Named Entities' Journey from UniteX to Wikidata



Ranka Stanković, University of Belgrade

Journées « Unitex », 15-16 Décembre 2021

Key topics

- Named Entities: recognition, annotation, extraction, linking with Wikidata
- Speeding up the process of data preparation and entry
- Tools to support the process:
- Use case:



- Serbian ELTeC collection of novels from period 1840-1920
- Wikidata data model and examples
- SPARQL queries of SrpELTeC Wikidata
- NEL (with Wikidata) in INCEpTION

The Serbian ELTeC collection of novels from period 1840-1920

- The COST action CA16204 began in 2017 and will end at the beginning of 2022.
- One of the most important goals of this action is the preparation of a multilingual corpus (called European Literary Text Collection ELTeC) <u>https://www.distant-reading.net/eltec/</u>
- When fully completed, ELTeC will contain 100 novels first published in the period 1840-1920. for a large number of European languages
 - {cze, deu, eng, fra, hun, pol por, rom, slv, srp} subcollections with 100+ novels
 - 7 subcollections are under development for level-1
- SrpELTeC collection contains
 - 100 novels that are part of ELTeC and
 - 13 novels from <u>SrpELTeC-ext</u> extended collection
- TEI-conformant ODD specification https://github.com/COST-ELTeC/Schemas

Distant [] Reading

SrpELTeC text formats: TEI - level 1

Text preparation

- <teiHeader> with rich information about author. 0 novel, publication and digital version
- In <text> structural elements: <div>. <head>. 0 <pb>,<milestone/>, <quote>, <l>, <foreign>, <ref>, <note>. <hi>....
- more in Special Issue of Infotheca "The Serbian 0 Part of the ELTeC Collection through the Magnifying Glass of Metadata" (Krstev, 2021)
- As TELXML: https://github.com/COST-ELTeC/ELTeC-srp/tree/master/level1
- As HTMI : https://distantreading.github.io/ELTeC/srp/index.html
- More on level-1 schema: https://distantreading.github.io/Schema/eltec-1.html

quote>

<1>Мито море, Митанче,</1> <1>Цркни ми, пукни, душманче </1> <1>Мито море, Мито, Мито!</1> </guote>

<р>И то "Мито" све прихватише и почеше са толико чежње, страсти, лудила да певају и изговарају, да и сама амамџика, раздрагана, смешећи се, мораде да оде, још јаче да притвори спољна амамска врата, да се тамо, у чаршији, не би чуло и око амама почели да скупљају мушки. </р>

<pb n="158"/>

<р>Симка већ беше обучена и села у чело астала. Испред ње је већ стајало стакло ракије, кисела грожћа, крушака и паприка. Око ње почеше и остале да се ређају, освежавајући се воћем и осталим слаткишима.

И Софка је грабила, да би била брзо готова. Сиђе и она, седе међ њих. Почеше и њу да служе, и, на изненађење свих, не само што поче јести, него испи и неколико чаша ракије, а једнако са оним задржаним осмехом на устима, с очима мало стиснутим и обрвама једва приметно набраним, као смејући се свима и све их сажаљевајући...

</div>

Kdiv type="chapter" xml:id="SRP19101 C16"> <head>**XVI**</head>

<р>У полумрак вратиле се из амама. И све док су пролазиле чаршијом, прелазиле онај мост, увијене, погнуте и кријући се једна уз другу,

ELTeC level-2 specification

- ELTeC level-2 is level-1 augmented by: lemmatisation, POS-tagging (UD tagset), NER annotation, morphosyntactic annotation optional, direct speech and sentiment analyses envisaged
- ELTeC level-2 includes all elements existing in level-1 and introduces some new ones:
 - <s> as the sentence tag, used for segmentation of text into sentences, and
 - <w> and <pc> for tokenization
 - mandatory linguistic attributes @pos, @lemma, and @join,
 - optional @xml:id, @msd
 - <rs> (referring string) for NER with @type attribute: **PERS**, **LOC**, **ORG**, **DEMO**, **WORKS**, **EVENTS**, **ROLE**
- The element <s> can contain a sequence of <w> elements, either directly or in the sub-paragraph elements <corr>, <emph>, <foreign>, <hi>, <label>, <title>.
- More about level-2: <u>https://distantreading.github.io/Schema/eltec-2.html</u>

ELTeC level-2 production

- Two main problems encountered in producing Serbian ELTeC level-2 were similar to those encountered for other languages:
 - 1) majority of morphosyntactic taggers do not work well with XML format and
 - 2) harmonization of NER and morphosyntactic annotations, which are performed separately with different tools.
- A solution was found in the TXM tool (<u>https://txm.gitpages.huma-num.fr/textometrie/</u>)



- an environment that enables tagging of XML files,
- solved the problem of alignment of NER and morphosyntactic tags.

The Serbian SrpELTeC Level-2 pipeline



SrpNER in Unitex

- Sentence splitting performed by Unitex transducer (Krstev 2008; Paumier 2021)
 - adapted to consider level-1 tags
 - outputs the start tag <s> at the beginning of a sentence and the end tag </s> at its end
- NER performed by the rule-based system SrpNER (Krstev et al. 2014), based on large-scale lexical resources for Serbian (Krstev, 2008), coupled with local grammars in the form of finite-state transducers (Vitas and Krstev, 2012).
- Since SrpNER works on Latin texts, it is necessary to transliterate Cyrilic texts to Latin.
- SrpNER recognises 11 classes of NEs:
 - o dates and time (moments and periods),
 - money and measurement expressions,
 - geopolitical names (countries, settlements, oronyms and hydronyms), and
 - personal names (last, first, nicknames, full, ...).
- Semi-automatic removal of the embedded NER tags from the SrpNER output

1.	<pre><role>carice</role></pre>
	<persname.full>Marije</persname.full>
-	Terezije
2.	<pre><role>Sekretar</role></pre>
	<persname.last>Živanović</persname.last>
3.	<org>Saborne crkve u</org>
	<top.gr>Beogradu</top.gr>
4.	<org>manastira <pers.spec>Sv.</pers.spec></org>
	Marka
5.	<pre><pers>Ali-<pers><persname.last>Milevu</persname.last></pers></pers></pre>
:	.last>
÷	
·····	↓
1.	<role>carice</role> <persname.full>Marije</persname.full>
-	Terezije
2.	<role>Sekretar</role>
-	<pre><persname.last>Zivanovic</persname.last></pre>
3.	<org>Saborne crkve u Beogradu</org>
4.	<org>manastira Sv. Marka</org>
5.	<pers>Ali-Milevu</pers>
- /	I Frank

The tagsets mapping in the NER&Beyond portal

http://nerbeyond.jerteh.rs/

- Level-2 tagset: PERS, ROLE, LOC, ORG, DEMO, EVENT, WORK
- An automatic procedure implemented as part of the NER&Beyond portal (Stanković et al. 2019; Šandrih Todorović et al. 2021) was developed and used to map SrpNER tags into the 7-categories ELTeC NER schema.

Distant [] Reading

Annotation tagset mapping and transliteration

Description of a procedure for harmonisation of two annotation tag sets:

1. User uploads a set of source annotated documents with any tagset with one or more document with gold (target) annotation tagset

2. System reads both tagsets and offer mappings (for each tag, some option should be selected)

3. System generates source documents in a new annotation scheme

4. User downloads the results

Upload a .zip file with following structure:

f.zip ______gold/ ______md1.ann ______mdm2.ann ______x1.ann ______x1.ann ______x2.xml ______x3.txt

Download Sample Archive

	DEMO	EVENT	LOC	ORG	PERS	ROLE	WORK	ignore	rem
back	O	0	0	ORO	0	OLL	O	®	
hody	õ	0	õ	0	õ	0	õ		0
demonum		õ	õ	õ	õ	0	õ	õ	ő
div	0	0	0	õ	0	0	0	۲	ē
event				0			õ	0	0
foreign		0	0	0	0		0	۲	ć
front	õ	õ	0	õ	0	0	0		
gan	õ	ő	õ	õ	õ	õ	õ		2
head	0	0	õ	0	õ	0	õ		6
hi	0	0		0	0				
1	0	0	0	0	0	0	0		6
milestone	õ	0	0	õ	õ	0		۲	6
note		0	0	0	0				
org	0	0	0	۲	0	0	0	0	
org pol									
p	0		0	0	0	0		۲	
pb		0		0				۲	
pers.spec					۲				
persName first	0	0	0	0	۲	0	0	0	
persName.full				0	۲			0	
persName.last	0	0	0	0	۲	0	0	0	
persName.name		0		0					
quote	0	0	0	0	0	0	0	۲	
ref			0	0				۲	
role	0	0	0	0	0	۲	0	0	
s	0	0	0	0	0	0	0	۲	0
text		0		0			0	۲	
title	0	0	0	0	0	0	۲	0	
top.deoGr		0	۲	0	0	0		0	
top.dr	0	0	۲	0	0		0	0	
top.geo		0	۲	0		0	0	0	
top.gr		0	۲	0					
top.hyd	0	0	۲	0	0	0	0	0	
top.oro	0	0	۲	0				0	
top.reg	0	0	۲	0	0	0	0	0	
top.supReg	0	0	۲	0			0	0	
tonal	0	0	0	0	0	0	0	0	1

The tagsets mapping

- The mapping procedure allows mapping, ignore or removal of XML elements.
- The following XML elements are ignored: back, body, div, foreign, front, gap, head, hi, I, milestone, note, p, pb, quote, ref, s, text
- Mapping is defined as follows:
 - pers.spec, persName.first, persName.full, persName.last, persName.name → PERS
 - top.deoGr, top.dr, top.geo, top.gr, top.hyd, top.oro, top.reg, top.supReg, top.ul → LOC
 - $\circ \qquad \text{demonym} \to \textbf{DEMO}$
 - \circ event \rightarrow **EVENT**
 - $\circ \qquad \text{org, org.pol} \rightarrow \textbf{ORG}$
 - $\circ \qquad \text{role} \to \textbf{ROLE}$
 - \circ title \rightarrow WORK

1.	<role>carice</role> <persname.full>Marije Terezije</persname.full>					
2.	<role>Sekretar</role> <persname.last>Živanović</persname.last>					
3.	<org>Saborne crkve u Beogradu</org>					
4.	<org>manastira Sv. Marka</org>					
5.	<pre><pers>Ali-Milevu</pers></pre>					
1.	<role>carice</role> <pers>Marije Terezije</pers>					
2.	<role>Sekretar</role> <pers>Živanović</pers>					
3.	<org>Saborne crkve u Beogradu</org>					
4.	<org>manastira Sv. Marka</org>					
5.	<pers>Ali-Milevu</pers>					

TXM: POS-tagging and lemmatisation

- TXM (Heiden, 2010) import option "XML/w+CSV"
 - CSV file with metadata
 - TreeTagger serbian model used for POS tagging and lemmatization (Stanković et al. 2020)
 - tokenization applied by a set of rules (not perfect)
- TreeTagger also requires a lexicon and Serbian morphological dictionaries were used as a lexicon for training
- POS tagset: Universal dependencies.

<s><PERS>

<w id="w SRP19101 29920"><txm:form>Sofka</txm:form> <txm:ana resp="#txm" type="#srpos">PROPN</txm:ana> <txm:ana resp="#txm" type="#srlemma">Sofka</txm:ana></w></PERS> <w id="w SRP19101 29921"><txm:form>zastade</txm:form> <txm:ana resp="#txm" type="#srpos">VERB</txm:ana> <txm:ana resp="#txm" type="#srlemma">zastati</txm:ana></w> <w id="w SRP19101 29922"><txm:form>.</txm:form> <txm:ana resp="#txm" type="#srpos">PUNCT</txm:ana> <txm:ana resp="#txm" type="#srlemma">.</txm:ana></w> </s>

Additional transformations

execXSL macro: txm-front-teitxm2xmlw.xsl adapted for level-2

- sentence counting
- use of required namespaces for the attributes xml:id, xml:lang
- removing some attributes
- mapping XML elements for NER tags

 - <rs>, with @type from {PERS, LOC, ORG, DEMO, ROLE, WORK, EVENT}

<xsl:template match="<br">"tei:PERS tei:LOC tei:ORG tei:DEMO tei:ROLE tei:WORK tei:EVENT"> <!-- produce a referring string element--> <xsl:element name="rs" namespace="http://www.tei-c.org/ns/1.0"> <xsl:element name="rs" namespace="http://www.tei-c.org/ns/1.0"> <xsl:element name="rs" namespace="http://www.tei-c.org/ns/1.0"> <xsl:element name="rs" namespace="http://www.tei-c.org/ns/1.0"> <xsl:element name="rs" namespace="http://www.tei-c.org/ns/1.0"> <xsl:element name="rs" namespace="http://www.tei-c.org/ns/1.0"> <xsl:attribute name="type"> <xsl:attribute name="type"> <xsl:attribute name="type"> <xsl:attribute> </xsl:attribute></xsl:attribute></xsl:attribute></xsl:attribute></xsl:element></xsl:element></xsl:element></xsl:element></xsl:element></xsl:element></xsl:template>
<s xml:id="s1333"><rs type="PERS"> <w <br="" lemma="Sofka" pos="PROPN">xml:id="w_SRP19101_29920">Sofka</w></rs> <w <br="" lemma="zastati" pos="VERB">xml:id="w_SRP19101_29921">zastade</w></s>

<s xml:id="s_SRP19101_1333"><rs type="PERS"></rs></s>
<w <="" lemma="Софка" pos="PROPN" td=""></w>
xml:id="w_SRP19101_29920">Софка
<w <="" lemma="3actatu" pos="VERB" td="" xml:id="w_SRP19101_29921"></w>
join="right">застаде
<pc pos="PUNCT" xml:id="w_SRP19101_29922">.</pc>

<w pos="PUNCT" lemma="." xml:id="w SRP19101 29922">.</w></s>

Number of occurrences of NER elements by class (in 65 novels of SrpELTeC)



Frequencies of PERS AND ROLE categories in 65 novels of **SrpELTeC**

PERS	*	F	*	NumN -	RelFKc *
bog			2211	0.52	358.50
Boža			641	0.57	112.82
Milan			601	0.2	60.60
Pera			536	0.24	50.66
Miloš			1203	0.19	41.46
Mara			732	0.24	40.94
Milana			335	0.18	36.19
Stojan			572	0.16	34.57
Srba			257	0.32	31.59
Danica			672	0.14	31.07
Marka			382	0.31	31.00
Jov			428	0.21	30.25
Ljubica			671	0.15	28.08
Sima			634	0.19	26.97
Jelena			366	0.1	26.77
Nikola			356	0.21	23.12
Darinka	1		765	0.08	22.19
Milica			244	0.12	19.09
Ana			563	0.13	17.94
Steva			481	0.14	17.37
Mari			230	0.22	17.33
Sava			265	0.24	16.39
Petar			253	0.24	15.73
Jova			193	0.17	15.00
Ivan			362	0.15	14.72

srlemma 💌	F		•	Num *	Re	elFK 🔻
gospodin		22	38	0.57	4	31.87
рор		12	57	0.44	1	56.53
gospođa		11	23	0.43	14	44.23
kapetan		10	75	0.31	14	42.40
učitelj		9	96	0.45	1	35.87
gazda		6	55	0.51	9	6.91
seljak		5	58	0.52	9	5.79
g.		6	50	0.40	9	0.66
gospodar		6	39	0.41	7	0.64
kmet		8	10	0.30	6	3.31
ministar		4	53	0.24	5	9.29
doktor		4	07	0.35	5	8.37
đak		5	90	0.36	5	0.28
gospođica		4	26	0.33	4	7.37
trgovac		3	12	0.47	4	6.51
predsednik		4	33	0.27	4	5.07
činovnik		2	32	0.40	3	6.38
sluga		2	89	0.43	3	4.89
lekar		2	83	0.35	2	7.26
car		2	94	0.30	2	7.22
pandur		2	72	0.28	2	5.79
knez		4	16	0.24	2	5.16
pisar		3	86	0.27	2	3.49
poslanik		1	78	0.25	2	1.04
radnik	1	1	92	0.33	1	9.78

Frequencies of DEMO AND LOC categories in 65 novels of **SrpELTeC**

DEMO	F	NumNo	RelFKol
Turčin	1788	0.37	159.95
srpski	1047	0.52	148.41
turski	674	0.41	71.44
Srbin	273	0.33	33.20
Ciganin	139	0.26	12.54
beogradski	166	0.27	11.98
nemački	93	0.25	9.82
Srpkinja	83	0.19	8.43
ruski	128	0.21	6.13
grčki	95	0.2	5.29
francuski	67	0.17	3.82
Arnautin	137	0.07	3.76
Grkinja	158	0.04	3.14
Hrvat	58	0.06	2.75
Grk	74	0.14	2.20
Nemac	43	0.14	1.96
Madžar	67	0.1	1.94
Arapin	84	0.09	1.76
užički	35	0.09	1.72
Vlah	26	0.11	1.69
Rus	79	0.11	1.69
carigradski	35	0.11	1.31
Srba	24	0.11	1.11
Bugarin	37	0.11	1.10
Ciganka	31	0.11	1.06

LOC	•	Freq	*	N	um 🔹	Re	IFK(-
Beograd		7	85		0.52	12	20.39
Srbija		594			0.42	5	7.25
Niš		1	18		0.20		7.44
Rusija		1	57		0.20	(5.98
Dunavo		1	18		0.24	(5.93
Kosovo		1	11		0.24	(5.78
Sava		1	05		0.20		5.24
Beč			99		0.19		5.10
Carigrad		1	07		0.18		4.45
Morava		18	81		0.17		4.39
Turska			89		0.19		3.79
Pariz			66		0.16		3.22
Kragujevac			62		0.16		3.03
Pešta	ešta 102			0.11		2.75	
Dunav			60		0.18		2.72
Stambol			50		0.13		2.33
Bosna			86		0.09		2.05
Užice			62		0.08		1.98
Kruševac			65		0.08		1.80
Zemun			30		0.14		1.65
Šumadija			30		0.15		1.59
Karlovci		74			0.07		1.51
Kalemegda	n	43			0.07		1.39
Srem			24		0.15		1.37
Rudnik			90		0.07		1.31

SrpELTeC data structure in Wikidata

Data preparation and mapping with Wikidata

A metadata header <TeiHeader> is required for each text annotated in accordance with TEI recommendations, so it is the case with all ELTeC corpus novels. The mandatory header elements are uniform for all collections and they must contain:

- 1. Description of the electronic edition, which includes the title of the work and the name of the author as well as the statements of responsibility (scanning, correction, annotation), date of publication, size (measured by the number of words). The author and the work can be joined by identifiers, such as viaf and wikidata.
- 2. A brief catalog description of the first edition and the edition used as the source for ELTeC (if different from the first edition).
- 3. Description of the text in terms of meeting the balance criteria (e.g. author's gender, size, date of publication,...).
- 4. Review of all changes to the digital edition since its first publication.

Wikidata - ELTeC data model: classes and instances



Wikidata properties	TEI XPath to element (attribute)	Name of column in prepared data	instance of
P214 (viaf id)	//fileDesc/titleSmt/title@ref	Title_ViafID	Q19832964 (VIAF ID)
P1476 (title)	//titleDesc.titleSmt/title	Title	Q783521 (title)
P50 (author)	//fileDesc/titleSmt/author	Author	Q482980 (author)
P214 (viaf id)	//fileDesc/titleSmt/author@ref	Author_ViafID	Q19832964 (VIAF ID)
P657 (number of words)	//fileDesc/extent/measure@unit	Words	Q8034324 (word count)
P1104 (number of pages)	//fileDesc/extent/measure@unit	Pages	Q1069725 (page)
P123 (publisher)	//fileDesc/publicationSmt/publisher	Publisher	Q105044823 (publisher)
P750 (distributed by)	//fileDesc/publicationSmt/distributor	Distibutor	Q60614978 (distributor)
P6216 (copyright status)	//fileDesc/publicationSmt/availability/licence@target	Licence	Q20007257 (CC BY 4.0)
P50(author)	//fileDesc/bible/author	FirstEdition_author	Q482980 (author)
P146 (title)	//fileDesc/bible/title	FirstEdition_title	Q783521 (title)
P291 (place of publication)	//fileDesc/bible/pubPlace	FirstEdition_pubPlace	Q1361759 (place of publication)
P123 (publisher)	//fileDesc/bible/publisher	FirstEdition_publisher	Q105044823 (publisher)
P577 (publication date)	//fileDesc/bible/date	FistEdition_date	Q1361758 (date of publication)
P407 (language of work or name)	//profileDesc/langUsage/language	Language	Q34770 (language)
P21 (sex or gender)	//textDesc/authorGender	authorGender	Q290 (sex)

Mapping between SrpELTeC metadata and Wikidata



Structure of SrpELTeC Wikidata items (example of ELTeC edition)



Wikidata data structure for ELTeC novels

the case of a novel: as a literary work (Q7725634)					the case of edition: as instance of edition (Q3331189) 1) first edition (Q10898227) and 2) electronic (ELTeC) edition (Q59466853)				
Ивкова слава :	приповетка	(0100006080)		Ивков	а слава : припов	етка : ELTeC издање 🕡	Q107648205)		
роман српског писца	приноветка	(0109330082)	🖍 уреди	ЕLTeC издањ и → На другим Конфигуриши	е српског романа језицима		🖋 уреди		
Конфигуриши				Језик	Ознака	Опис	Псеудоними		
Језик Ознака српски / srpski Ивкова енглески Ознака	слава : приповетка није дефинисана	Опис роман српског писца Опис није дефинисан	Псеудоними	српски / srps енглески	ki Ивкова слава : приповет ELTeC издање Ознака није дефинисане	ка : ELTeC издање српског романа Опис није дефинисан			
издање	Ивкова слава : прили језик дела или имена датум издавања	оветка : ELTeC издање српски језик 2021	🖋 уреди	број речи	€ 42.963 начин одрег ~ 0 рефере	Hansa word count enecku	🖉 уроди		
	место издавања	Београд					+ додај референцу		
	• 0 референие						+ додај вредност		
	🗧 Ивкова слава : прили	оветка	+ додај референцу ✔ уреди	пун текст дост	упан на 👸 https://distan	reading.github.io/ELTeC/srp/SRP18950.html)	🖉 уреди		
	језик дела или имена Датум издавања место издавања	а српски језик 1895 Београд			* 0 рефере	વાલ	+ додај референцу		
	- 0 референце	Баоград			https://udalje	nocitanje.unilib.rs/pregled/15/strana/1/	🖍 уреди		
			+ додај референцу				+ додај референцу		
			+ додај вредност						

Transforming natural language into Wikidata

Ivkova slava is a literary work written by Stevan Sremac.	Ivkova slava (Q107648205) is (P31) a literary work (Q7725634) written by (P50) Stevan Sremac (Q559989).	Q107648205 P31 Q7725634;
Ivkova Slava was published in Belgrade in 1899.	Ivkova Slava (Q107648205) was published in (P291) Belgrade (Q3711) in (P577) "1899".	P50 Q559989; P291 Q3711; P577 "1899".
Stevan Sremac was born on 23rd November 1855 in Senta, and he died on 26th august 1906. He was a writer and belonged to realism. He's VIAF ID is 66526515.	Stevan Sremac (Q559989) was born on (P569) "23rd November 1855". in (P19) Senta (Q571136). and he died on (P570) ,,26th august 1906". He (Q559989) was (P106) a writer (Q36180) and belonged to (P135) realism (Q667661) . He (Q559989) has VIAF ID (P214) "66526515".	Q559989 P569 "23rd november 1855"; P19 Q571136; P570 "26th august 1906"; P106 Q36180: P135 Q667661; P214 "66526515".

SrpELTeC Wikidata entry and enrichment automation

The manual population of individual data (in our case ELTeC editions of novels) into Wikidata is often a time-consuming task. In this way 54 novels from the SrpELTeC sub-collection were described in Wikidata by students at the University of Belgrade through different activities. The control of manual entries revealed that some of the entries were incomplete or contained incorrect information, such as incomplete novel title, an author's VIAF entered as a work's VIAF , connecting wrong people as authors (e.g. a football player Dušan Đurić (Q116994) with the same name as a writer (Q108986248)) or wrong year of the first edition.





Automatization with web resources

Since the automation of the process of data preparation and import was envisaged (using <TeiHeader> data), the different solutions were analysed and finally synergy of OpenRefine and QuickStatements tools was chosen as the best choice.

Tutorial for automatization using OpenRefine and QuickStatement





Tutorial for automatization using OpenRefine and QuickStatement

... Step by step, ilustrated tutorial with SrpELTeC text collection

nRefine Novels_V5 Permalink							c NOVEIS_V5 Permalink "https://distantreading.github.io/EL										
ter Undo / Redo 28 / 28 🔇	21	row	s				do / Redo 50 / 50	<	21 rows								
Reset All Remove All Show as: rows records Show: 5 10 25 50 rows						S	Reset All Remo	eset All Remove All		now as: rows records Show: 5 10 25 50 rows							
idament change		u	TD ID		author	▼ full_text_on		change	T All	💌 ID	 author 	v full_text_on					
t by: name count	☆	Image: State		https://distantreading.github.io/ELTeC/srp/SRP19060			☆ 딕 1.	SRP19060	Q108983669 Choose new match	https://distantreading.github.io/ELTeC/srp/SRP19060.html new Choose new match							
e counts	☆	☆ 🖓 2. SRP18791			Search for match Коста Барунчий	https://distantreading.github.io/ELTeC/srp/SRP18791	l.		숣 듸 2.	SRP18791	Q108983662 Choose new match	https://distantreading.github.io/ELTeC/srp/SRP18791.html new Choose new match	Пастир к приповет				
est candidate's score change reset value present.				00040000	Search for match		ent.		☆ 🦪 3.	SRP19090	Ivo Ćipiko Choose new match	https://distantreading.github.io/ELTeC/srp/SRP19090.html edit					
r: judgment change	24	-4 ·	o. :	24619090	Search for match	nitps.//distantreading.gtmub.to/ELTeC/srp/SRP1909L	it	change	숤 듸 4.	SRP18820	Jakov Ignjatović Choose new match	Search to yearch https://github.com/COST-ELTeC/ELTeC-srp- ext/blob/master/level1/SRP18820_JakovI_PoetaAdvokat.xml	Поета и : Choose nev				
t by: name count	☆	9 4	4. 5	SRP18820	Jаков Игњатовић	https://github.com/COST-ELTeC/ELTeC-srp- ext/blob/master/level1/SRP18820_JakovI_PoetaAdv	count			00040000	Paris	Create new item Search for match	Gaustino				
e counts	☆	9	5. 5	SRP19022	Борисав Станковић Э Э Награда Борисав	https://distantreading.github.io/ELTeC/srp/SRP19022	didate's chang	ge reset	ж ц э.	SRP19022	Stanković Choose new match	Create new item	Choose new				
r: best candidate's change reset					CTAHKOBUN (42)				☆ 디 6.	SRP18941	Stevan Sremac Choose new match	https://distantreading.github.io/ELTeC/srp/SRP18941.html	Пол Ћир издање Choose nev				
0	\$	9	3. 9	SRP18941	Стеван Сремац Стеван Сремац (37) Стеван Сремац (37)	https://distantreading.github.io/ELTeC/srp/SRP18941			☆ 딕 7.	SRP19000	Svetolik Ranković Choose new match	https://distantreading.github.io/ELTeC/srp/SRP19000.html	Choose new				
32.5 — 41.52 □ Non-numeric ☑ Blank □ Error					 ✓ ✓ ОШ "Стеван Сремац" Сурчин (35) ✓ ✓ Награда Стеван Сремац за књигу прозе (27) 	Œ	— 41.52 meric 🗹 Blank 🗌	Error	☆ 디 8.	SRP19110	Milutin Uskoković Choose new match	https://distantreading.github.io/ELTeC/srp/SRP19110.html					
0 17 0					 Основна школа Стеван Сремац у Борчи (27) Награда Стеван Сремац 	Q	gment	change	☆ 더 9.	SRP19091	Сретен Динић Choose new match	https://distantreading.github.io/ELTeC/srp/SRP19091.html					
					за новинску причу (26) У Зграда Основне школе		e count		☆ 듸 10	SRP1920a	Ivo Andrić Choose new match	https://github.com/COST-ELTeC/ELTeC-srp- ex/tbb/master/level1/SRP1920a_IvoA_AlijaGj.xml	Пут Алиј Choose ner				

When creating new item "Create new item" manual option was chosen for each cell, to check the results. There is an option to automatically create a new item, for each cell in the column.

Examples of reconciling cells for ELTeC edition of novels

- 1. **Title** (P1476) to an entity of type *edition, version, or translation* (Q3331189)
- 2. Author (P50) to an entity of type *human* (Q5) and then search for match
- 3. Language of work or name (P407) to an entity of type *language* (Q34770)
- 4. **Number of pages** (P1104) to an entity of type *natural number* (Q21199)
- 5. Number of words (P6570) to an entity of type *natural number* (Q21199)
- 6. **Published in** (P1433) to an entity of type *text corpus* (Q461183)
- 7. **VIAF ID** (P214) to an entity of type *VIAF ID* (Q19832964)
- 8. Full work available at URL (P953) to an entity of type URL (Q42253)
- 9. **Publication date** (P577) to an entity of type *calendar year* (Q3186692)
- 10. **Place of publication** (P291) to an entity of type *city* (Q515)
- 11. **Volume** (ID of novel) (P478) to an entity of type *volume* (Q1238720) For all novels their identification (ID) assigned as the property volume of a book. e.g. SRP19012

Finally... OpenRefine

Creating a Wikidata input set schema defines predicates (properties) that will connect subjects and objects in RDF triples

Each statement for a subject has a property and value that can be a Wikidata item, external URL, or literal (string)

As presented in Table on slide 18, the property from the first column is related to content (values: items or literals) in the third column.

The property from the first column is related to content (values: items or literals) in the third column

After editing and saving the Wikidata schema it was exported as a QuickStatements file.

In the final stage the prepared file was exported in the QuickStatements tool and Wikidata items were automatically created.

Wikidata properties	TEI XPath to element (attribute)	Name of column in prepared data	instance	
P214 (viaf id)	//fileDesc/titleSmt/title@ref	Title_ViafID	Q19832964 (VIAF	
P1476 (title)	//titleDesc.titleSmt/title	Title	Q783521 (title)	
P50 (author)	//fileDesc/titleSmt/author	Author	Q482980 (author)	
P214 (viaf id)	//fileDesc/titleSmt/author@ref	Author_ViafID	Q19832964 (VIAF	
P657 (number of words)	//fileDesc/extent/measure@unit	Words	Q8034324 (word o	
P1104 (number of pages)	//fileDesc/extent/measure@unit	Pages	Q1069725 (page)	
P123 (publisher)	//fileDesc/publicationSmt/publisher	Publisher	Q105044823 (pub)	
P750 (distributed by)	//fileDesc/publicationSmt/distributor	Distibutor	Q60614978 (distri	
P6216 (copyright status)	//fileDesc/publicationSmt/availability/licence@target	Licence	Q20007257 (CC B	

CREATE		
LAST	Lsr "By	л-Марикина прикажња : приповетка : ELTeC издање"
LAST	Dsr "EL	ТеС издање романа српског писца"
LAST	P31 Q33	31189
LAST	P1433	Q106927517
LAST	P1433	Q106936149
LAST	P1476	sr:"Ђул-Марикина прикажња : приповетка : ELTeC издање"
LAST	P50 Q36	25974
LAST	P407	Q9299
LAST	P577	+2021-00-00T00:00:00Z/9
LAST	P291	Q3711
LAST	P1104	107
LAST	P6570	20244
LAST	P953	"https://distantreading.github.io/ELTeC/srp/SRP19012.html"
LAST	P478	"SRP19012"

Overview of SrpELTeC@Wikidata by SPARQL queries









Црна Гора

Kosovë /



https://w.wiki/4LuC

Visualization of an author (Stevan Sremac) and his properties













Layers can be reduced (hiden)...

\leftarrow	\rightarrow C	A Neb	ezbedno	inception.jer	teh.rs/p/27/annota	ite?19#!d=	=259&f=389								to	Ġ	5≡ 5	Ð 🙎	
INC	EpTION		rojects	📕 Dash	board						1)	🛛 Help	💥 Adminis	stration	💄 adm	in I	🕩 Log o	ut 🛈 2	29 min
	milica.il	konic: Sr	ELTeC-	Wiki-NEL/	/SRP19101_N	ecistaK	rv.tsv				389-393	/ 3912 se	ntences [doc 1 ,	/ 1]	Annotati	on			
	Image: Constraint of the state of																		
	389 Posle dođe rat , i oslobođenje , nestanak turske vlasti i gospodstva , pa i nestanak Sofkinog oca , efendi - <u>Mita PERS</u> Mite .										Text Mite				•				
	DEMO ROLE LOC 390 Sa Turcima i begovima i on prebegao , i tamo u Turskoj počeo da se bavi , tobož trgujući s njima . No link 391 Retko bi otuda ovamo prelazio . Identifie									identifier	or rela	notation.	nect to t	his					
	Sofka PERS 392 I što je Sofka bivala veća , on je sve ređe dolazio ; u godini dana jedan put i to obično noću . 393 Ostane po dva i tri dana , ali nikuda iz kuće ne izlazi .									Character from novel Impure blood				b					
															PERS				

Export to NIF

"NLP Interchange Format (NIF):

- an RDF/OWL-based format that allows to combine and chain several **NLP** tasks in a flexible, light-weight way.
- The core of **NIF** consists of a vocabulary, which can represent Strings as RDF resources.
- A special URI Design is used to pinpoint annotations to a part of a document."

<file:/srv/...#offset_39967_39972>

a	<pre>nif:Word , nif:EntityOccurrence , nif:OffsetBasedString ;</pre>
nif:anchorOf	"Sofki" ;
nif:beginIndex	"39967"^^xsd:nonNegativeInteger ;
nif:endIndex	"39972"^^xsd:nonNegativeInteger ;
nif:lemma	"Sofka" ;
nif:nextWord	<file: #offset="" 39973="" 39974="" srv=""> ;</file:>
nif:posTag	"PROPN" ;
nif:previousWord	<file: #offset="" 39965="" 39966="" srv=""> ;</file:>
nif:referenceContext	<file: #offset="" 0="" 378682="" srv=""> ;</file:>
nif:sentence	<file: #offset="" 39885="" 40052="" srv=""> ;</file:>
itsrdf:taClassRef	<pre><pers> ;</pers></pre>
itsrdf:taIdentRef	<http: entity="" q109693861="" www.wikidata.org=""> .</http:>

SrpELTeC @Wikidata status

Wikidata has 110 novels from SrpELTeC collection with associated items for the first edition and electronic SrpELTeC editions.

That means that we automatically added more than **2500** statements.



Conclusion

- Presented examples can be seen as succesful show cases
- Similar approach can be used in many other similar cases
- Places and main characters from novels will be further added to Wikidata and linked in text
- <u>Web annotation</u> Data Model, W3C Recommendation (2017), implementation is planned

Step by step authors: <u>https://scribehow.com/shared/OpenRefine__0K5TU3J9SmalqRBYePhdhQ</u>

Step by step novels: <u>https://scribehow.com/shared/OpenRefine_AOPgpu7dRqmxUkYBtVbGgQ</u>