¹Istituto di Linguistica Computazionale CNR, Pisa, Italy
 ²CRH, EHESS Paris, France
 ³Institute of Polish Language, Polish Academy of Sciences
 ⁴LaTTiCe CNRS, Université Sorbonne Nouvelle - Paris 3
 ⁵Linguateca & University of Oslo, Norway
 ⁶University of Belgrade, Serbia

Named Entity Recognition for Distant Reading in ELTeC

Francesca Frontini¹, Carmen Brando², Joanna Byszuk³, Ioana Galleron⁴, Diana Santos⁵, Ranka Stanković⁶

CLARIN annual conference 5-7 October 2020

Overview of the talk



Introduction

Manual annotation

Annotation guidelines

Software infrastructure

BRAT

NER & Beyond

Comparing manual and automatic annotation

Testing Automatic NER

Results

Analysis

Links to CLARIN

Links



Distant [■] Reading

"Distant Reading for European Literary History (COST Action CA16204) is a project aiming to create a vibrant and diverse network of researchers jointly developing the resources and methods necessary to change the way European literary history is written. Grounded in the Distant Reading paradigm (i.e. using computational methods of analysis for large collections of literary texts), the Action will create a shared theoretical and practical framework to enable innovative, sophisticated, data-driven, computational methods of literary text analysis across at least 10 European languages."

https://www.distant-reading.net

Motivation



- Most NER systems and studies are
 - not specifically concerned with literary texts
 - not dealing with 19th century texts
 - not multilingual enough
- ▶ We wanted to explore the issues in COST texts
 - what kind of entities
 - what kind of problems
- ► Create an infrastructure to evaluate current NER systems for literary texts in the languages of COST

Exercise



- Languages: Czech, German, English, French, Hungarian, Norwegian, Portuguese, Serbian and Slovene
- "Lightweight semantic annotation" defined: marking persons, demonyms, professions and other roles, works, places, facilities and organizations
- ► 5 random passages of 400 white space-delimited tokens taken from 20 novels from each language, manually annotated using brat



EVENT: Pussiness
EVENT: Bussiness
EVENT: Natural_disaster
EVENT: Political
EVENT: Culture
EVENT: Sport
EVENT: Religious

NE annotation



Table: Current manually NE-annotated corpus¹.

	DEMO	EVENT	LOC	ORG	OTHER	PERS	ROLE	WORK
cze	163	5	275	0	0	1150	454	0
deu	66	2	323	12	0	973	458	4
eng	56	7	198	37	0	1184	203	25
fra	77	3	262	22	128	900	244	18
hun	29	7	152	20	0	1091	367	7
nor	4	8	83	25	3	990	201	10
por1	17	9	351	19	0	940	490	54
por2	34	1	256	30	7	1059	347	7
slv	133	54	336	37	0	1230	620	2
srp	121	18	185	11	0	985	301	4
	700	114	2425	213	138	10514	3685	131

¹ Portuguese the exercise was done twice: with canonical, modernized ortography works (por1), and with non-canonical, old-ortography works (por2)

Annotation guidelines



Literary text show a broad variety NEs with respect to more standardised non fictional texts.

- ▶ People are often referred to by profession or their origin, as well as by family relations only ("Maman"). . .
- ► Fictional characters may be present and animals and objects may have proper names.
- ▶ Additional types of entities, which are particularly interesting for the purposes of literary and cultural analysis, require annotation, such titles of works of art, books, publications; literary movements, . . .

Classification is sometimes difficult and there is a risk of proliferation of categories and sub-categories. E.g. "Vive la **Réforme**! à bas Giuzot!" (here "**Réforme**" refers to a proposal of reform of electoral law which the politician Giuzot opposed).

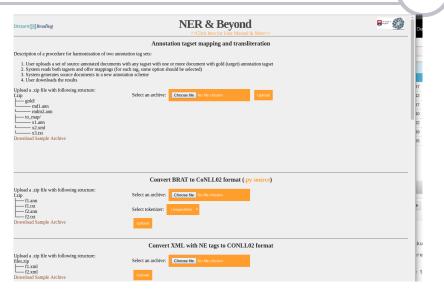
Manual annotation infrastructure: BRAT



←	NER Guidelines NER&Beyond Path: /eltec-gold/por2/POR0046_FraGAmo_Selvagens_sample
	PERS PERS
	despenharem-se outra vez no abysmo da ignorância! Meu Deus, meu Deus, compadece-te d'elles!
6	<pre><pn="por0046398">— N\u00e40 desesperes assim !</pn="por0046398"></pre>
7	</td
	i i
	PERS
8	n="POR0046400">Decorrera um anno depois que Manuel Félix e sua mulher se tinham separado do padre. Este continuava quasi entrevado,
	PLACE
	em Santarém; os seus passeios limitavam-se a andar encostado a um pau, á roda do cercado ou quintal povoado de bananeiras, araticús,
	eni Santaleni, os seus passelos liniliavani-se a anual encostado a um pau, a roda do cercado ou quintal povoado de bahaneiras, aradicus,
	ROLE
	mangueiras e, laranjeiras, que havia na casa do irmão. Ás horas do calor atavam-se as redes nas arvores; o padre deitava-se n'uma d'ellas,
	DEMO! DEES!
	lamentando sempre o seu estado e a perda dos seus Índios das cachoeiras : e Romualdo sentava-se n'outra, consolando-o com palavras
	aíTectuosas, ou lendo-lhe, n'algum livro favorito, cousas de religião e de moral.
	PLACE
9	<
	EVENT PLACE
	com as noticias da revolução de 1820, em Portugal. Quem podia lembrar-se, vendo abalados pela base os alicerces das velhas sociedades
	PLACE PLACE DEMO
	europeas, de que havia um rio no Brazil chamado Tapajós, onde uns pobres índios necessitavam de quem os ajudasse a completar a obra da
	sua redempção ?! A sabedoria humana trabalhava em grande n'esse momento; não podia descer a miudezas; occupava-se do geral, e não do
	particular ; pensava em reformar nacões, e não em civilisar aldeias. O vento que andava no mundo sacudia as cabecas, e fazia cair d'ellas
	opiniões que espantavam os próprios que as enunciavam. Era 1820 o precursor das novas ideias, que annos depois transformavam a colónia
	ROLE
	em paiz independente e a metrópole em terra de homens livres. Rei e governos, diante da gravidade das circumstancias, tratavam da própria
	PERS
	salvação, 6 não da alheia. O padre Félix comprehendeu isto ; deu a sua causa por perdida, e, sem se resignar nem esquecer d'ella, deixou de
	esperar o remédio, que sempre suspeitou que lhe não mandariam.
10	

Web portal NER & Beyond





Web portal NER & Beyond: tools



At University of Belgrade and Jerteh, http://nerbeyond.jerteh.rs/

- ► Conversion tools among different annotation formats
 - ► BRAT to CoNLL02
 - ► XML with NE tags to CONLL02
 - ► XML with NE tags to BRAT
 - ► BRAT (ann & txt) to XML
 - ► CONLL02 to BRAT
- Automatic Named entity recognition (and annotation)
 - spaCy (eng, fra, ita, nld, ger, por, srp, spa, multi)
 - Stanford (eng, ger, srp)
 - CVNER (ser)
- lacktriangle Computation of statistics based on BRAT .ann format
- ► NER evaluation with the Gemini tool
 - Precision and recall tables
 - ► Visual comparison in HTML
- ► Annotation tagset harmonization (mapping and transliteration)

Example of comparison of gold (manual, blue) and automatic (SpaCy, pink) annotation

<sample>And at this moment entered the room the young nobleman whom we have before mentioned, accompanied by an individual who was approaching perhaps the termination of his fifth lustre but whose general air rather betokened even a less experienced time of life. Tall, with a well-proportioned figure and a graceful carriage, his countenance touched with a sensibility that at once engages the affections. Charles Egremont was not only admired by that sex, whose approval generally secures men enemies among their fellows, but was at the same time the favourite of his own."Ah, Egremont! come and sit here," exclaimed more than one banqueter. "I saw you waltzing with the little Bertie , old fellow," said Lord Fitzheron , "and therefore did not stay to speak to you, as I thought we should meet here. I am to call for you, mind." "How shall we all feel this time to-morrow?" said Egremont, smiling. "The happiest fellow at this moment must be Cockie Graves ," said Lord Milford . "He can have no suspense. I have been looking over his book, and I defy him, whatever happens, not to lose." "Poor Cockie ." said Mr Berners; "the has asked me to dine with him at the Clarendon on Saturday." Annotation <PERS> "Cockie is a very good Cockie", said Lord Milford, and Text ENG18450 Disraeli sample.xml any gentleman sportsman present wishes to give seven to two, I will take him to any ar Text File2: ANN FILE "My book is made up." said Egremont; "and I stand or fall by

Testing Automatic NER



- Case study on four collections only, in English, French, Portuguese and Serbian
- ► For each collection, we tested two tools: one common for all (spaCy) and another one language specific
- ▶ BRAT outputs were compared to annotations produced by these tools
- Evaluation of string detection was strict (segments must match exactly)

Testing Automatic NER - results



	Cat	Correct	Missing	Spurious	Precision	Recall	Excess
CEMC	LOC	73	112	. 84	0.465	0.395	0.535
SEM-fra	PERS	82	512	115	0.416	0.138	0.584
CDACY (LOC	103	78	468	0.180	0.569	0.820
SPACY-fra	PERS	329	194	297	0.526	0.629	0.474
DAL AV/DAC1	LOC	223	63	44	0.835	0.780	0.165
PALAVRAS-por1	PERS	816	90	86	0.905	0.901	0.095
CDACV1	LOC	225	84	440	0.338	0.728	0.662
SPACY-por1	PERS	465	256	374	0.554	0.645	0.446
PALAVRAS-por2	LOC	151	67	91	0.624	0.693	0.376
PALAVKA3-por2	PERS	857	133	285	0.750	0.866	0.250
SPACY-por2	LOC	157	57	396	0.284	0.734	0.716
SPACT-porz	PERS	569	236	393	0.591	0.707	0.409
Stanford-eng	LOC	98	100	126	0.438	0.495	0.563
Stanioru-eng	PERS	649	535	399	0.619	0.548	0.381
SPACY-eng	LOC	98	100	170	0.366	0.495	0.634
SPACT-eng	PERS	536	648	240	0.691	0.453	0.309
SrpNER-srp	LOC	107	78	19	0.849	0.578	0.151
21hirriv-21h	PERS	718	267	158	0.820	0.729	0.180
SPACY-srp	LOC	57	128	104	0.354	0.308	0.646
SFACT-STP	PERS	553	432	315	0.637	0.561	0.363

Table: Results of the strict evaluation, per language and category.

Testing Automatic NER - analysis



- ► A strict evaluation of detection is often penalising for PERS, because of honorifics which we chose to include in our annotation
- ▶ It is further complicated by the fact that the XML annotated input was processed as such by tools which often expect plain text
- ▶ In most cases, LOC seems to be less problematic for the pre-trained models.

Links to CLARIN



- ▶ The NE corpus, as part of the larger annotated ELTeC corpus, will have to adhere to the FAIR requirements, for long term preservation, visibility, accessibility. A deposit in a repository (possibly Textgrid) is envisaged at the end of the project.
- ► Accessibility could be further enhanced by making the corpus part of the CLARIN federated content search
- ► Finally, CLARIN collaboration is welcome for the creation of better adapted NE tools which can overcome the problems that we have identified for state of the art tools.

Useful links



- ► The COST Action web page: https://www.distant-reading.net
- ► The ELTeC corpus: https://distantreading.github.io/ELTeC/index.html
- ► The current NE annotated sub-corpus http://brat.jerteh.rs/#/eltec-simplified/

Acknowledgements



Presented research is supported by COST Action 16204 – Distant Reading for European Literary History

* * *

We thank Lou Burnard for creating the samples, Tajda Liplin, Zala Vidic, Karolina Zgaga, Michael Preminger, Tonje Vold, Emma Takács, Silvie Cinková, Klára Macháčková, and Cvetana Krstev for annotating the texts, and Branislava Šandrih for help in creating the sites at the University of Belgrade, and all our colleagues at the COST action Distant Reading for European Literary History.

Thank you!