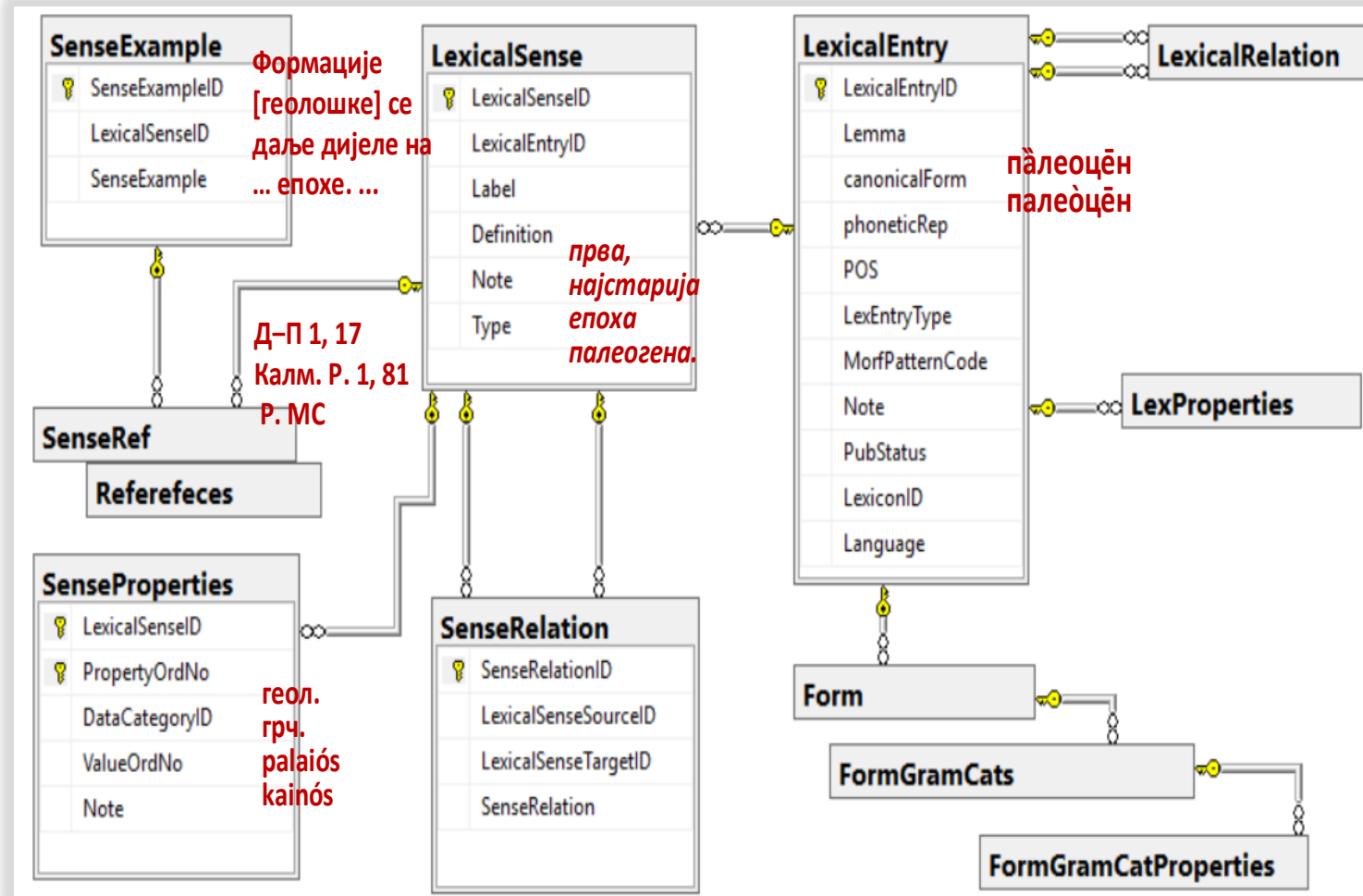
*TEI compliant tagging of a dictionary article*

*Lexical database*

SQL Server

**The Dictionary of the Serbian Academy: from the Text to the Lexical Database**

# Transformation:
# dictionary article text form -> lexical database

# Results and Evaluation <span style="color:orange">1st (1959) and 19th (2014) volumes</span>

**EURALEX**

1) **Automatic procedure recognized, structured, annotated and stored in the lexical database <u>15,988</u> dictionary articles from 1st and <u>11,153</u> from 19th volume**

2) **Digitization of all previously published volumes is in progress**

3) **In future, the linear processing of entries may be abandoned which would possibly accelerate the production process.**

4) **The supplements to the already published volumes could also be produced.**

5) **Once the DSA is fully populated, users with different levels of accessibility, will be able to search through its lexical data base.**

The Dictionary of the Serbian Academy: from the Text to the Lexical Database

1st (1959) and
19th (2014) volumes



Part of speech

| | null | noun | verb | adjective | adverb | other |
|---|---|---|---|---|---|---|
| Volume 1 | 598 | 10633 | 1364 | 2645 | 670 | 160 |
| Volume 19 | 520 | 7808 | 1192 | 1391 | 259 | 35 |

Grammatical gender

| | м -masculine | ф - feminine | с - neuter | combination |
|---|---|---|---|---|
| Volume 1 | 4410 | 4012 | 1311 | 60 |
| Volume 19 | 3127 | 2749 | 795 | 49 |

Verbal aspect

| | несвр. - imperfective | свр. - perfective | несвр. свр. | гл. |
|---|---|---|---|---|
| Volume 1 | 105 | 872 | 383 | 4 |
| Volume 19 | 526 | 519 | 146 | 1 |

EURALEX

The Dictionary of the Serbian Academy: from the Text to the Lexical Database