

# The CLaDA-BG Dictionary Creation System: Specifics and Perspectives

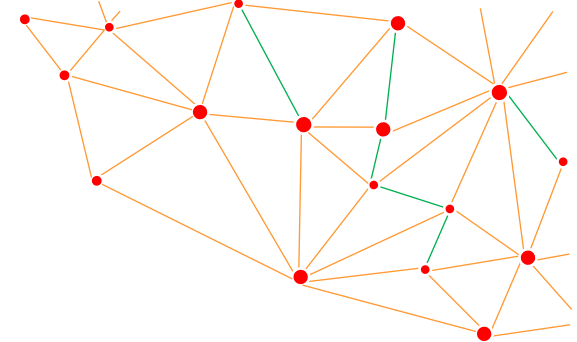
Zhivko Angelov, Kiril Simov, Petya Osenova, Zara Kancheva  
IICT, BAS

CLARIN Annual Conference  
10-12 September 2022

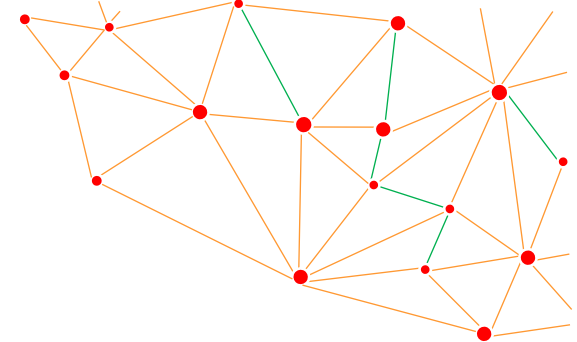


# Plan of the Talk

- Introduction
- Related Work
- System Specifics and Functionalities
- Supporting Language Resources
- Conclusions

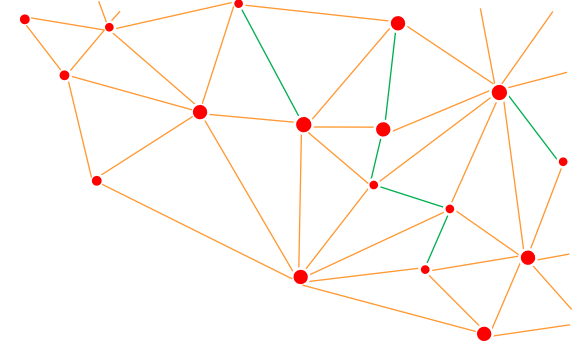


# Introduction: Aim

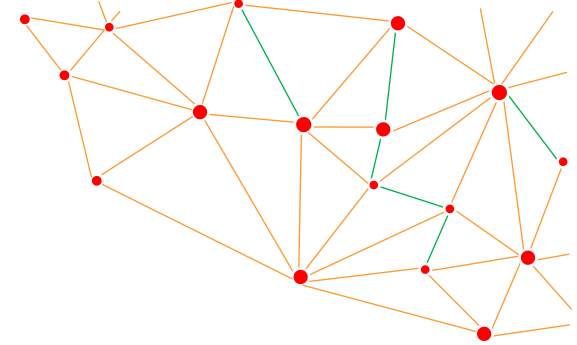


- We present the main principles and perspectives behind the CLaDA-BG Dictionary Creation System - [CLaDA-BG-Dict](#)
- The ultimate goal is:
  - to support the compilation of new dictionaries by individuals or collaborators with respect to a certain task and through the usage of all the available resources
  - within the Bulgarian National Interdisciplinary Research e-Infrastructure for Resources and Technologies in favor of the Bulgarian Language and Cultural Heritage, part of the EU infrastructures CLARIN and DARIAH - CLaDA-BG

# Introduction: Our Idea

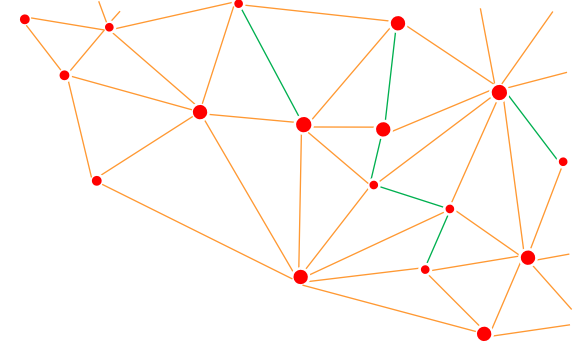


- At the heart of this system lies the **Bulgarian BulTreeBank WordNet (BTB-WN)**
- It has been developed as an aggregator of semantic knowledge around which other dictionaries and sources of information (including grammatical, encyclopedic, etc.) have been organized in the form of a(n) (inter)linked knowledge network



# Introduction: Our Motivation

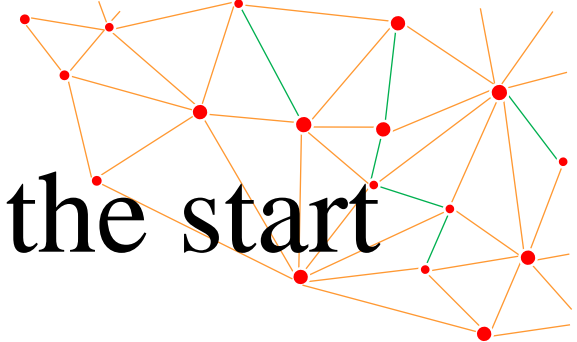
- Better control on the consistency in the creation of lexical language resources
- User friendly and communicative collaborative environment
- Better connections among the available resources
- We build on these best practices – systems for construction of other wordnets like GermaNet, Polish Wordnet and BulNet. However, we needed a system that reflects our rhythm of work.
- Functionality to provide links to grammatical paradigms, valences, links to Wikipedia, mappings to other wordnets, appropriate examples, etc.



# Related Work

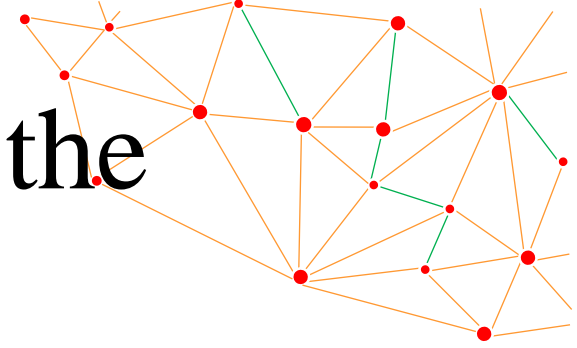
- We follow the approaches described in two existing wordnet editing systems: (Henrich and Hinrichs, 2010) and (Naskret et al., 2018)
- Many efforts have already been invested in dictionary creation systems:
  - ELEXIS - lexicographic workflow tools and crowdsourcing and gamification tools
  - DARIAH-ERIC – dictionary representation standard TEI-Lex0
  - CLARIN-ERIC - Lexica Resource Family overviews 89 lexica

# System Specifics and Functionalities: the start



- The existing version of the BTB-WN was initiated in an XML format within the CLaRK System
- However: the lexicographers had only a local view (without an underlying database) over the existing Bulgarian synsets
- Support was needed to mapping the Open English WordNet (OEW)
- Support was needed for the better integration of BTB-WN with other language and knowledge resources for Bulgarian

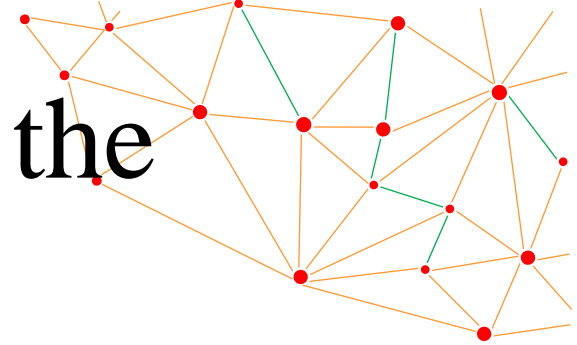
# System Specifics and Functionalities: the technicalities



- The system is a client-server web-based editor using a thick client model
- The database is installed on a server and it is accessed online via the web
- The data is stored in a relational database
- Work tracking system available
- Ticketing system available

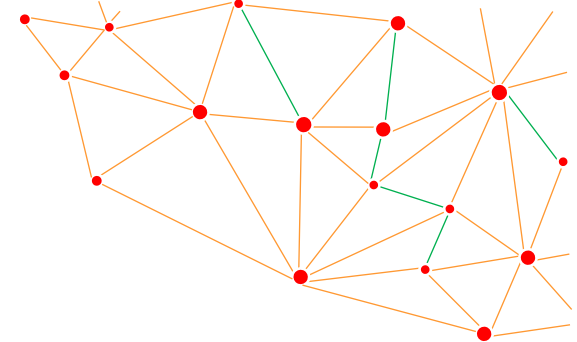


# System Specifics and Functionalities: the Utility



- The system provides information about a selected lemma:
  - its meanings (synsets) and associated examples;
  - its internal relations as well as the mappings to the OEW;
  - it also provides the ratio among the used relations
- In case of equivalent synsets between BTB-WN and OEW, the Bulgarian synset inherits all the relations from the English synsets
- The main form is literal-based (in other words, string-based)
- There is a possibility to consult a graphical representation, related information in the available dictionaries in the system

# User Interface



CLaDA-BG-Dict

Действия Форми Инструменти Помощ История

Списък от лемми

Начало на лемите Част на речта **Намерени са 120 лемми. (0.354)**

къ (и)

Категория от Релация към

без

Само лемите с отворени въпроси

Темата на въпроса да съдържа низа

Въпросът да съдържа низа

| №   | Лема            | Част на речта | Еквивалент | Тикет |
|-----|-----------------|---------------|------------|-------|
| 97  | късо съединение | n             | E          |       |
| 98  | късоврат        | a             |            |       |
| 99  | късовълнов      | a             |            |       |
| 100 | късоглед        | a             | E          |       |
| 101 | късоглед        | n             | E          |       |
| 102 | късоглед        | s             | E          |       |
| 103 | късогледство    | n             | e          |       |
| 104 | късокрак        | a             |            |       |
| 105 | късопаметен     | a             |            |       |
| 106 | късопръст       | a             |            |       |
| 107 | късче           | n             | E          |       |
| 108 | кът             | n             | E          |       |
| 109 | кътче           | n             | E          |       |
| 110 | къч             | n             | E          |       |
| 111 | къшей           | n             |            |       |
| 112 | къшла           | n             | E          |       |
| 113 | къшлак          | n             | E          |       |
| 114 | къща            | n             | E          | 1     |
| 115 | къщен           | a             | E          |       |
| 116 | къщи            | r             | E          |       |
| 117 | къщица          | n             | E          |       |
| 118 | къщичка         | n             | E          |       |
| 119 | къщурка         | n             | E          |       |
| 120 | къщя            | n             | E          |       |

Лема: къща

Синонимно гнездо

| Част | Категория     | Дефиниция  |
|------|---------------|--|
| n    | noun.artifact | Вид сграда, жилище, дом на един или повече етажи, в който живеят постоянно или временно хора от един или повече семейства. |
| n    | noun.location | Жилището, в което някой живее.   |

Лексикална единица

3 примера, 3 от които към лема

| Лема    | Пример   | # Примери | Идентификация |
|---------|--|-----------|---------------|
| къща    | @@@ Къщите @@@ имат сравнително ниска етажност | 3         | 9686          |
| къщица  |  | 0         | 199851        |
| къщичка |  | 0         | 199852        |
| къщурка |  | 0         | 9813          |

Примери

Лема

Концептуални релации

Надлонятия / подлонятия

Допълнителна информация

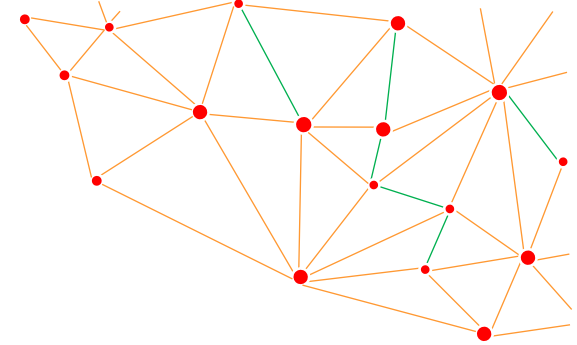
Проблеми / въпроси

Отворен въпрос

Временни бележки

Еквивалентните

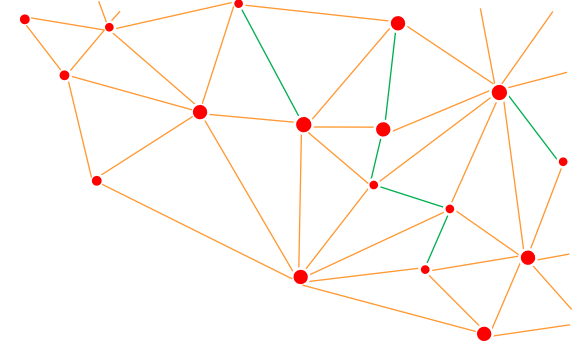
A conceptual relations diagram for the word 'къща' (house). The diagram shows a central node 'къща, къщ' connected to various related concepts. The nodes are arranged in a grid-like structure with red arrows indicating relationships. The nodes include: 'къща, къщ', 'жилище, об...', 'living ass...', 'конструк...', 'артефакт', 'цялост, ця...', 'вещ, пред...', 'физическа...', 'свъщност, ...', 'house', 'home, dwell', 'constructio...', 'artifact, ar...', 'whole, unl...', 'object, phy...', 'physical en...', 'entity', 'edifice, bu...', and 'строеж, зд...'.



# Supporting Language Resources

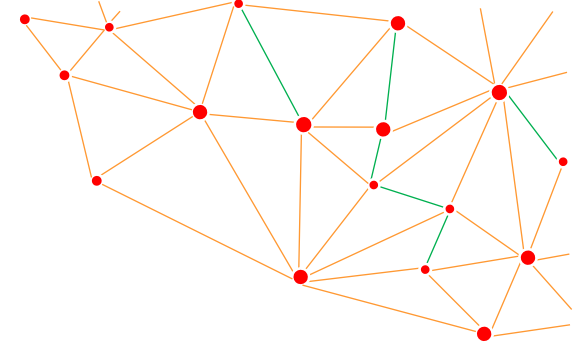
- In our view the necessary minimum of functionalities of such a system would include:
  - an editor of lexical entries supporting different structures of interrelated elements
  - access to existing dictionaries and corpora
  - concordances and other materials
- At the moment the integration of BTB-WN has been made with OEW, but a coverage comparison and linking between corresponding senses is planned to be done with several dictionaries such as
  - Bulgarian Explanatory Dictionary
  - Bulgarian Inflectional Lexicon (currently partially done)
- The system supports mapping to Wikipedia via the inclusion of Wikipedia article URI to the corresponding synset

# Conclusions



- **CLaDA-BG-Dict** is an editor, which could be used both for creating lexical databases like wordnets as well as traditional types of dictionaries
- It provides possibilities of linking the available data in many ways depending on the goal
- **CLaDA-BG-Dict** has already been successfully used for editing of more than 19 000 synsets that were created at earlier stages in an XML format, and for the addition of around 13 000 synsets together with appropriate examples

# Future Plans



- Our vision for future is to enhance replicability and re-usage of dictionary compilation as much as possible.
- In its current beta-version the system uses its own format for uploading corpora and other digitally-born or digitized dictionaries.
- However, it is planned to conform to the common standards such as TEI, TEI LEX0, Lemon, etc.
- All the participating resources will be made available through the CLaDA-BG repository and dedicated web services.