# Integrating Research & Research Infrastructures into Teaching

Marko Simonović (Uni Graz)
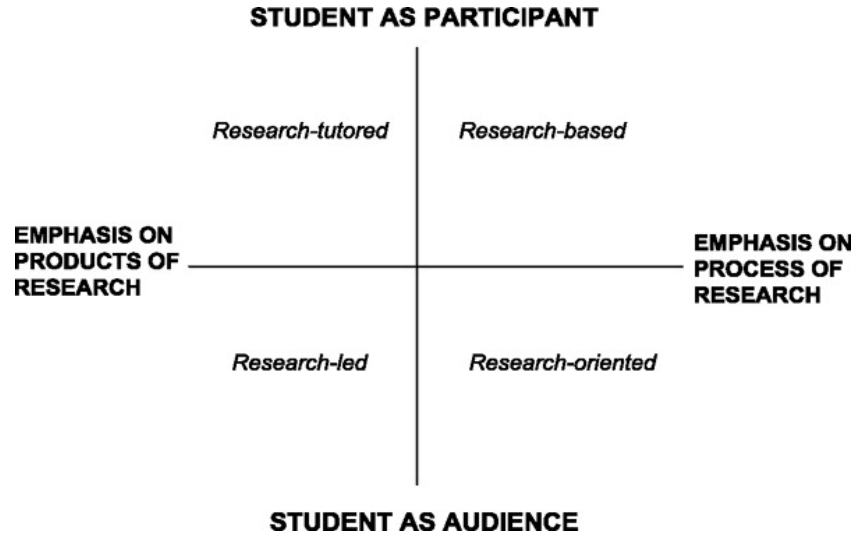Iulianna van der Lek (CLARIN ERIC)

Graz Multiplier Event, 5 July 2022

# Overview

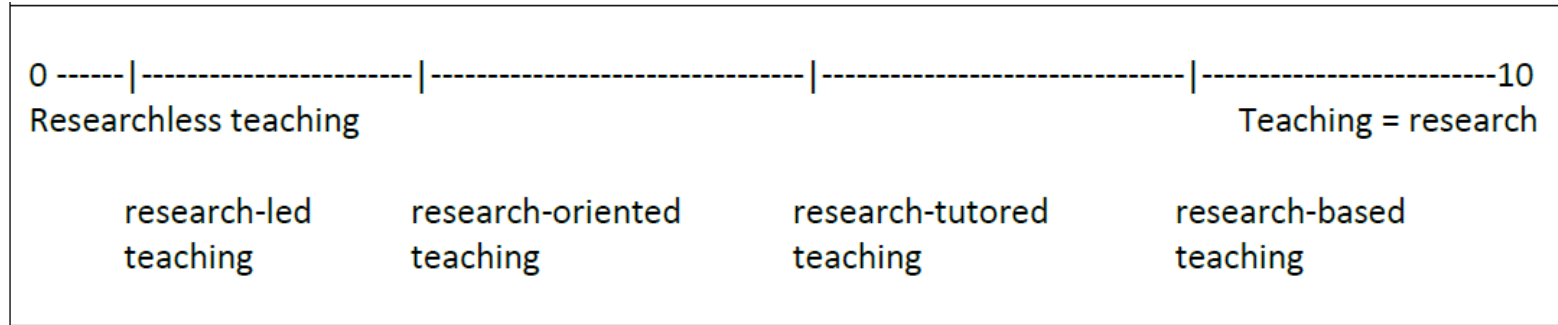Research-Based Teaching: Guidelines and Best Practices

Integrating Research Infrastructures into Teaching: The Case of CLARIN

# What is research-based teaching?



**STUDENT AS PARTICIPANT**

*Research-tutored*      *Research-based*

**EMPHASIS ON PRODUCTS OF RESEARCH**      **EMPHASIS ON PROCESS OF RESEARCH**

*Research-led*      *Research-oriented*

**STUDENT AS AUDIENCE**

Healey's (2005, Jenkins et. al 2007) adapted model of the research-teaching nexus (from Visser-Wijnveen et al. 2010).

# What is research-based teaching?

0 ------|-----------------------|--------------------------------|---------------------------------|------------------------10
Researchless teaching                                                                                    Teaching = research

    research-led          research-oriented          research-tutored          research-based
    teaching              teaching                   teaching                  teaching

Dekker & Wolf (2016)'s scale of the research-teaching nexus.

# Why research-based teaching?

Advantages for for students

- direct **preparation for a research career**, but also
- stimulate **curiosity**
- practice **problem solving**
- practice **forward thinking**
- gain **independence**
- gain **experience**

Advantages for lecturers

- integrate **research and teaching**
- **new insights** about own research
- **more interaction** with students
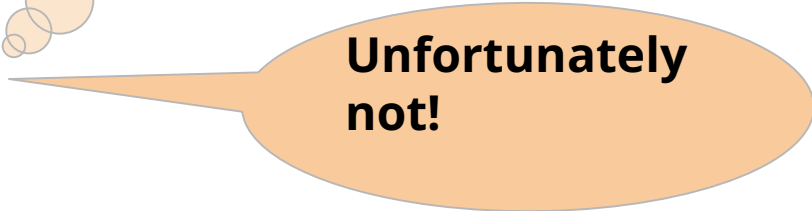
# Why research-based teaching?

Advantages for for students

- direct **preparation for a**
- stimulate **curiosity**
- practice **problem s**
- practice **forward th**
- gain **independence**
- gain **experience**

**Isn't this what university lecturers do all the time?**

Advantages for lecturers

- integrate **research and teaching**
- **new insights** about own research
- **more interaction** with students

# Why research-based teaching?

Advantages for for students

- direct **preparation for a**
- stimulate **curiosity**
- practice **problem s**
- practice **forward th**
- gain **independence**
- gain **experience**

Advantages for lecturers

- integrate **research and teaching**
- **new insights** about own research
- **more interaction** with students

**Isn't this what university lecturers do all the time?**

**Unfortunately not!**

# Research-teaching nexus

**Teaching** and **research** are typically planned, performed and evaluated separately.

The connection between them, usually discussed in the literature under the rubric of the **research-teaching nexus**, receives little attention in the reality of most academic institutions (for a notable exception, and a rich collection of experiences and ideas, see Bastiaens, van Tilburg & van Merriënboer 2017).

Visser-Wijnveen (2009: 141): "academics' conceptions of the research-teaching nexus are related to their conceptions of teaching and **not to their conceptions of research and knowledge**".

# Why we don't do more research-based teaching?

- The results of first research-based courses may be viewed as unsatisfactory by the lecturer.
  - Course design turns out to be time consuming.
  - Teaching on own research turns out to be more challenging than originally assumed.

Still, as transpired from the survey we held within the preliminary needs analysis of our project, a bit over 50% of lecturers indicated that they integrate research into teaching, while around 75% indicated that they would be happy to follow dedicated training focusing on this.

# Filling the gap: Our guidelines and best practices

Co-funded by the
Erasmus+ Programme
of the European Union

SKILLS

**Research-based teaching: Guidelines and Best Practices**

*UPSKILLS Intellectual output 2.1*

Compiled by:

Marko Simonović*, Boban Arsenijević*, Stavros Assimakopoulos†, Darja
Fišer*, Iulianna van der Lek*, Maja Miličević Petrović‡, Lonneke van der Plas†,
Margherita Pallottino*, Genoveva Puskas*, Tanja Samardžić*

* University of Graz
† University of Malta
*CLARIN ERIC
‡ University of Bologna
*University of Geneva
*University of Zurich

# Guidelines for integrating research into teaching

Featuring:
- Instructions for choosing/developing a course subject
- Detailed pick-and-choose list of research-related topics
- Detailed pick-and-choose list of research-related learning outcomes
- Making instructions for students
- Organising and supervising the work
- Evaluation and grading
- After the course

- Course description template
- Course Examples
  - Current trends in Phonology
  - Multilingualism
  - Deverbal nominalisation in West South Slavic
- Annexes
  - The structure of a research report
  - Assessment survey

# Integrating Research Infrastructures into Teaching

Iulianna van der Lek

The Case of CLARIN

# The Language Resource Life Cycle

**Data Sharing**

- Repositories
- Catalogues
- Citation
- Legal Issue**s**

**Data Acquisition and Collection**

- Research Infrastructures
- Repositories
- Data Catalogues
- IPR & Legal Issues
- Citation

**Data Archiving**

- Deposit LRs in a institutional or domain-specific FAIR repository
- Describe LRs using metadata standards
- IPR and legal issues for handling sensitive data; licenses

**Data Curation and Annotation**

- Standards and formats
- Annotation: (semi) automatic, manual

**Linguistic Analysis and Research**

- Querying metadata, LRs and their annotations
- Analysis & visualisation; combining data from different sources -> exchange standards

The Language Resource Life Cycle

# A Guide and a Course for Teachers

**Course:** Introduction to Language Data: Standards & Repositories
- 4-5 ECTS

**Unit 1:** Intro to Language Resources and Infrastructures for Language Resources and Technologies
**Unit 2:** Finding, accessing and using LRs
**Unit 3:** Tools for Linguistic Processing
**Unit 4:** Archiving and Sharing LRs

**Available for Piloting in September 2022**

# Finding Language Resources

**CLARIN Virtual Language Observatory**

- Metadata search in the CLARIN language repositories across Europe

- Does not contain data, only refers to it using PIDs

- Useful as a first step to identify which language resources are available

# Finding Language Resources

**Content Search**

- To search for specific patterns across corpora at the same time
- The search results from a specific corpus can be downloaded in various formats and processed in other tools for further analysis
- Possible to retrieve the citation information
- Useful as a first step to discover where interesting lg resources are hosted

# Finding Language Resources

**Discipline-specific repositories**

- Access to curated LRs
- Datasets can be cited via PIDs
- E.g. CLARIN B-Centres, The Open Language Archive, Meta-Share



https://arche.acdh.oeaw.ac.at/browser/

# Find and Query Large Collections of Corpora

CLARIN Resource Families

- User-friendly overviews of well-curated corpora and tools
- Download or query with concordancers, e.g. Korp, Corpuscle and KonText

**Corpora**
- Computer-mediated communication corpora
- Corpora of academic texts
- Historical corpora
- L2 learner corpora
- Literary corpora
- Manually annotated corpora
- Multimodal corpora
- Newspaper corpora
- Parallel corpora
- Parliamentary corpora
- Reference corpora
- Spoken corpora

**Lexical Resources**
- Lexica
- Dictionaries
- Conceptual Resources
- Glossaries
- Wordlists

**Tools**
- Normalization
- Named entity recognition
- Part-of-speech tagging and lemmatization
- Tools for sentiment analysis

# Find and Query Large Collections of Corpora

**Croatian Twitter training corpus ReLDI-NormTagNER-hr 2.0**

Croatian

This corpus contains Tweets. The corpus is morphosyntactically tagged following the MULTEXT-East Version 4 tagset.

The corpus is available through the concordancers KonText and noSketchEngine and for download from the CLARIN.SI repository.

For the relevant publication, see Miličević and Ljubešić (2016)

KonText

noSketch

Download

**Size:** 89,000 tokens

**Annotation:** tokenisation, sentence segmentation, word normalisation, morphosyntactic tagging, lemmatisation and Named Entity recognition

**Licence:** CC BY 4.0

Query the ReLDI corpus via NoSketch engine on CLARIN.SI

---

🔍 🗄 ReLDI-hr (manually tagged Croatian tweets)

My jobs

User guide 🔗

**ReLDI-hr (manually tagged Croatian tweets)** ❓

Croatian tweets with manualy normalised (standardised), morphosyntactically tagged and lemmatised words and named entities ReLDI-hr v2.1

| Counts | | General info | | Lexicon sizes | | Tags legend | |
|---|---|---|---|---|---|---|---|
| Tokens | 89,104 | Corpus description | Document | word | 27,289 | Noun | N.* |
| Words | 71,768 | Language | Croatian | norm ❓ | 25,395 | Noun proper | Np.* |
| Sentences | 7,939 | Encoding | UTF-8 | lempos | 17,270 | Noun common | Nc.* |
| Documents | 3,871 | Compiled | 09/11/2019 15:41:23 | tag ❓ | 694 | Verb | V.* |
| | | Tagset | Description | ud_pos ❓ | 90 | Adjective | A.* |
| | | | | ud_feats ❓ | 764 | Pronoun | P.* |
| | | | | diff ❓ | 5 | Adverb | R.* |
| | | | | lc ❓ | 25,219 | Preposition | S.* |
| | | | | lemma | 16,675 | Conjunction | C.* |
| | | | | lemma_lc | 16,020 | | |

**Structures and attributes**

| text 3,871 | ≫ |
|---|---|
| name 5,883 | ≫ |

# Collecting and Citing Language Resources

## Virtual Collection Registry

- Collect the datasets discovered in the VLO or other data catalogues into a Virtual Collection
- Share the collection with the others
- Cite the collection
- Process the collection with Switchboard tools

# Processing Text Collections

## The Language Resource Switchboard

- A service that allows you to find a matching tool to analyse plain text files
- Taggers, lemmatizers, named entity recognizers, chunking tools etc.

# Processing Text Collections

Example of morphological analysis in WebLicht

# Archiving Language Resources

- Institutional repository
- Domain-specific repository
- Infrastructural repository: Depositing Services | CLARIN ERIC
- Support with data preparation, license selection, depositing process

Item submission

① Basic Info — ② Who's involved — ③ Describe — ④ Upload — ⑤ License — ⑥ Note — ⑦ Review — ⑧ Complete

Example of  data submission lifecycle in CLARIN.SI

# Sharing Language Resources

- Select appropriate licenses
- The repository assigns unique identifiers
- Others can cite the resources



Automatically stress labelled morphological lexicon Sloleks 1.2, version 1.1

BIBTEX  CMDI

Please use the following text to cite this item or export to a predefined format:

Krsnik, Luka; Robnik-Šikonja, Marko; Šef, Tomaž and Krek, Simon, 2018, *Automatically stress labelled morphological lexicon Sloleks 1.2, version 1.1*, Slovenian language resource repository CLARIN.SI, ISSN 2820-4042, http://hdl.handle.net/11356/1186.

Share:

# Guidance on the Use of Standards and Formats

[CLARIN Standards Information System](#)

- Data deposition formats
- Language-technology related standards

# Guidance on Licenses and Legal Issues in Data Reuse

On the CLARIN Legal Information Platform, you will find useful information about:

- [Introduction to Copyright and Related Rights](#)
- [Licensing Practice](#)
- [Personal Data Protection](#)

# Questions?

**Email**: iulianna@clarin.eu