

Extended Overview of HIPE-2022: Named Entity Recognition and Linking in Multilingual Historical Documents

Maud Ehrmann¹, Matteo Romanello², Sven Najem-Meyer¹, Antoine Doucet³ and Simon Clematide⁴

¹Digital Humanities Laboratory, Ecole Polytechnique Fédérale de Lausanne (EPFL), Lausanne, Switzerland

²University of Lausanne, Lausanne, Switzerland

³University of La Rochelle, La Rochelle, France

⁴Department of Computational Linguistics, University of Zurich, Zurich, Switzerland

Abstract

This paper presents an overview of the second edition of HIPE (Identifying Historical People, Places and other Entities), a shared task on named entity recognition and linking in multilingual historical documents. Following the success of the first CLEF-HIPE-2020 evaluation lab, HIPE-2022 confronts systems with the challenges of dealing with more languages, learning domain-specific entities, and adapting to diverse annotation tag sets. This shared task is part of the ongoing efforts of the natural language processing and digital humanities communities to adapt and develop appropriate technologies to efficiently retrieve and explore information from historical texts. On such material, however, named entity processing techniques face the challenges of domain heterogeneity, input noisiness, dynamics of language, and lack of resources. In this context, the main objective of HIPE-2022, run as an evaluation lab of the CLEF 2022 conference, is to gain new insights into the *transferability* of named entity processing approaches across languages, time periods, document types, and annotation tag sets. Tasks, corpora, and results of participating teams are presented. Compared to the condensed overview [1], this paper contains more refined statistics on the datasets, a break down of the results per type of entity, and a discussion of the ‘challenges’ proposed in the shared task.

Keywords

Named entity recognition and classification, Entity linking, Historical texts, Information extraction, Digitised newspapers, Classical commentaries, Digital humanities

1. Introduction

Through decades of massive digitisation, an unprecedented amount of historical documents became available in digital format, along with their machine-readable texts. While this represents a major step forward in terms of preservation and accessibility, it also bears the potential for new ways to engage with historical documents’ contents. The application of machine reading to


CLEF 2022: Conference and Labs of the Evaluation Forum, September 5–8, 2022, Bologna, Italy

✉ maud.ehrmann@epfl.ch (M. Ehrmann); matteo.romanello@unil.ch (M. Romanello); sven.najem-meyer@epfl.ch (S. Najem-Meyer); antoine.doucet@univ-lr.fr (A. Doucet); simon.clematide@uzh.ch (S. Clematide)

🆔 0000-0001-9900-2193 (M. Ehrmann); 0000-0002-7406-6286 (M. Romanello); 0000-0002-3661-4579 (S. Najem-Meyer); 0000-0001-6160-3356 (A. Doucet); 0000-0003-1365-0662 (S. Clematide)



© 2022 Copyright for this paper by its authors. Use permitted under Creative Commons License Attribution 4.0 International (CC BY 4.0).

 CEUR Workshop Proceedings (CEUR-WS.org)

historical documents is potentially transformative and the next fundamental challenge is to adapt and develop appropriate technologies to efficiently search, retrieve and explore information from this ‘big data of the past’ [2]. Semantic indexing of historical documents is in great demand among humanities scholars, and the interdisciplinary efforts of the digital humanities (DH), natural language processing (NLP), computer vision and cultural heritage communities are progressively pushing forward the processing of facsimiles, as well as the extraction, linking and representation of the complex information enclosed in transcriptions of digitised collections [3]. In this regard, information extraction techniques, and particularly named entity (NE) processing, can be considered among the first and most crucial processing steps.

Yet, the recognition, classification and disambiguation of NEs in historical texts is not straightforward, and performances are not on par with what is usually observed on contemporary well-edited English news material [4]. In particular, NE processing on historical documents faces the challenges of domain heterogeneity, input noisiness, dynamics of language, and lack of resources [5]. Although some of these issues have already been tackled in isolation in other contexts (with e.g., user-generated text), what makes the task particularly difficult is their simultaneous combination and their magnitude: texts are severely noisy, and domains and time periods are far apart.

Motivation and Objectives. As the first evaluation campaign of its kind on multilingual historical newspaper material, the CLEF-HIPE-2020 edition¹ [6, 7] proposed the tasks of NE recognition and classification (NERC) and entity linking (EL) in ca. 200 years of historical newspapers written in English, French and German. HIPE-2020 brought together 13 teams who submitted a total of 75 runs for 5 different task bundles. The main conclusion of this edition was that neural-based approaches can achieve good performances on historical NERC when provided with enough training data, but that progress is still needed to further improve performances, adequately handle OCR noise and small-data settings, and better address entity linking. HIPE-2022 attempts to drive further progress on these points, and also confront systems with new challenges. An additional point is that in the meantime several European cultural heritage projects have prepared additional NE-annotated text material, thus opening a unique window of opportunity to organize a second edition of the HIPE evaluation lab in 2022.

HIPE-2022² shared task focuses on named entity processing in historical documents covering the period from the 18th to the 20th century and featuring several languages. Compared to the first edition, HIPE-2022 introduces several novelties:

- the addition of a new type of document alongside historical newspapers, namely classical commentaries³;
- the consideration of a broader language spectrum, with 5 languages for historical newspapers (3 for the previous edition), and 3 for classical commentaries;
- the confrontation with heterogeneous annotation tag sets and guidelines.

¹<https://impresso.github.io/CLEF-HIPE-2020>

²<https://hipe-eval.github.io/HIPE-2022/>

³Classical commentaries are scholarly publications dedicated to the in-depth analysis and explanation of ancient literary works. As such, they aim to facilitate the reading and understanding of a given literary text.

Overall, HIPE-2022 confronts participants with the challenges of dealing with more languages, learning domain-specific entities, and adapting to diverse annotation schemas. The objectives of the evaluation lab are to contribute new insights on how best to ensure the transferability of NE processing approaches across languages, time periods, document and annotation types, and to answer the question whether one architecture or model can be optimised to perform well across settings and annotation targets in a cultural heritage context. In particular, the following research questions are addressed:

1. How well can general prior knowledge transfer to historical texts?
2. Are in-domain language representations (i.e. language models learned on the historical document collections) beneficial, and under which conditions?
3. How can systems adapt and integrate training material with different annotations?
4. How can systems, with limited additional in-domain training material, (re)target models to produce a certain type of annotation?

Recent work on NERC showed encouraging progress on several of these topics: Beryozkin et al. [8] proposed a method to deal with related, but heterogeneous tag sets. Several researchers successfully applied meta-learning strategies to NERC to improve transfer learning: Li et al. [9] improved results for extreme low-resource few-shot settings where only a handful of annotated examples for each entity class are used for training; Wu et al. [10] presented techniques to improve cross-lingual transfer; and Li et al. [11] tackled the problem of domain shifts and heterogeneous label sets using meta-learning, proposing a highly data-efficient domain adaptation approach.

The remainder of this paper is organized as follows. Sections 2 and 3 present the tasks and the material used for the evaluation. Section 4 details the evaluation framework, with evaluation metrics and the organisation of system submissions around tracks and challenges. Section 5 introduces the participating systems, while Section 6 presents and discusses their results. Finally, Section 7 summarizes the benefits of the task and concludes.⁴

2. Task Description

HIPE-2022 focuses on the same tasks as CLEF-HIPE-2020, namely:

Task 1: Named Entity Recognition and Classification (NERC)

- **Subtask NERC-Coarse:** this task includes the recognition and classification of high-level entity types (person, organisation, location, product and domain-specific entities, e.g. mythological characters or literary works in classical commentaries).
- **Subtask NERC-Fine:** includes the recognition and classification of entity mentions according to fine-grained types, plus the detection and classification of nested entities of depth 1. This subtask is proposed for English, French and German only.

⁴For space reasons, the discussion of related work is included in the extended version of this overview [12].

Table 1

Overview of HIPE-2022 datasets with an indication of which tasks they are suitable for according to their annotation types.

Dataset alias	Document type	Languages	Suitable for
hipe2020	historical newspapers	de, fr, en	NERC-Coarse, NERC-Fine, EL
newseye	historical newspapers	de, fi, fr, sv	NERC-Coarse, NERC-Fine, EL
sonar	historical newspapers	de	NERC-Coarse, EL
letemps	historical newspapers	fr	NERC-Coarse, NERC-Fine
topres19th	historical newspapers	en	NERC-Coarse, EL
ajmc	classical commentaries	de, fr, en	NERC-Coarse, NERC-Fine, EL

Task 2: Named Entity Linking (EL) This task corresponds to the linking of named entity mentions to a unique item ID in Wikidata, our knowledge base of choice, or to a NIL value if the mention does not have a corresponding item in the knowledge base (KB). We will allow submissions of both end-to-end systems (NERC and EL) and of systems performing exclusively EL on gold entity mentions provided by the organizers (EL-only).

3. Data

HIPE-2022 data consists of six NE-annotated datasets composed of historical newspapers and classic commentaries covering ca. 200 years. Datasets originate from the previous HIPE-2020 campaign, from HIPE organisers' previous research project, and from several European cultural heritage projects which agreed to postpone the publication of 10% to 20% of their annotated material to support HIPE-2022. Original datasets feature several languages and were annotated with different entity tag sets and according to different annotation guidelines. See Table 1 for an overview.

3.1. Original Datasets

Historical newspapers. The historical newspaper data is composed of several datasets in English, Finnish, French, German and Swedish which originate from various projects and national libraries in Europe:

- **HIPE-2020 data** corresponds to the datasets of the first HIPE-2020 campaign. They are composed of articles from Swiss, Luxembourgish and American newspapers in French, German and English (19C-20C) that were assembled during the *impresso* project⁵ [13]. Together, the train, dev and test hipe2020 datasets contain 17,553 linked entity mentions, classified according to a fine-grained tag set, where nested entities, mention components and metonymic senses are also annotated [14].

⁵<https://impresso-project.ch>

- **NewsEye data** corresponds to a set of NE-annotated datasets composed of newspaper articles in French, German, Finnish and Swedish (19C-20C) [15]. Built in the context of the *NewsEye* project⁶, the newseye train, dev and test sets contain 36,790 linked entity mentions, classified according to a coarse-grained tag set and annotated on the basis of guidelines similar to the ones used for hipec2020. Roughly 20% of the data was retained from the original dataset publication and is published for the first time for HIPE-2022, where it is used as test data (thus the previously published test set became a second dev set in HIPE-2022 data distribution).
- **SoNAR data** is an NE-annotated dataset composed of newspaper articles from the Berlin State library newspaper collections in German (19C-20C), produced in the context of the *SoNAR* project⁷. The sonar dataset contains 1,125 linked entity mentions, classified according to a coarse-grained tag set. It was thoroughly revised and corrected on NE and EL levels by the HIPE-2022 organisers. It is split in a dev and test set – without providing a dedicated train set.
- **Le Temps data**: a previously unpublished, NE-annotated diachronic dataset composed of historical newspaper articles from two Swiss newspapers in French (19C-20C) [4]. This dataset contains 11,045 entity mentions classified according to a fine-grained tag set similar to hipec2020.
- **Living with Machines data** corresponds to an NE-annotated dataset composed of newspaper articles from the British Library newspapers in English (18C-19C) and assembled in the context of the *Living with Machine* project⁸. The topres19th dataset contains 4,601 linked entity mentions, exclusively of geographical types annotated following their own annotation guidelines [16]. Part of this data has been retained from the original dataset publication and is used and released for the first time for HIPE-2022.

Historical commentaries. The classical commentaries data originates from the *Ajax Multi-Commentary* project and is composed of OCRed 19C commentaries published in French, German and English [17], annotated with both universal NEs (person, location, organisation) and domain-specific NEs (bibliographic references to primary and secondary literature). In the field of classical studies, commentaries constitute one of the most important and enduring forms of scholarship, together with critical editions and translations. They are information-rich texts, characterised by a high density of NEs.

These six datasets compose the HIPE-2022 corpus. They underwent several preparation steps, with conversion to the tab-separated HIPE format, correction of data inconsistencies, metadata consolidation, re-annotation of parts of the datasets, deletion of extremely rare entities (esp. for topres19th), and rearrangement or composition of train and dev splits⁹.

⁶<https://www.newseye.eu/>

⁷<https://sonar.fh-potsdam.de/>

⁸<https://livingwithmachines.ac.uk/>

⁹Additional information is available online by following the links indicated for each datasets in Table 1.

Table 2

Entity types used for NERC tasks, per dataset and with information whether nesting and linking apply. *: these types are not present in letemps data. **: linking applies, unless the token is flagged as *InSecondaryReference*.

Dataset	Coarse tag set	Fine tag set	Nesting	Linking	
hipe2020 letemps	pers	pers.ind	yes	yes	
		pers.coll			
	org*	pers.ind.articleauthor			
		org.adm	yes	yes	
		org.ent			
	prod*	org.ent.pressagency			
		prod.media	no	yes	
	time*	prod.doctr			
		time.date.abs	no	no	
	loc	loc.adm.town	loc.adm.town	yes	yes
			loc.adm.reg		
			loc.adm.nat		
		loc.adm.sup	loc.adm.sup		
			loc.phys.geo	yes	yes
		loc.phys.hydro	loc.phys.hydro		
			loc.phys.astro		
loc.oro		yes	yes		
loc.fac		yes	yes		
loc.add.phys		yes	yes		
loc.add.elec					
loc.unk	no	no			
newseye	pers	pers.articleauthor	yes	yes	
	org	-	yes	yes	
	humanprod	-	yes	yes	
	loc	-	no	yes	
topres19th	loc	-	no	yes	
	building	-	no	yes	
	street	-	no	yes	
ajmc	pers	pers.author	yes	yes**	
		pers.editor			
		pers.myth			
		pers.other			
	work	work.primlit	yes	yes**	
		work.seclit			
		work.fragm			
	loc	-	yes	yes**	
	object	object.manuscr	yes	no	
		object.museum			
date	-	yes	no		
	scope	-	yes	no	
sonar	pers	-	no	yes	
	loc	-	no	yes	
	org	-	no	yes	

3.2. Corpora Characteristics

Overall, the HIPE-2022 corpus covers five languages (English, French, Finnish, German and Swedish), with a total of over 2.3 million tokens (2,211,449 for newspapers and 111,218 for classical commentaries) and 78,000 entities classified according to five different entity typologies and linked to Wikidata records. Detailed statistics about the datasets are provided in Table 3, 4 and 5

The datasets in the corpus are quite heterogeneous in terms of annotation guidelines. Two datasets – `hipe2020` and `letemps` – follow the same guidelines [14, 18], and `newseye` was annotated using a slightly modified version of these guidelines. In the `sonar` dataset, persons, locations and organisations were annotated, whereas in `topres19th` only toponyms were considered. Compared to the other datasets, `ajmc` stands out for having being annotated according to domain-specific guidelines [19], which focus on bibliographic references to primary and secondary literature. This heterogeneity of guidelines leads to a wide variety of entity types and sub-types for the NERC task (see Table 2 and 5). Among these types, only persons, locations and organisations are found in all datasets (except for `topres19th`), thus constituting a set of “universal” entity types. Certain entity types are under-represented in some datasets (e.g. objects, locations and dates in `ajmc`) and, as such, constitute good candidates for the application of data augmentation strategies. Moreover, while nested entities are annotated in all datasets except `topres19th` and `sonar`, only `hipe2020` and `newseye` have a sizable number of such entities.

Detailed information about entity mentions that are affected by OCR mistakes is provided in `ajmc` and `hipe2020` (only for the test set for the latter). As OCR noise constitutes one of the main challenges of historical NE processing [5], this information can be extremely useful to explain differences in performance between datasets or between languages in the same dataset. For instance, looking at the percentage of noisy mentions for the different languages in `ajmc`, we find that it is three times higher in French documents than in the other two languages.

HIPE-2022 datasets show significant differences in terms of lexical overlap between train, dev and test sets. Following the observations of Augenstein et al. [20] and Taillé et al. [21] on the impact of lexical overlap on NERC performance, we computed the percentage of mention overlap between data folds for each dataset, based on the number of identical entity mentions (in terms of surface form) between train+dev and test sets (see Table 6). Evaluation results obtained on training and test sets with low mention overlap, for example, can be taken as an indicator of the ability of the models to generalise well to unseen mentions. We find that `ajmc`, `letemps` and `topres19th` have a mention overlap which is almost twice that of `hipe2020`, `sonar` and `newseye`.

Finally, regarding entity linking, it is interesting to observe that the percentage of NIL entities (i.e. entities not linked to Wikidata) varies substantially across datasets. The Wikidata coverage is drastically lower for `newseye` than for the other newspaper datasets (44.36%). Conversely, only 1.45% of the entities found in `ajmc` cannot be linked to Wikidata. This fact is not at all surprising considering that commentaries mention mostly mythological figures, scholars of the past and literary works, while newspapers mention many relatively obscure or unknown individuals, for whom no Wikidata entry exists.

Table 3

Overview of newspaper corpora statistics (hipe-2022 release v2.1). NIL percentages are computed based on linkable entities (i.e., excluding time entities for hipe2020).

Dataset	Lang.	Fold	Docs	Tokens	Mentions				
					All	Fine	Nested	%noisy	%NIL
hipe2020	de	Train	103	86,446	3,494	3,494	158	-	15.70
		Dev	33	32,672	1,242	1,242	67	-	18.76
		Test	49	30,738	1,147	1,147	73	12.55	17.40
		Total		185	149,856	5,883	5,883	298	-
	en	Train	-	-	-	-	-	-	-
		Dev	80	29,060	966	-	-	-	44.18
		Test	46	16,635	449	-	-	5.57	40.28
	Total		126	45,695	1,415	-	-	-	42.95
	fr	Train	158	166,218	6,926	6,926	473	-	25.26
		Dev	43	37,953	1,729	1,729	91	-	19.81
		Test	43	40,855	1,600	1,600	82	11.25	20.23
		Total		244	245,026	10,255	10,255	646	-
Total			555	440,577	17,553	16,138	944	-	22.82
newseye	de	Train	7	374,250	11,381	21	876	-	51.07
		Dev	12	40,046	539	5	27	-	22.08
		Dev2	12	39,450	882	4	64	-	53.74
		Test	13	99,711	2,401	13	89	-	48.52
		Total		44	553,457	15,203	43	1,056	-
	fi	Train	24	48,223	2,146	15	224	-	40.31
		Dev	24	6,351	223	1	25	-	40.36
		Dev2	21	4,705	203	4	22	-	42.86
		Test	24	14,964	691	7	42	-	47.47
		Total		93	74,243	3,263	27	313	-
	fr	Train	35	255,138	10,423	99	482	-	42.42
		Dev	35	21,726	752	3	29	-	30.45
		Dev2	35	30,457	1,298	10	63	-	38.91
		Test	35	70,790	2,530	34	131	-	44.82
		Total		140	378,111	15,003	146	705	-
	sv	Train	21	56,307	2,140	16	110	-	32.38
		Dev	21	6,907	266	1	7	-	25.19
		Dev2	21	6,987	311	1	20	-	37.30
		Test	21	16,163	604	0	26	-	35.43
		Total		84	86,364	3,321	18	163	-
Total			361	1,092,175	36,790	234	2,237	-	44.36
letemps	fr	Train	414	379,481	9,159	9,159	69	-	-
		Dev	51	38,650	869	869	12	-	-
		Test	51	48,469	1,017	1,017	12	-	-
	Total		516	466,600	11,045	11,045	93	-	-
topres19th	en	Train	309	123,977	3,179	-	-	-	18.34
		Dev	34	11,916	236	-	-	-	13.98
		Test	112	43,263	1,186	-	-	-	17.2
	Total		455	179,156	4,601	-	-	-	17.82
Total			455	179,156	4,601	-	-	-	17.82
sonar	de	Train	-	-	-	-	-	-	-
		Dev	10	17,477	654	-	-	-	22.48
		Test	10	15,464	471	-	-	-	33.33
	Total		20	32,941	1,125	-	-	-	27.02
Total			20	32,941	1,125	-	-	-	27.02
Grand Total (newspapers)			1,907	2,211,449	71,114	27,417	3,274	30.23	

Table 4
Corpus statistics for the ajmc dataset (HIPE-2022 release v2.1).

Dataset	Lang.	Fold	Docs	Tokens	Mentions				
					All	Fine	Nested	%noisy	%NIL
ajmc	de	Train	76	22,694	1,738	1,738	11	13.81	0.92
		Dev	14	4,703	403	403	2	11.41	0.74
		Test	16	4,846	382	382	0	10.99	1.83
	Total		106	32,243	2,523	2,523	13	13.00	1.03
	en	Train	60	30,929	1,823	1,823	4	10.97	1.66
		Dev	14	6,507	416	416	0	16.83	1.70
		Test	13	6,052	348	348	0	10.34	2.61
	Total		87	43,488	2,587	2,587	4	11.83	1.79
	fr	Train	72	24,670	1,621	1,621	9	30.72	0.99
		Dev	17	5,426	391	391	0	36.32	2.56
		Test	15	5,391	360	360	0	27.50	2.80
	Total		104	35,487	2,372	2,372	9	31.16	1.52
Grand Total (ajmc)			297	111,218	7,482	7,482	26		1.45

Table 5
Entity counts by coarse type (HIPE-2022 release v2.1). Although they appear under the same label, identical types present in different data sets may be annotated differently.

		hipe2020			letemps	newseye				sonar	topsres19th	ajmc		
		de	en	fr	fr	de	fi	fr	sv	de	en	de	en	fr
Universal	pers	1849	558	3706	4086	4061	1212	6201	1132	399		910	844	839
	loc	2923	565	4717	6367	6620	1338	5502	1446	477	3727	43	45	24
	org	652	194	1125	592	3584	350	1758	230	249				
Space	building										563			
	street										316			
Time	time	236	46	397										
	date											2	20	5
Man-made	prod	223	52	310										
	humanprod					6620	1338	5502	1446					
	object											12	3	10
	work											465	678	557
	scope											1091	997	937

Table 6

Overlap of mentions between test and train (plus dev) sets as percentage of the total number of mentions.

Dataset	Lang.	% overlap	Folds
ajmc	de	31.43	train+dev vs test
	en	30.50	train+dev vs test
	fr	27.53	train+dev vs test
	Total	29.87	
hipe2020	de	16.22	train+dev vs test
	en	6.22	dev vs test
	fr	19.14	train+dev vs test
	Total	17.12	
letemps	fr	25.70	train+dev vs test
sonar	de	10.13	dev vs test
newseye	fr	14.79	train+dev vs test
	de	20.77	train+dev vs test
	fi	6.63	train+dev vs test
	sv	10.36	train+dev vs test
	Total	16.18	
topres19th	en	32.33	train+dev vs test

3.3. HIPE-2022 Releases

HIPE-2022 data is released as a single package consisting of the neatly structured and homogeneously formatted original datasets. The data is released in IOB format with hierarchical information, similarly to CoNLL-U¹⁰, and consists of UTF-8 encoded, tab-separated values (TSV) files containing the necessary information for all tasks (NERC-Coarse, NERC-Fine, and EL). There is one TSV file per dataset, language and split. Original datasets provide different document metadata with different granularity. This information is kept in the files in the form of metadata blocks that encode as much information as necessary to ensure that each document is self-contained with respect to HIPE-2022 settings. Metadata blocks use namespacing to distinguish between mandatory shared task metadata and dataset-specific metadata.

HIPE-2022 data releases are published on the HIPE-eval GitHub organisation repository¹¹ and on Zenodo¹². Various licences (of type CC-BY and CC-BY-NC-SA) apply to the original datasets – we refer the reader to the online documentation.

¹⁰<https://universaldependencies.org/format.html>

¹¹<https://github.com/impresso/CLEF-HIPE-2020/tree/master/data>

¹²<https://doi.org/10.5281/zenodo.6579950>

4. Evaluation Framework

4.1. Task Bundles, Tracks and Challenges

To accommodate the different dimensions that characterise the HIPE-2022 shared task (languages, document types, entity tag sets, tasks) and to foster research on transferability, the evaluation lab is organised around **tracks** and **challenges**. Challenges guide participation towards the development of approaches that work across settings, e.g. with documents in at least two different languages or annotated according to two different tag sets or guidelines, and provide a well-defined and multi-perspective evaluation frame.

To manage the total combinations of datasets, languages, document types and tasks, we defined the following elements (see also Figure 1):

- **Task bundle:** a task bundle is a predefined set of tasks as in HIPE-2020 (see bundle table in Fig. 1). Task bundles offer participating teams great flexibility in choosing which tasks to compete for, while maintaining a manageable evaluation frame. Concretely, teams were allowed to submit several ‘submission bundles’, i.e. a triple composed of dataset/language/taskbundle, with up to 2 runs each.
- **Track:** a track corresponds to a triple composed of dataset/language/task and forms the basic unit for which results are reported.
- **Challenge:** a challenge corresponds to a predefined set of tracks. A challenge can be seen as a kind of tournament composed of tracks.

HIPE-2022 specifically evaluates 3 challenges:

1. **Multilingual Newspaper Challenge (MNC):** This challenge aims at fostering the development of multilingual NE processing approaches on historical newspapers. The requirements for participation in this challenge are that submission bundles consist only of newspaper datasets and include at least two languages for the same task (so teams had to submit a minimum of two submission bundles for this challenge).
2. **Multilingual Classical Commentary Challenge (MCC):** This challenge aims at adapting NE solutions to domain-specific entities in a specific digital humanities text type of classic commentaries. The requirements are that submission bundles consist only of the `ajmc` dataset and include at least three languages for the same task.
3. **Global Adaptation Challenge (GAC):** Finally, the global adaptation challenge aims at assessing how efficiently systems can be retargeted to any language, document type and guidelines. Bundles submitted for this challenge could be the same as those submitted for MNC and MCC challenges. The requirements are that they consist of datasets of both types (commentaries and newspaper) and include at least two languages for the same task.

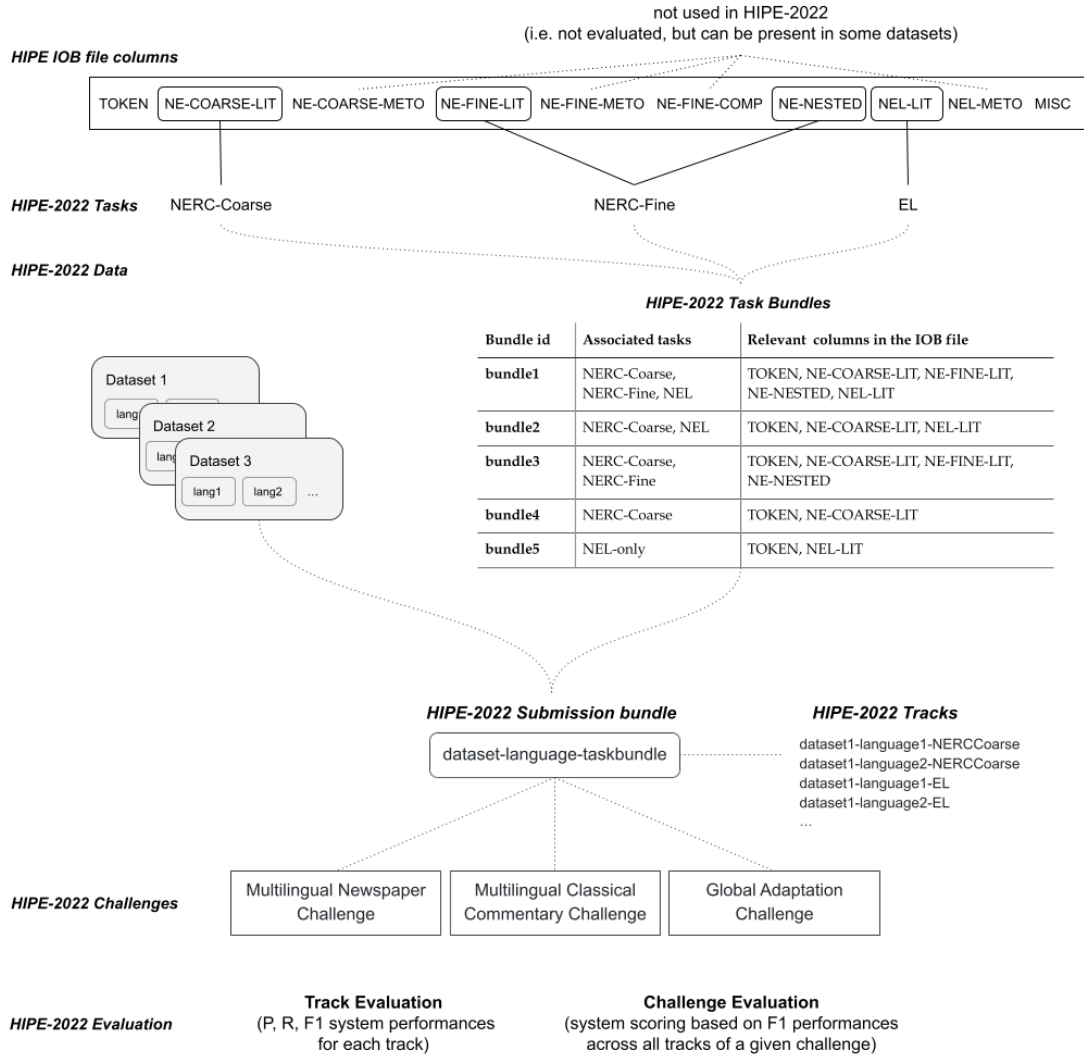


Figure 1: Overview of HIPE-2022 evaluation setting.

4.2. Evaluation Measures

As in HIPE-2020, NERC and EL tasks are evaluated in terms of Precision, Recall and F-measure (F1-score) [22]. Evaluation is carried out at entity level according to two computation schemes: micro average, based on true positives, false positives, and false negative figures computed over all documents, and macro average, based on averages of micro figures per document. Our definition of macro differs from the usual one: averaging is done at document level and not at entity type level. This allows to account for variance in document length and entity distribution within documents and avoids distortions that would occur due to the unevenly distributed entity classes.

Both NERC and EL benefit from strict and fuzzy evaluation regimes, depending on how

strictly entity type and boundaries correctness are judged. For NERC (Coarse and Fine), the strict regime corresponds to exact type and boundary matching, and the fuzzy to exact type and overlapping boundaries. It is to be noted that in the strict regime, predicting wrong boundaries leads to a ‘double’ punishment of one false negative (entity present in the gold standard but not predicted by the system) and one false positive (entity predicted by the system but not present in the gold standard). Although it penalizes harshly, we keep this metric to be consistent with CoNLL and refer to the fuzzy regime when boundaries are of less importance.

The definition of strict and fuzzy regimes differs for entity linking. In terms of boundaries, EL is always evaluated according to overlapping boundaries in both regimes (what is of interest is the capacity to provide the correct link rather than the correct boundaries). EL strict regime considers only the system’s top link prediction (NIL or Wikidata QID), while the fuzzy regime expands system predictions with a set of historically related entity QIDs. For example, “Germany” QID is complemented with the QID of the more specific “Confederation of the Rhine” entity and both are considered as valid answers. The resource allowing for such historical normalization was compiled by the task organizers for the entities of the test data sets (for *hipe2020* and *ajmc* datasets), and are released as part of the HIPE-scorer. For this regime, participants were invited to submit more than one link, and F-measure is additionally computed with cut-offs @3 and @5 (meaning, counting a true positive if the ground truth QID can be found within the first 3 or 5 candidates).

4.3. System Evaluation, Scorer and Evaluation Toolkit

Teams were asked to submit system responses based on submission bundles and to specify at least one challenge to which their submitted bundles belong. Micro and macro scores were computed and published for each track, but only micro figures are reported here.

The evaluation of challenges, which corresponds to an aggregation of tracks, was defined as follows: given a specific challenge and the tracks submitted by a team for this challenge, the submitted systems are rewarded points according to their F1-based rank for each track (considering only the best of the submitted runs for a given track). The points obtained are summed over all submitted tracks, and systems/teams are ranked according to their total points. Further details on system submission and evaluation can be found in the HIPE Participation Guidelines [23].

The evaluation is performed using the **HIPE-scorer**¹³. Developed during the first edition of HIPE, the scorer has been improved with minor bug fixes and additional parameterisation (input format, evaluation regimes, HIPE editions). Participants could use the HIPE-scorer when developing their systems. After the evaluation phase, a complete **evaluation toolkit** was also released, including the data used for evaluation (v2.1), the system runs submitted by participating teams, and all the evaluation recipes and resources (e.g. historical mappings) needed to replicate the present evaluation¹⁴.

¹³<https://github.com/hipe-eval/HIPE-scorer>

¹⁴<https://github.com/hipe-eval/HIPE-2022-eval>

5. System Descriptions

In this second HIPE edition, 5 teams submitted a total of 103 system runs. Submitted runs do not cover all of the 35 possible tracks (dataset/language/task combinations), nevertheless we received submission for all datasets, with most of them focusing on NERC-Coarse.

5.1. Baselines

As a neural baseline (**NEUR-BSL**) for NERC-Coarse and NERC-fine, we fine-tuned separately for each HIPE-2022 dataset XLM-R_{BASE}, a multilingual transformer-based language representation model pre-trained on 2.5TB of filtered CommonCrawl texts [24]. The models are implemented using HuggingFace¹⁵ [25]. Since transformers rely primarily on subword-based tokenisers, we chose to label only the first subwords. This allows to map the model outputs to the original text more easily. Tokenised texts are split into input segments of length 512. For each HIPE-2022 dataset, fine-tuning is performed on the train set (except for sonar and hipecoarse-en which has only dev sets) for 10 epochs using the default hyperparameters (Adam $\epsilon = 10e-8$, Learning rate $\alpha = 5e-5$). The code of this baseline (configuration files, scripts) is published in a dedicated repository on the HIPE-eval Github organisation¹⁶, and results are published in the evaluation toolkit.

For entity linking in EL-only setting, we provide the NIL baseline (**NIL-BSL**), where each entity link is replaced with the NIL value.

5.2. Participating Systems

The following system descriptions are compiled from information provided by the participants. More details on the implementation and results can be found in the system papers of the participants [26].

Team **L3I**, affiliated with *La Rochelle University* and with the *University of Toulouse*, France, successfully tackled an impressive amount of multilingual newspaper datasets with strong runs for NERC-coarse, NERC-fine and EL. For the classical commentary datasets (ajmc) the team had excellent results for NERC¹⁷. For NERC, L3I – the winning team in HIPE’s 2020 edition – builds on their transformer-based approach [27]. Using transformer-based adapters [28], parameter-efficient fine-tuning in a hierarchical multitask setup (NERC-coarse and NERC-fine) has been shown to work well with historical noisy texts [27]. The innovation for this year’s submission lies in the addition of context information in the form of external knowledge from two sources (inspired by [29]). First, French and German Wikipedia documents based on dense vector representations computed by a multilingual Sentence-BERT model [30], including a k-Nearest-Neighbor search functionality provided by ElasticSearch framework. Second, English Wikidata knowledge graph (KG) embeddings that are combined with the first paragraph of English Wikipedia pages (Wikidata5m) [31]. For the knowledge graph embeddings, two methods are tested on the HIPE-2022 data: 1) the one-stage KG Embedding Retrieval Module that retrieves

¹⁵<https://github.com/huggingface/transformers/>

¹⁶<https://github.com/hipecoarse/HIPE-2022-baseline/>

¹⁷The EL results for ajmc were low, probably due to some processing issues.

top- k KG “documents” (in this context, a document is an Elasticsearch retrieval unit that consists of an entity identifier, an entity description and an entity embedding) via vector similarity on the dense entity embedding vector space; 2) the two-stage KG Embedding Retrieval Module that retrieves the single top similar document first and then in a second retrieval step gets the k most similar documents based on that first entity. All context enrichment techniques work by simply concatenating the original input segment with the retrieved context segments and processing the contextualized segments through their “normal” hierarchical NER architecture. Since the L3i team’s internal evaluation on HIPE-2022 data (using a multilingual BERT base pre-trained model) indicated that the two-stage KG retrieval was the best context generator overall, it was used for one of the two officially submitted runs. The other “baseline” run did not use any context enrichment techniques. Both runs additionally used stacked monolingual BERT embeddings for English, French and German, for the latter two languages in the form of Europeana models that were built from digitized historical newspaper text material. Even with improved historical monolingual BERT embeddings, the context-enriched run was consistently better in terms of F1-score in NERC-Coarse and -Fine settings.

Team **histeria**, affiliated with the *Bayerische Staatsbibliothek München*, Germany, the *Digital Philology* department of the University of Vienna, Austria and the *NLP Expert Center, Volkswagen AG Munich*, Germany, focused on the *ajmc* dataset for their NERC-coarse submission (best results for French and English, second best for German), but also provides experimental results for all languages of the *newseye* datasets¹⁸. Their NER tagging experiments tackle two important questions:

a) How to build an optimal multilingual pre-trained BERT language representation model for historical OCRized documents? They propose and release *hmBERT*¹⁹, which includes English, Finnish, French, German and Swedish in various model and vocabulary sizes, and specifically apply methods to deal with OCR noise and imbalanced corpus sizes per language. In the end, roughly 27GB of text per language is used in pre-training.

b) How to fine-tune a multilingual pre-trained model given comparable NER annotations in multiple languages? They compare a single-model approach (training models separately for each language) with a one-model approach (training only one model that covers all languages). The results indicate that, most of the time, the single-model approach works slightly better, but the difference may not be large enough to justify the considerably greater effort to train and apply the models in practice.

histeria submitted two runs for each *ajmc* datasets, using careful hyperparameter grid search on the dev sets in the process. Both runs build on the one-model approach in a first multilingual fine-tuning step. Similar to [29], they build monolingual models by further fine-tuning on language-specific training data²⁰. Run 1 of their submission is based on *hmBERT* with vocabulary size 32k, while run 2 has a vocabulary size of 64k. Somewhat unexpectedly, the larger vocabulary does not improve the results in general on the development set. For

¹⁸Note that these experiments are evaluated using the officially published *Newseye* test sets [15] (released as dev2 dataset as part of HIPE-2022) and not the HIPE-2022 *newseye* test sets, which were unpublished prior to the HIPE 2022 campaign.

¹⁹For English data, they used the Digitised Books. c. 1510 - c. 1900, all other languages use Europeana newspaper text data.

²⁰This improves the results by 1.2% on average on the HIPE-2022 data.

the test set, though, the larger vocabulary model is substantially better overall. Similar to the team L3I, HISTERIA also experimented with context enrichment techniques suggested by [29]. However, for the specific domain of classical commentaries, general-purpose knowledge bases such as Wikipedia could not improve the results. Interestingly, L3I also observed much less improvement with Wikipedia context enrichment on *ajmc* in comparison to the *hipe2020* newspaper datasets. In summary, HISTERIA outperformed the strong neural baseline by about 10 F1-score percentage points in strict boundary setting, thereby demonstrating the importance of carefully constructed domain-specific pre-trained language representation models.

Team **AAUZH**, affiliated with *University of Zurich*, Switzerland and *University of Milan*, Italy, focused on the multilingual newspaper challenge in NERC-coarse setting and experimented with 21 different monolingual and multilingual, as well as contemporary and historical transformer-based language representation models available on the HuggingFace platform. For fine-tuning, they used the standard token classification head of the transformer library for NER tagging with default hyperparameters and trained each dataset for 3 epochs. In a preprocessing step, token-level NER IOB labels were mapped onto all subtokens. At inference time, a simple but effective summing pooling strategy for NER for aggregating subtoken-level to token-level labels was used [32]. Run 2 of AAUZH are the predictions of the best single model. Run 1 is the result of a hard-label ensembling from different pre-trained models: in case of ties between O and B/I labels, the entity labels were preferred. The performance of the submitted runs varies strongly in comparison with the neural baseline: for German and English it generally beats the baseline clearly for *hipe2020* and *sonar* datasets, but suffers on French *hipe2020* and German/Finnish *newseye* datasets. This again indicates that in transfer learning approaches to historical NER, the selection of pre-trained models has a considerable impact. The team also performed some post-submission experiments to investigate the effect of design choices: Applying soft-label ensembling using averaged token-level probabilities turned out to improve results on the French *newseye* datasets by 1.5 percentage point in micro average and 2.4 points in macro average (F1-score). For all languages of the *newseye*, they also tested a one-model approach with multilingual training. The best multilingual *dbmdz* Europeana BERT model had a better performance on average (58%) than the best monolingual models (56%). However, several other multilingual pre-trained language models had substantially worse performance, resulting in 57% ensemble F1-score (5 models), which was much lower than 67% achieved by the monolingual ensemble.

Team **SBB**, affiliated with the *Berlin State Library*, Germany, participated exclusively in the EL-only subtask, but covered all datasets in English, German and French. Their system builds on models and methods developed in the HIPE-2020 edition [33]. Their approach uses Wikipedia sentences with an explicit link to a Wikipedia page as textual representations of its connected Wikidata entity. The system makes use of the metadata of the HIPE-2022 documents to exclude entities that were not existing at the time of its publication. Going via Wikipedia reduces the amount of accessible Wikidata IDs, however, for all datasets but *ajmc* the coverage is still 90%. Given the specialised domain of *ajmc*, a coverage of about 55% is to be accepted. The entity linking is done in the following steps: a) A candidate lookup retrieves a given number of candidates (25 for submission run 1, 50 for submission run 2) using a nearest neighbour index based on word embeddings of Wikipedia page titles. An absolute cut-off value is used to limit the retrieval (0.05 for submission 1 and 0.13 for submission 2). b) A probabilistic candidate

sentence matching is performed by pairwise comparing the sentence with the mention to link and a knowledge base text snippet. To this end, a BERT model was fine-tuned on the task of whether or not two sentences mention the same entity. c) The final ranking of candidates includes the candidate sentence matching information as well as lookup features from step (a) and more word embedding information from the context. A random forest model calculates the overall probability of a match between the entity mention and an entity linking candidate. If the probability of a candidate is below a given threshold (0.2 for submission run 1 and 2), it is discarded. The random forest model was trained on concatenated training sets of the same language across datasets.

There are no conclusive insights from HIPE-2022 EL-only results whether run 1 or 2 settings are preferable. Post-submission experiments in their system description paper investigate the influence of specific hyperparameter settings on the system performances.

6. Results and Discussion

We report results for the best run of each team and consider micro Precision, Recall and F1-score exclusively. Results for NERC-Coarse and NERC-Fine for all languages and datasets according to both evaluation regimes are presented in Table 7 and 8 respectively. Table 10 reports performances for EL-only, with a cut-off @1. We refer the reader to the HIPE-2022 website and the evaluation toolkit for more detailed results²¹, and to the extended overview paper for further discussion of the results [12].

6.1. General Observations

All systems now use transformer-based approaches with strong pre-trained models. The choice of the pre-trained model – and the corresponding text types used in pre-training – have a strong influence on performance.

The quality of available multilingual pre-trained models for fine-tuning on NER tasks proved to be competitive compared to training individual monolingual models. However, to get the maximum performance out of it, the multilingual fine-tuning in a first phase must be complemented by a monolingual second phase.

NERC. In general, the systems demonstrated a good ability to adapt to heterogeneous annotation guidelines. They achieved their highest F1-scores for the NERC-Coarse task on *ajmc*, a dataset annotated with domain-specific entities and of relatively small size compared to the newspaper datasets, thus confirming the ability of strong pre-trained models to achieve good results when fine-tuned on relatively small datasets. The good results obtained on *ajmc*, however, may be partly due to the relatively high mention overlap between train and test sets (see Section 3.2). Moreover, it is worth noting that performances on the French subset of the *ajmc* dataset do not substantially degrade despite the high rate of noisy mentions (three times higher than English and German), which shows a good resilience of transformer-based models to OCR noise on this specific dataset.

²¹See <https://hipe-eval.github.io/HIPE-2022> and <https://github.com/hipe-eval/HIPE-2022-eval>

Table 7

Results for NERC-Coarse (micro P, R and F1-score). Bold font indicates the highest, and underlined font the second-highest value.

	Strict			Fuzzy			Strict			Fuzzy			Strict			Fuzzy		
	P	R	F	P	R	F	P	R	F	P	R	F	P	R	F	P	R	F
hipe2020																		
French						German						English						
AAUZH	.718	.675	.696	.825	.776	.800	.716	.735	<u>.725</u>	.812	.833	<u>.822</u>	.538	.490	<u>.513</u>	.726	.661	<u>.692</u>
L3i	.786	.831	.808	.883	.933	.907	.784	.805	.794	.865	.888	.876	.624	.617	.620	.793	.784	.788
NEUR-BSL	.730	.785	<u>.757</u>	.836	.899	<u>.866</u>	.665	.746	.703	.750	.842	.793	.432	.532	.477	.564	.695	.623
Letemps																		
French						German						English						
AAUZH	.589	.710	<u>.644</u>	.642	.773	<u>.701</u>	.512	.548	.529	.655	.741	.695	.816	.760	.787	.869	.810	.838
NEUR-BSL	.595	.744	.661	.639	.800	.711	.267	.361	.307	.410	.554	.471	.747	.782	<u>.764</u>	.798	.836	<u>.816</u>
sonar																		
French						German						English						
HISTeria	.834	.850	.842	.874	.903	.888	.930	.898	.913	.938	.953	.945	.826	.885	.854	.879	.943	.910
L3i	.810	.842	.826	.856	.889	.872	.946	.921	.934	.965	.940	.952	.824	.876	.850	.868	.922	.894
NEUR-BSL	.707	.778	.741	.788	.867	.825	.792	.846	.818	.846	.903	.873	.680	.802	.736	.766	.902	.828
ajmc																		
French						German						English						
AAUZH	.655	.657	.656	.785	.787	.786	.395	.421	.408	.480	.512	.495						
NEUR-BSL	.634	.676	.654	.755	.805	.779	.429	.537	.477	.512	.642	.570						
newseye																		
French						German						English						
AAUZH	.655	.657	.656	.785	.787	.786	.395	.421	.408	.480	.512	.495						
NEUR-BSL	.634	.676	.654	.755	.805	.779	.429	.537	.477	.512	.642	.570						
ajmc																		
Finnish						Swedish						English						
AAUZH	.618	.524	.567	.730	.619	.670	.686	.604	.643	.797	.702	.746						
NEUR-BSL	.605	.687	.644	.715	.812	.760	.588	.728	.651	.675	.836	.747						

Table 8

Results for NERC-Fine and Nested (micro P, R and F1-score).

	French						German						English					
	Strict			Fuzzy			Strict			Fuzzy			Strict			Fuzzy		
	P	R	F	P	R	F	P	R	F	P	R	F	P	R	F	P	R	F
hipe2020 (Fine)																		
L3i	.702	.782	.740	.784	.873	.826	.691	.747	.718	.776	.840	.807	-	-	-	-	-	-
NEUR-BSL	.685	.733	.708	.769	.822	.795	.584	.673	.625	.659	.759	.706	-	-	-	-	-	-
hipe2020 (Nested)																		
L3i	.390	.366	.377	.416	.390	.403	.714	.411	.522	.738	.425	.539	-	-	-	-	-	-
ajmc (Fine)																		
L3i	.646	.694	.669	.703	.756	.728	.915	.898	.906	.941	.924	.933	.754	.848	.798	.801	.899	.847
NEUR-BSL	.526	.567	.545	.616	.664	.639	.819	.817	.818	.866	.864	.865	.600	.744	.664	.676	.839	.749

EL-only. Entity linking on already identified mentions appears to be considerably more challenging than NERC, with F1-scores varying considerably across datasets. The linking of toponyms in topres19th is where systems achieved the overall best performances. Conversely, EL-only on historical commentaries (ajmc) appears to be the most difficult, with the lowest

Table 9

NERC-Coarse **recall** by entity type. Figures are from the same runs reported in the generic NERC-Coarse table (although recall scores from another run submitted by a team may be better in some cases.)

Dataset	Lang.	Type	NEUR-BSL	AAUZH	L3i	HISTeria
hipe2020	fr	pers	.934	.882	.970	-
		loc	.930	.760	.954	-
		org	.593	.580	.724	-
		prod	.603	.721	.929	-
		time	.810	.739	.803	-
	de	pers	.936	.884	.974	-
		loc	.874	.881	.929	-
		org	.407	.560	.645	-
		prod	.579	.607	.722	-
		time	.882	.930	.906	-
	en	pers	.872	.846	.872	-
		loc	.867	.646	.790	-
		org	.128	.478	.637	-
		prod	-	-	.722	-
		time	.538	.938	1.000	-
letemps	fr	pers	.841	.816	-	-
		loc	.832	.827	-	-
newseye	de	pers	.529	.486	-	-
		loc	.658	.561	-	-
		org	.338	.330	-	-
		humanprod	-	.100	-	-
	fi	pers	.863	.899	-	-
		loc	.836	.698	-	-
		org	.324	.491	-	-
		humanprod	.429	.789	-	-
	fr	pers	.889	.888	-	-
		loc	.785	.818	-	-
		org	.539	.540	-	-
		humanprod	.600	.704	-	-
sv	pers	.684	.765	-	-	
	loc	.898	.751	-	-	
	org	.326	.533	-	-	
	humanprod	.571	.889	-	-	
sonar	de	pers	.394	.656	-	-
		loc	.808	.836	-	-
		org	.128	.496	-	-
topres19th	en	loc	.892	.860	-	-
		building	.600	.738	-	-
		street	.661	.685	-	-
ajmc	de	pers	.953	-	.938	.930
		loc	-	-	1.000	.500
		work	.729	-	.958	.985
		scope	.861	-	.939	.950
	en	pers	.821	-	.901	.945
		loc	-	-	.800	1.000
		work	.926	-	.947	.947
		scope	.887	-	.921	.921
	fr	pers	.878	-	.928	.906
		loc	-	-	.444	.556
		work	.700	-	.831	.910
		scope	.876	-	.883	.879

the correctly predicted entity links (true positives) correspond to abbreviated mentions, which nevertheless represent about 47% of all linkable mentions.

6.2. Observations on Challenges

A complementary way of looking at the results is to consider them in light of the specific challenges raised by historical documents. Such challenges are one of the novelties of HIPE-2022 and were introduced as “thematic aggregations” of system rankings across datasets and languages (see Section 4).²²

Multilingual Newspaper Challenge (MNC). Overall, four teams participated in this challenge: three teams in the NERC-Coarse task, two in EL-only and one team in end-to-end EL. The top-ranked team, AAUZH, was able to tackle five different newspaper datasets in five languages in total, with performances above the strong neural baseline in 6 out of 10 cases. However, it should be noted that they used two different systems, one based on the best fine-tuned model (selected by development set performance) and another one being an ensemble of all their fine-tuned models. Thus, despite leading in the MNC challenge, their work does not answer the question of which single system works best across datasets and in different languages. Conversely, the second-ranked system by L3I team submitted for fewer datasets and languages, but showed an overall higher quality of predictions both for NERC-Coarse and EL-only for languages they covered. It would be interesting to see how this system performs on the remaining languages and datasets, especially in comparison with the baseline. In general, one aspect of the MNC challenge which remained unexplored is entity linking beyond a standard set of languages such as English, German and French, as no runs for EL-only on Finnish and Swedish newspapers were submitted.

Multilingual Classical Commentary Challenge (MCC). This challenge had in total three participants: two teams participated in this challenge in the NERC-Coarse task, and one team participated in the end-to-end EL and EL-only. Given the overlap between MCC and GAC challenges in terms of languages and datasets, teams participating in the latter were also considered for the former.

The NERC-Coarse results showed how a BERT-based multilingual model pre-trained on large corpora of historical documents and fine-tuned on domain-specific data performs better or on-par with a system implementing a more complex transformer-based architecture. An interesting insight which emerged from this challenge is that methods employing context enrichment techniques which rely on lexical information from Wikipedia do not yield performance improvements as they do when applied to other document types, such as newspapers.

Regarding EL, MCC exemplified well some characteristics an EL system needs to have for it to be applied successfully across different domains. In particular, assumptions about which entities are to be retained when constructing a knowledge base for this task need to be relaxed. Moreover, an aspect that emerged from this challenge and will deserve more research in the future is the linking of abbreviated entities (e.g. mentions of literary works in commentaries), which proved to be challenging for participating systems.

²²The challenge leaderboards can be found at the HIPE-2022 results page in the *Challenge Evaluation Results section*, see <https://hipe-eval.github.io/HIPE-2022/results>.

Global Adaptation Challenge (GAC). The high level of difficulty entailed by this challenge is reflected by the number of participants: one team for the NERC-Coarse task and one team for end-to-end EL and EL-only. Both teams had already participated in the first edition of HIPE with very good results, and tackled this year the challenge of adapting their systems to work in a multi-domain and multilingual scenario. In general, the results of this challenge confirm that the systems proposed by L3I and SBB respectively for the tasks of NERC-Coarse and EL are suitable to be applied to data originating from heterogeneous domains. They also show that EL across languages and domains remains a more challenging task than NERC, calling for more future research on this topic. Moreover, no team has worked on adapting annotation models to be able to combine different NER training datasets with sometimes incompatible annotations and benefit from a larger dataset overall. This data augmentation strategy to global adaptation, which could be beneficial for underrepresented entity types (e.g. dates or locations in the *ajmc* dataset), remains to be explored in future work.

7. Conclusion and Perspectives

From the perspective of natural language processing, this second edition of HIPE provided the possibility to test the robustness of existing approaches and to experiment with transfer learning and domain adaptation methods, whose performances could be systematically evaluated and compared on broad historical and multilingual data sets. Besides gaining new insights with respect to domain and language adaptation and advancing the state of the art in semantic indexing of historical material, the lab also contributed an unprecedented set of multilingual and historical NE-annotated datasets that can be used for further experimentation and benchmarking.

From the perspective of digital humanities, the lab's outcomes will help DH practitioners in mapping state-of-the-art solutions for NE processing of historical texts, and in getting a better understanding of what is already possible as opposed to what is still challenging. Most importantly, digital scholars are in need of support to explore the large quantities of digitised text they currently have at hand, and NE processing is high on the agenda. Such processing can support research questions in various domains (e.g. history, political science, literature, historical linguistics) and knowing about their performance is crucial in order to make an informed use of the processed data.

From the perspective of cultural heritage professionals, who increasingly focus on advancing the usage of artificial intelligence methods on cultural heritage text collections [34, 35], the HIPE-2022 shared task and datasets represent an excellent opportunity to experiment with multilingual and multi-domain data of various quality and annotation depth, a setting close to the real-world scenarios they are often confronted with.

Overall, HIPE-2022 has contributed to further advance the state of the art in semantic indexing of historical documents. By expanding the language spectrum and document types and integrating datasets with various annotation tag sets, this second edition has set the bar high, and there remains much to explore and experiment.

8. Acknowledgments

The HIPE-2022 team expresses her greatest appreciation to the HIPE-2022 partnering projects, namely AjMC, impresso-HIPE-2020, *Living with Machines*, NewsEye, and SoNAR, for contributing (and hiding) their NE-annotated datasets. We particularly thank Mariona Coll-Ardanuy (LwM), Ahmed Hamdi (NewsEye) and Clemens Neudecker (SoNAR) for their support regarding data provision, and the members of the HIPE-2022 advisory board, namely Sally Chambers, Frédéric Kaplan and Clemens Neudecker.

References

- [1] M. Ehrmann, M. Romanello, S. Najem-Meyer, A. Doucet, S. Clemenide, Overview of HIPE-2022: Named Entity Recognition and Linking in Multilingual Historical Documents, in: A. Barrón-Cedeño, G. Da San Martino, M. Degli Esposti, F. Sebastiani, C. Macdonald, G. Pasi, A. Hanbury, M. Potthast, G. Faggioli, N. Ferro (Eds.), *Experimental IR Meets Multilinguality, Multimodality, and Interaction. Proceedings of the Thirteenth International Conference of the CLEF Association (CLEF 2022)*, Lecture Notes in Computer Science (LNCS), Springer, 2022.
- [2] F. Kaplan, I. di Lenardo, Big Data of the Past, *Frontiers in Digital Humanities* 4 (2017) 1–21. URL: <https://www.frontiersin.org/articles/10.3389/fdigh.2017.00012/full>. doi:10.3389/fdigh.2017.00012, publisher: Frontiers.
- [3] M. Ridge, G. Colavizza, L. Brake, M. Ehrmann, J.-P. Moreux, A. Prescott, The past, present and future of digital scholarship with newspaper collections, in: *DH 2019 book of abstracts*, Utrecht, The Netherlands, 2019, pp. 1–9. URL: <http://infoscience.epfl.ch/record/271329>.
- [4] M. Ehrmann, G. Colavizza, Y. Rochat, F. Kaplan, Diachronic evaluation of NER systems on old newspapers, in: *Proceedings of the 13th Conference on Natural Language Processing (KONVENS 2016)*, Bochumer Linguistische Arbeitsberichte, Bochum, 2016, pp. 97–107. URL: <https://infoscience.epfl.ch/record/221391>.
- [5] M. Ehrmann, A. Hamdi, E. L. Pontes, M. Romanello, A. Doucet, Named Entity Recognition and Classification on Historical Documents: A Survey, arXiv:2109.11406 [cs] (2021 (to appear in ACM journal *Computing Surveys* in 2022)). URL: <http://arxiv.org/abs/2109.11406>, arXiv: 2109.11406.
- [6] M. Ehrmann, M. Romanello, A. Flückiger, S. Clemenide, Overview of CLEF HIPE 2020: Named entity recognition and linking on historical newspapers, in: A. Arampatzis, E. Kanoulas, T. Tsirikika, S. Vrochidis, H. Joho, C. Lioma, C. Eickhoff, A. Névéal, L. Cappellato, N. Ferro (Eds.), *Experimental IR Meets Multilinguality, Multimodality, and Interaction*, Lecture Notes in Computer Sciences, Springer International Publishing, Cham, 2020, pp. 288–310.
- [7] M. Ehrmann, M. Romanello, A. Flückiger, S. Clemenide, Extended Overview of CLEF HIPE 2020: Named Entity Processing on Historical Newspapers, in: L. Cappellato, C. Eickhoff, N. Ferro, A. Névéal (Eds.), *Working Notes of CLEF 2020 - Conference and Labs of the Evaluation Forum*, volume 2696, CEUR-WS, Thessaloniki, Greece, 2020, p. 38. URL: <https://infoscience.epfl.ch/record/281054>. doi:10.5281/zenodo.4117566.

- [8] G. Beryozkin, Y. Drori, O. Gilon, T. Hartman, I. Szpektor, A joint named-entity recognizer for heterogeneous tag-sets using a tag hierarchy, in: Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics, Florence, Italy, 2019, pp. 140–150. URL: <https://aclanthology.org/P19-1014>.
- [9] J. Li, B. Chiu, S. Feng, H. Wang, Few-shot named entity recognition via meta-learning, IEEE Transactions on Knowledge and Data Engineering (2020) 1–1.
- [10] Q. Wu, Z. Lin, G. Wang, H. Chen, B. F. Karlsson, B. Huang, C. Lin, Enhanced meta-learning for cross-lingual named entity recognition with minimal resources, CoRR abs/1911.06161 (2019). URL: <http://arxiv.org/abs/1911.06161>.
- [11] J. Li, S. Shang, L. Shao, Metaner: Named entity recognition with meta-learning, in: Proceedings of The Web Conference 2020, WWW '20, Association for Computing Machinery, New York, NY, USA, 2020, p. 429–440. URL: <https://doi.org/10.1145/3366423.3380127>. doi:10.1145/3366423.3380127.
- [12] M. Ehrmann, M. Romanello, S. Najem-Meyer, A. Doucet, S. Clemenide, Extended Overview of HIPE-2022: Named Entity Recognition and Linking on Multilingual Historical Documents, in: G. Faggioli, N. Ferro, A. Hanbury, M. Potthast (Eds.), Working Notes of CLEF 2022 - Conference and Labs of the Evaluation Forum, CEUR-WS, 2022.
- [13] M. Ehrmann, M. Romanello, S. Clemenide, P. B. Ströbel, R. Barman, Language Resources for Historical Newspapers: The *Impresso* Collection, in: Proceedings of the 12th Language Resources and Evaluation Conference, European Language Resources Association, Marseille, France, 2020, pp. 958–968.
- [14] M. Ehrmann, M. Romanello, A. Flückiger, S. Clemenide, *Impresso* Named Entity Annotation Guidelines, Annotation Guidelines, Ecole Polytechnique Fédérale de Lausanne (EPFL) and Zurich University (UZH), 2020. URL: <https://zenodo.org/record/3585750>. doi:10.5281/zenodo.3604227.
- [15] A. Hamdi, E. Linhares Pontes, E. Boros, T. T. H. Nguyen, G. Hackl, J. G. Moreno, A. Doucet, A Multilingual Dataset for Named Entity Recognition, Entity Linking and Stance Detection in Historical Newspapers, in: Proceedings of the 44th International ACM SIGIR Conference on Research and Development in Information Retrieval, SIGIR '21, Association for Computing Machinery, New York, NY, USA, 2021, pp. 2328–2334. doi:10.1145/3404835.3463255.
- [16] M. Coll Ardanuy, D. Beavan, K. Beelen, K. Hosseini, J. Lawrence, Dataset for Toponym Resolution in Nineteenth-Century English Newspapers, 2021. URL: <https://doi.org/10.23636/b1c4-py78>. doi:10.23636/b1c4-py78.
- [17] M. Romanello, N.-M. Sven, B. Robertson, Optical Character Recognition of 19th Century Classical Commentaries: the Current State of Affairs, in: The 6th International Workshop on Historical Document Imaging and Processing (HIP '21), Association for Computing Machinery, Lausanne, 2021. URL: <https://doi.org/10.1145/3476887.3476911>. doi:10.1145/3476887.3476911.
- [18] S. Rosset, Grouin, Cyril, Zweigenbaum, Pierre, Entités Nommées Structurées : Guide d'annotation Quaero, Technical Report 2011-04, LIMSI-CNRS, Orsay, France, 2011.
- [19] M. Romanello, S. Najem-Meyer, Guidelines for the Annotation of Named Entities in the Domain of Classics, 2022. URL: <https://doi.org/10.5281/zenodo.6368101>. doi:10.5281/zenodo.6368101.
- [20] I. Augenstein, L. Derczynski, K. Bontcheva, Generalisation in named entity recognition: A

- quantitative analysis, *Computer Speech & Language* 44 (2017) 61–83. URL: <http://www.sciencedirect.com/science/article/pii/S088523081630002X>. doi:10.1016/j.csl.2017.01.012.
- [21] B. Taillé, V. Guigue, P. Gallinari, Contextualized Embeddings in Named-Entity Recognition: An Empirical Study on Generalization, in: J. M. Jose, E. Yilmaz, J. Magalhães, P. Castells, N. Ferro, M. J. Silva, F. Martins (Eds.), *Advances in Information Retrieval, Lecture Notes in Computer Science*, Springer International Publishing, Cham, 2020, pp. 383–391. doi:10.1007/978-3-030-45442-5_48.
- [22] J. Makhoul, F. Kubala, R. Schwartz, R. Weischedel, Performance measures for information extraction, in: *In Proceedings of DARPA Broadcast News Workshop*, 1999, pp. 249–252.
- [23] M. Ehrmann, M. Romanello, A. Doucet, S. Clematide, HIPE 2022 Shared Task Participation Guidelines, Technical Report, Zenodo, 2022. URL: <https://zenodo.org/record/6045662>. doi:10.5281/zenodo.6045662.
- [24] A. Conneau, K. Khandelwal, N. Goyal, V. Chaudhary, G. Wenzek, F. Guzmán, E. Grave, M. Ott, L. Zettlemoyer, V. Stoyanov, Unsupervised Cross-lingual Representation Learning at Scale, 2020. [arXiv:1911.02116](https://arxiv.org/abs/1911.02116).
- [25] T. Wolf, L. Debut, V. Sanh, J. Chaumond, C. Delangue, A. Moi, P. Cistac, T. Rault, R. Louf, M. Funtowicz, J. Davison, S. Shleifer, P. von Platen, C. Ma, Y. Jernite, J. Plu, C. Xu, T. L. Scao, S. Gugger, M. Drame, Q. Lhoest, A. M. Rush, Transformers: State-of-the-art natural language processing, in: *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing: System Demonstrations*, Association for Computational Linguistics, Online, 2020, pp. 38–45. URL: <https://www.aclweb.org/anthology/2020.emnlp-demos.6>.
- [26] G. Faggioli, N. Ferro, A. Hanbury, M. Potthast (Eds.), *Working Notes of CLEF 2022 - Conference and Labs of the Evaluation Forum*, CEUR-WS, 2022.
- [27] E. Boros, A. Hamdi, E. Linhares Pontes, L. A. Cabrera-Diego, J. G. Moreno, N. Sidere, A. Doucet, Alleviating digitization errors in named entity recognition for historical documents, in: *Proceedings of the 24th Conference on Computational Natural Language Learning*, Association for Computational Linguistics, Online, 2020, pp. 431–441. doi:10.18653/v1/2020.conll-1.35.
- [28] N. Houlsby, A. Giurgiu, S. Jastrzebski, B. Morrone, Q. De Laroussilhe, A. Gesmundo, M. Attariyan, S. Gelly, Parameter-efficient transfer learning for NLP, in: K. Chaudhuri, R. Salakhutdinov (Eds.), *Proceedings of the 36th International Conference on Machine Learning*, volume 97 of *Proceedings of Machine Learning Research*, PMLR, 2019, pp. 2790–2799. URL: <https://proceedings.mlr.press/v97/houlsby19a.html>.
- [29] X. Wang, Y. Shen, J. Cai, T. Wang, X. Wang, P. Xie, F. Huang, W. Lu, Y. Zhuang, K. Tu, W. Lu, Y. Jiang, DAMO-NLP at SemEval-2022 task 11: A knowledge-based system for multilingual named entity recognition, 2022. URL: <https://arxiv.org/abs/2203.00545>. doi:10.48550/ARXIV.2203.00545.
- [30] N. Reimers, I. Gurevych, Making monolingual sentence embeddings multilingual using knowledge distillation, in: *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP)*, Association for Computational Linguistics, Online, 2020, pp. 4512–4525. URL: <https://aclanthology.org/2020.emnlp-main.365>. doi:10.18653/v1/2020.emnlp-main.365.
- [31] X. Wang, T. Gao, Z. Zhu, Z. Zhang, Z. Liu, J. Li, J. Tang, Kepler: A unified model for

- knowledge embedding and pre-trained language representation (2019). URL: <https://arxiv.org/pdf/1911.06136.pdf>. arXiv:1911.06136.
- [32] J. Ács, Á. Kádár, A. Kornai, Subword pooling makes a difference, in: Proceedings of the 16th Conference of the European Chapter of the Association for Computational Linguistics: Main Volume, Association for Computational Linguistics, Online, 2021, pp. 2284–2295. URL: <https://aclanthology.org/2021.eacl-main.194>. doi:10.18653/v1/2021.eacl-main.194.
- [33] K. Labusch, C. Neudecker, Named entity disambiguation and linking on historic newspaper OCR with BERT, in: Working Notes of CLEF 2020 - Conference and Labs of the Evaluation Forum, number 2696 in CEUR Workshop Proceedings, CEUR-WS, 2020. URL: http://ceur-ws.org/Vol-2696/paper_163.pdf.
- [34] T. Padilla, Responsible Operations: Data Science, Machine Learning, and AI in Libraries, Technical Report, OCLC Research, USA, 2020. doi:10.25333/xk7z-9g97.
- [35] M. Gregory, C. Neudecker, A. Isaac, G. Bergel, Others, AI in Relation to GLAMs Task Force - Report and Recommendations, Technical Report, Europeana Network Association, 2021. URL: <https://pro.europeana.eu/project/ai-in-relation-to-glams>.