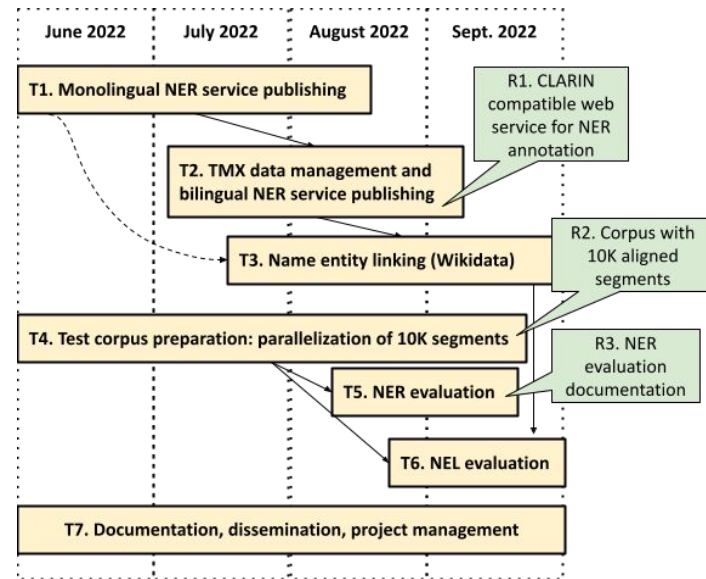# It-Sr-NER: CLARIN compatible NER and geoparsing web services for parallel texts: case study Italian and Serbian

*Olja Perišić & JeRTeh team*
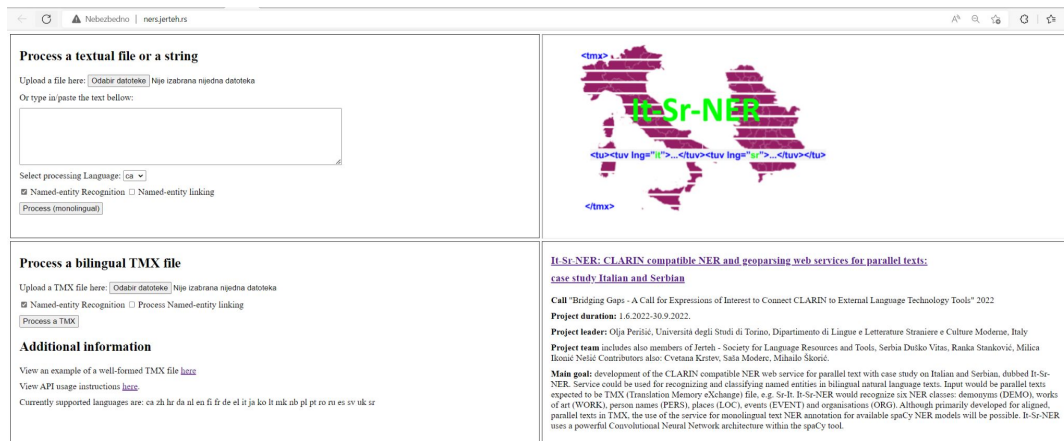*20.10.2022, Jerteh seminar*

http://jerteh.rs/

# Project intro

- Project title: *It-Sr-NER: CLARIN compatible NER and geoparsing web services for parallel texts: case study Italian and Serbian*

- grantNo: CE-2022-2070, funding type: EU - CLARIN Bridging Gaps project.

- Project URL: https://github.com/rankastankovic/It-Sr-NER/

- Demo data URL: https://github.com/rankastankovic/It-Sr-NER/tree/main/corpus

- Service URL: http://ners.jerteh.rs/

- Project tasks T1-T7, Project outcomes: R1-R3

- **Project leader: *Olja Perišić,*** Università degli Studi di Torino, Dipartimento di Lingue e Letterature Straniere e Culture Moderne, Italy

- **Project team** Duško Vitas, Ranka Stanković, Milica Ikonić, Cvetana Krstev, Saša Moderc, Mihailo Škorić.

# T1. Monolingual NER service publishing

- It-Sr-NER-ws web service is a CLARIN-compatible web service accessible via a convenient web interface, with options to either upload monolingual, text files with output containing NE annotations. Main web service site is http://ners.jerteh.rs/

- The user can input text or upload the text file in one of the offered languages.

- If the short text is pasted into the corresponding text area, the application will show the input text containing named entity annotations. In the case of a text file, the same text file but containing NE-annotated text will be returned.

- The annotation is possible for languages with spaCy NER models on the backhand.

- For web application and service development, Python programming language was used. The guidelines published on https://weblicht.sfs.uni-tuebingen.de/weblic htwiki/index.php/Developer_Manual will be followed.

- It-Sr-NER-ws is hosted at JERTEH infrastructure and integrated into the CLARIN ERIC Switchboard, with consultation with the central CLARIN developer's office.

# T2. TMX management and bilingual NER services publishing

It-Sr-NER-ws function for processing TMX files as input, then providing appropriate named entity annotations, retaining the same TMX format.

Here is an example translation unit (<tu>) with annotated translation unit variants (<tuv>), where each aligned segment (<seg>) is annotated using the appropriate NER model

(e.g. Italian or Serbian)

```
<tu>
    <tuv xml:lang="it" creationid="n12 " creationdate="20211202T203105Z">
        <seg><PERS>Phileas  Fogg</PERS>,  <DEMO>inglese</DEMO>  certamente,  non  era,  forse
<DEMO>londinese</DEMO>. </seg>
    </tuv>
    <tuv xml:lang="sr" creationid="n12 " creationdate="20211202T203105Z">
        <seg>Iako je na prvi pogled bio <DEMO>Englez</DEMO>, <PERS>Fileas Fog</PERS> verovatno
nije bio <DEMO>Londonac</DEMO>.</seg>
    </tuv>
</tu>
```

# T3. Named entity disambiguation and linking  (NEL) (geolocation and other entity types) using Wikidata

NEL is based on spacyopentapioca

For example the person name <PERS>Phileas Fogg</PERS> found in the Italian sentence, and <PERS>Fileas Fog</PERS> in the Serbian sentence, would be both mapped to Wikidata Q2587533 that can be useful for many applications.

Namely, these NEs are co-indexed to refer to the same person, providing named-entity identification by referring to the ID in a common knowledge base Wikidata and using common QID= Q2587533 for the entry for Phileas Fogg on Wikidata https://www.wikidata.org/wiki/Q2587533.

Service provide two options:

```
<tu>
    <tuv xml:lang="it" creationid="n12 " creationdate="20211202T203105Z">
    <seg><PERS ref="https://www.wikidata.org/wiki/Q2587533"  desc="character created by Jules
Verne">Phileas Fogg</PERS>, <DEMO>inglese</DEMO> certamente, non era, forse,
<DEMO>londinese</DEMO>. </seg>
    </tuv>
    <tuv xml:lang="sr" creationid="n12 " creationdate="20211202T203105Z">
        <seg>Iako je na prvi pogled bio <DEMO>Englez</DEMO>, <PERS
ref="https://www.wikidata.org/wiki/Q2587533"  desc="character created by Jule
Fog</PERS> verovatno nije bio <DEMO>Londonac</DEMO>.</seg>
    </tuv>
 </tu>
```

1. NEL relies on entities that are previously annotated with NER tags, and linking is performed just on annotated segments

```
<tu>
    <tuv xml:lang="it" creationid="n12 " creationdate="20211202T203105Z">
    <seg><WDT ref="https://www.wikidata.org/wiki/Q2587533" label="PERSON" desc="character
created by Jules Verne">Phileas Fogg</WDT>, inglese certamente, non era, forse,
<DEMO>londinese</DEMO>. </seg>
    </tuv>
    <tuv xml:lang="sr" creationid="n12 " creationdate="20211202T203105Z">
        <seg>Iako je na prvi pogled bio Englez, <WDT ref="https://www.wikidata.org/wiki/Q2587533"
label="PERSON" desc="character created by Jules Verne">Fileas Fog</WDT> verovatno nije bio
Londonac.</seg>
    </tuv>
</tu>
```

2. NEL is applied on input text without NER annotation, using only opentapioca annotation

# T4. Test corpus preparation: parallelization of 10K segments

- Corpus consists of 10000 aligned segments (sentences) from several novels.
- Novels are represented with samples in which segments are shuffled in order to avoid copyright problems.
- The novels are aligned and prepared following the previously used pipeline for parallelization (Utvić et al. 2008; Krstev and Vitas 2011).
- The corpus is published at the ILC4CLARIN B Centre, so that it will be visible via the VLO (Virtual Language Observatory) https://vlo.clarin.eu/.
- Corpus is also accessible from working github location  https://github.com/rankastankovic/It-Sr-NER/tree/main/corpus

**Novels**:
- Umberto Eco: The Name of the Rose;
- Carlo Collodi: The Adventures of Pinocchio;
- Elena Ferrante: Those Who Leave and Those Who Stay,
- Luigi Pirandello: One, No one and One Hundred Thousand,
- Jules Verne: Around the World in Eighty Days;
- Ivo Andrić: Legends of Anika and The Bridge on the Drina,
- Borisav Stanković Impure Blood,
- Branislav Nušić: Municipal child: a novel of an infant.

| Italian (it) | Serbian (sr) |
| --- | --- |
| **n1** Durante una vacanza a Ischia si innamora di Nino Sarratore, per il quale lascia il marito. | **n1** Tokom jednog letovanja na Iskiji zaljubljuje se u Nina Saratorea, zbog koga ostavlja muža. |
| **n2** Insieme a suo padre Fernando, e grazie a Lila e al denaro di Stefano Carracci, mette su il calzaturificio Cerullo. | **n2** Zajedno sa ocem Fernandom, i zahvaljujući Lili i novcu Stefana Karačija, otvara obućarsku radnju „Čerulo". |
| **n3** Dopo le elementari Elena continua a studiare con crescente successo; | **n3** Nakon osnovne škole Elena nastavlja da se školuje s rastućim uspehom; |
| **n4** Lavora nella salumeria di famiglia prima, e nel negozio di scarpe poi. | **n4** Isprva radi u porodičnoj delikatesnoj radnji, a zatim u obućarskoj radnji. |
| **n5** È fidanzato con Marisa Sarratore e diventa il responsabile del negozio di scarpe di piazza dei Martiri. | **n5** Verio se s Marizom Saratore i postaje odgovoran za obućarsku radnju na Trgu mučenika. |

# T4. Test corpus preparation: parallelization of 10K segments

Apart for download, corpus is available also for search on the digital library Bibliša for search with possibilities of morphological and semantic query expansion for Serbian.

The option for browsing of aligned collections is available with a limited number of segments for all and with a larger number of segments for authorized users.

# T5. NER evaluation

Automatically annotated 1000 sentence pairs imported to the INCEPTION tool http://inception.jerteh.rs for manual correction

**Visual comparison of automatic and manual annotation**

Manual annotation in blue and automatic in orange
For each name entity type one HTML file per language is generated
and published as Italian set  and Serbian set .

⚙Language Resource Switchboard    Upload    Tool Inventory    Help                    CLARIN

# Resources

| Vern-test-it-sr-TMX.tmx.xml  94.18 KiB  🗑 | Mediatype | Language |
|---|---|---|
| | application/xml        ∨ | Select language        ∨ |

**➕ Add another resource**

# Matching Tools                          Group by task ☑  Search for tool

∨ Named Entity Recognition

    ⟩  spaCy bilingual NEL (for TMX)  **Open** 🔗

    ⟩  spaCy bilingual NER (for TMX)  **Open** 🔗

    ⟩  spaCy bilingual NER and Geoparsing (for TMX)  **Open** 🔗

    ⟩  spaCy binlingual NER and NEL (for TMX)  **Open** 🔗

# Summary of the project results

R1. CLARIN compatible web service for NER annotation of 1) bilingual texts tested on it-sr with 6 NE classes 2) monolingual texts with supporting documentation (open source), tested on Serbian and Italian.

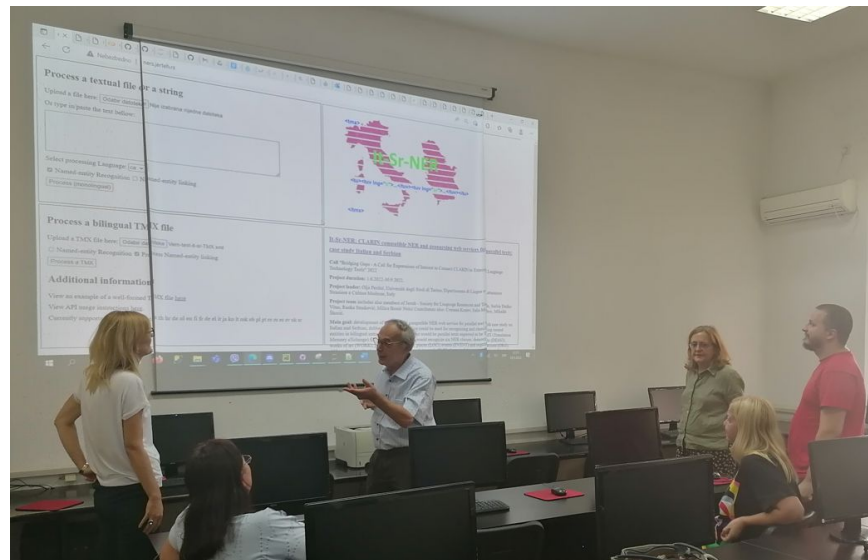R2. Corpus with 10000 aligned segments to be used for testing (open data)

Automatically linked  LOC-related entities to the corresponding pages on Wikidata, followed by an analysis of the NEL success on a validation subset of parallel sentences.

R3. NER evaluation documentation that will include manually corrected 1000 segments with metrics and visualization of differences.

We will conduct an exhaustive evaluation on a subset of the parallel sentences. These results will be published and would enable users to gain insight into the model's performance and to understand on what kind of input these models yield the best results.

# Belgrade workshop

The meeting in Belgrade was held in September 2022 for final discussions, web service testing and evaluation of results. Apart from quantitative analysis, qualitative analysis was discussed and possibilities for improvement of the result of automatic annotation.

# The CLARIN Bazaar 2022

# Discussion and future plans

- POS tagging and lemmatisation of TMX format enabling CQL bilingual queries over NER annotated text: the corpus will be published on https://noske.jerteh.rs/ .
- Textometry using TXM tool will be used for the analysis of entities on both sides (it-sr).
- The bilingual corpus augmentation (set of novels are already collected for preparation and alignment).
- Further research will explore new technologies for NER and NEL: training of a new model for Serbian that would include training set augmentation and use of gazettiers.
- The students of Italian in University of Belgrade as well as students of Serbian in Torino, will benefit and use developed resources in future teaching.
  - Results are open and so available for other students and researchers as well.
  - Parallel corpora are central to translation studies and contrastive linguistics and easy-to-use concordancers considerably facilitates the study of interlinguistic phenomena.
  - New types of exercises could be produced: automatically generated tests where sentences can have a lema form so that they make a grammatically correct sentence.
  - Also, the analysis of multiple translations will enlarge the interpretive process and perspectives that students draw from the text. The parallel corpus will
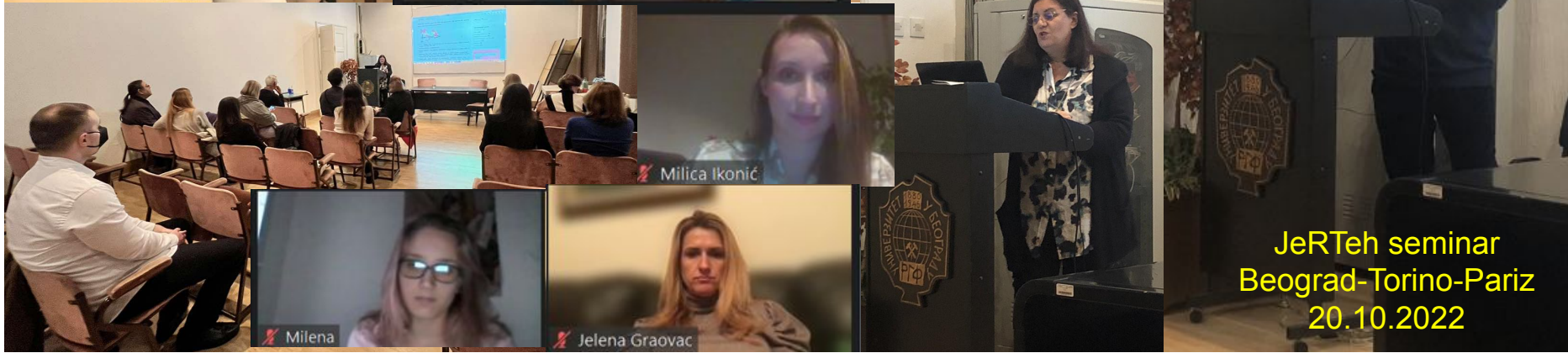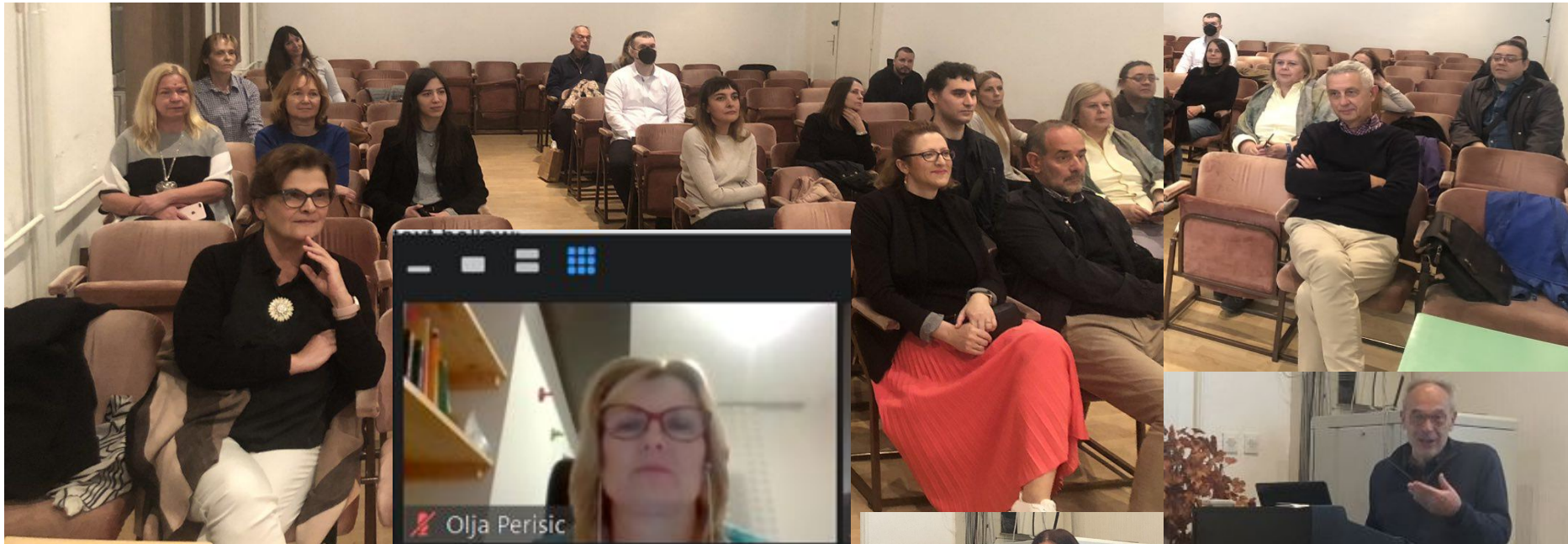- The team will look for possibilities for the new projects and frameworks to cooperate.

Olja Perisic

Milica Ikonić

Milena

Jelena Graovac

JeRTeh seminar
Beograd-Torino-Pariz
20.10.2022