



CLEF-HIPE-2020

Named Entity Recognition and Linking on Historical Newspapers

Maud Ehrmann (EPFL)
Matteo Romanello (EPFL)
Alex Flückiger (UZH)
Simon Clematide (UZH)

CLEF 2020 (online), 23.09.2020

EPFL

 **Media Monitoring**
impresso of the Past

Universität Zürich^{UZH}

CLEF-HIPE-2020 in a nutshell

- **HIPE:** Identifying Historical People, Places and other Entities
- 1st NE processing shared task on historical documents
- Tasks:
 - NE recognition and classification
 - NE linking
- Participating teams: 13



Why HIPE?

New data:

Emergence of large-scale
archives of digitized contents

New needs:

Content retrieval by
humanities scholars

Challenge: NLP on historical texts is hard

- Spelling variations
- Noisy OCR
- Multilingualism
- Data sparsity
- Limited resources or KB coverage

→ Objectives

1. strengthen the **robustness** of approaches;
2. enable **performance comparison**;
3. foster **efficient semantic indexing** of digitized cultural heritage collections.

Background

impresso project
mining 200 years of historical newspapers



Digital Humanities Laboratory (DHLAB)
Ecole Polytechnique Federale de Lausanne,
Switzerland



University of
Zurich ^{UZH}

Institute of Computational Linguistics Zurich
University, Switzerland



Centre for Contemporary and Digital History
(C2DH) Luxembourg University,
Luxembourg.



The screenshot shows the impresso project's web interface. At the top, there is a navigation bar with links for "Search", "Newspapers", "Topics", "Inspect & Compare", "Text reuse", and "Collections". Below the navigation is a search bar with placeholder text "search for ...". To the right of the search bar are three buttons: "SEARCH ARTICLES", "SEARCH IMAGES", and "NGRAMS". Underneath the search bar is a section titled "IMPRESSO DATA RUNDOWN" which provides the following statistics:

- 76 newspapers collected,
- 600,930 issues,
- 5,422,655 pages scanned,
- 47,758,465 content items identified,
- 3,492,799 images,
- 12,493,358,703 words.

Below this, it says "2 countries of publication" and "530,086 named entities disambiguated". A link to the blog is provided: "More? Check on our [blog](#)". Further down, there are links for "info @ impresso-project [dot] ch", "project website: impresso-project.ch", "github: [impresso-project](https://github.com/impresso-project/impresso-project)", and "twitter: [gain full access](https://twitter.com/@impresso-project">@impresso-project". There are also buttons for "LINES: OFF" and "DARK MODE: ON". On the right side of the interface, there is a sidebar with the text "Mining 200 years of historical newspapers" and "How can newspapers help understand the past? How to explore them?". At the bottom, there is a note about legal access: "For legal reasons not all content is available in Open Access. To <a href=)" and a button for "DOWNLOAD NON-DISCLOSURE AGREEMENT FORM".



Swiss National Library, SNL



National Library of Luxembourg,
BNL



State Archives of Valais, AEV.



Swiss Economic Archives, SWA.



Le Temps



Neue Zürcher Zeitung, NZZ.



History department, University
of Lausanne, UNIL.



infoclio



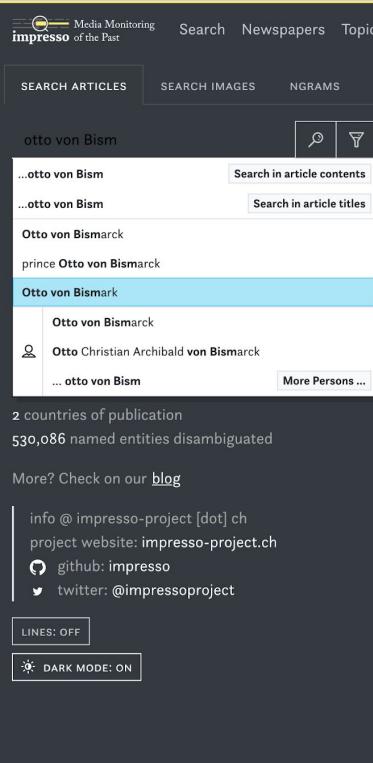
SWISS NATIONAL SCIENCE FOUNDATION

Project: <https://impresso-project.ch>

Interface: <https://impresso-project.ch/app/>

Semantic indexation of historical newspapers

Search NEs
(among others)
over 47M articles



The screenshot shows the impresso - Media Monitoring of the Past website. The search bar at the top contains "otto von Bism". Below the search bar, there are three tabs: "SEARCH ARTICLES", "SEARCH IMAGES", and "NGRAMS". The "SEARCH ARTICLES" tab is selected. The search results for "otto von Bism" are displayed in a list. The first two results are "...otto von Bism" (with options to "Search in article contents" or "Search in article titles"). The third result is "Otto von Bismarck". The fourth result is "prince Otto von Bismarck". The fifth result is "Otto von Bismark" (highlighted with a blue background). Below these results, there is a small icon of a person and the text "Otto Christian Archibald von Bismarck". At the bottom of the list, there is a link "... Otto von Bism" and a button "More Persons ...". To the left of the list, it says "2 countries of publication" and "530,086 named entities disambiguated". Below the list, there is a link "More? Check on our [blog](#)". At the bottom of the page, there is contact information: "info @ impresso-project [dot] ch" and "project website: [impresso-project.ch](#)". There are also links to GitHub ("github: [impresso](#)") and Twitter ("twitter: [@impressoproject](#)"). There are buttons for "LINES: OFF" and "DARK MODE: ON". On the right side of the page, the text "Mining 200 years of historical newspapers" is followed by a small info icon. Below that, the text "How can newspapers help understand the past? How to explore them?" is preceded by a small sunburst icon. A yellow sidebar on the right contains the text "For legal reasons not all content is available in Open Access. To gain full access: [DOWNLOAD NON-DISCLOSURE-AGREEMENT FORM](#) ... and return the signed form to info@impresso-project.ch". At the very bottom right, there is a small circular profile picture of a person and the text "Maud Ehrmann STAFF".

Media Monitoring of the Past

Mining 200 years of historical newspapers ⓘ

How can newspapers help understand the past? How to explore them?

For legal reasons not all content is available in Open Access.
To gain full access:

[DOWNLOAD NON-DISCLOSURE-AGREEMENT FORM](#)

... and return the signed form to info@impresso-project.ch

5

CLEF-HIPE-2020 Overview - M. Ehrmann, M. Romanello, A. Flückiger, S. Clematide

Semantic indexation of historical newspapers

Visualize facsimile, OCR and entity mentions

Media Monitoring
impresso of the Past

Search ... Newspapers Topics Inspect & Compare Text reuse Collections

FAQ ⓘ Maud Ehrmann STAFF

TABLE OF CONTENTS

43 articles in 4 pages (Personal use (no export))

Add filters from your current search query * 3 search filters can't be applied.

Search words...
 show only matching articles (no results)

La Journée politique

M. de Brazza est arrivé lundi à bord du Stamboul, courrier de la côte occidentale d'Afrique. M. de Brazza est parti hier pour Alg [article](#) [short text](#) p.1

Continue reading: [TRANSCRIPT](#)

LOCATIONS France ⓘ, Congo Free State ⓘ, Westminster ⓘ, Rome ⓘ, Naples ⓘ, Italy ⓘ, Bulgaria ⓘ, Russia ⓘ

PEOPLE Pierre Savorgnan de Brazza ⓘ, John Redmond ⓘ, Ferdinand I of Bulgaria ⓘ

A la Chambre Séance du 12 février Présidence de M. de Wacquani, président[...]

A la Chambre Séance du 12 février Présidence de M. de Wacquani, président. Décidément cei tains de nos députés trouvent qu'ils ont le temps pour s'occ [article](#) [short text](#) pp.2,3 (2 pages)

Continue reading: [TRANSCRIPT](#)

LOCATIONS Laval, Mayenne ⓘ

PEOPLE Victor de Tornaco ⓘ

Sans titre

Postas. — Un député de notre ville a tonné il y a quelques jours à la Chambre contre la distribution des lettres, qui depuis plusieurs années est fait [article](#) [short text](#) p.3

L'indépendance luxembourgeoise · Wednesday, February 13, 1895 [14 p. of 4] ⓘ

Faux départ

LOCATIONS Marcel Dallamagne ⓘ

PEOPLE Otto von Bismarck ⓘ, Guillaume Cale ⓘ

TOPICS

42.5% FR gouvernement - parti - ministre - président - politique ⓘ 18.6% FR roi - prince - empereur - comte - reine ⓘ 12.2% FR question - point - pays - conseil - gouvernement ⓘ

9.2% FR vie - esprit - pays - travail - effort ⓘ

[ADD TO COLLECTION ...](#)

CURRENT SELECTION X

Otto von Bismarck PERSON

12 apply current search filters (3 filters)

317 results in total (from 1857 to 2010)

[ADD AS SEARCH FILTER](#) [EXCLUDE FROM CURRENT SEARCH](#)

 Otto von Bismarck (Q8442)
HUMAN Schönhausen, 1815 - Friedrichsruh, 1898

[W] German statesman and Chancellor (1815-1898)
SOURCE: WIKIDATA [W] Q8442

MORE...

NO TEXT REUSE
PASSEGES AVAILABLE ⓘ

Faux

L'opinion publique assez vivement préoccupée de politique qui imposent en dépit des déments

en Allemagne est réoccupée du changement politique et qu'entraîne, dans les entretiens officieux, le chan-

gement de personnalité de l'empereur. Elle note les symptômes qui lui permettent de porter un jugement sur les tendances actuellement prévalentes en haut lieu. Les débats du Reichstag, ceux surtout de la commission chargée d'examiner le projet contre les menées subversives, ont déjà jeté un jour assez vif sur ce point. Bien que le prince de Hohenlohe, fidèle aux traditions d'une vie presque tout entière passée dans la diplomatie

ⓘ

Semantic indexation of historical newspapers

Otto von Bismarck PERSON

OVERVIEW 317 RELATED ARTICLES MENTIONED 334 TIMES

OPEN IN SEARCH PAGE... "German statesman and Chancellor (1815-1898)" (wikidata)



W
Otto von Bismarck
German statesman and Chancellor (1815-1898)
Schönhausen, 1815 - Friedrichsruh, 1898

SOURCE: W/Q8442



COUNTRY	LANGUAGE	TYPE	Count
Switzerland	German	article	244
Luxembourg	French	page article	66
		unclassified content	5
		advertisement	2

PERSON	LOCATION	TOPIC	Count
Otto von Bismarck	Berlin	DE regierung - paris - frankreich - minister - kammer	113
Wilhelm II, German Emperor	Germany	FR roi - prince - empereur - comte - reine	87
Herbert von Bismarck	Auch	DE könig - kaiser - königin - prinz - prinzessin	68
Karl Marx	Paris	DE verlag - buch - band - geschichte - werk	62
Adolf Hitler	Switzerland	DE krieg - deutschland - frankreich - welt - friede	59
Richard Wagner	France	DE welt - leben - mensch - wort - art	43
Debbie Reynolds	England	DE mann - hand - kopf - nacht - gesicht	42
François Duvalier	Lage	FR gouvernement - parti - ministre - président - politique	42
Frederick III, German Emperor	Marcel Dallemagne	FR guerre - paix - pays - peuple - politique	40
Frederick II of Prussia	Italy	FR vie - monde - mort - foi - peuple	39

PARTNER	ACCESSRIGHT	COLLECTION	Count
Swiss National Library	Personal use	tide	199
NZT	Personal use (no export)		80

Overview of named entities

Semantic indexation of historical newspapers

The screenshot displays the impresso platform's search results for the query "Otto von Bismarck". The interface includes a navigation bar with links to "Search", "Newspapers", "Topics", "Inspect & Compare", "Text reuse", and "Collections". A user profile for "Maud Ehrmann STAFF" is visible on the right.

The main search results page shows a sidebar titled "BROWSE 530,086 ENTITIES" with a "filter entities" dropdown and a "Most mentioned in different articles" button. Below this is a list of locations and their article counts:

- Lausanne **location**: 2,918,317 ARTICLES, 4,642,691 MENTIONS
- Suisse, Moselle **location**: 2,651,532 ARTICLES, 4,268,837 MENTIONS
- Switzerland **location**: 2,390,778 ARTICLES, 4,727,170 MENTIONS
- Fribourg **location**: 2,285,590 ARTICLES, 4,264,647 MENTIONS
- Paris **location**: 2,132,568 ARTICLES, 3,551,470 MENTIONS
- France **location**: 2,082,458 ARTICLES, 4,036,841 MENTIONS
- Gare de Cornavin **location**: 1,858,798 ARTICLES, 2,942,990 MENTIONS
- Lake Neuchâtel **location**: 1,791,813 ARTICLES, 2,680,593 MENTIONS
- La Chaux-de-Fonds **location**: 1,411,002 ARTICLES, 2,104,854 MENTIONS
- Zürich **location**: 1,112,003 ARTICLES, 2,257,052 MENTIONS
- Italy **location**: 1,037,830 ARTICLES, 1,067,989 MENTIONS

The main content area for "Otto von Bismarck" shows the following details:

- OVERVIEW**: 317 RELATED ARTICLES
- MENTIONED 334 TIMES**
- OPEN IN SEARCH PAGE...**
- ORDER BY PUBLICATION DATE, OLDEST FIRST ▾**

Sample snippets of news articles containing mentions of Otto von Bismarck include:

- Blättert liest man folgendes „ Gedicht an **Otto von Bismarck** , ehemaligen Reichskanzler : Der alte Otto ! Vr
- Vermischtes Bündner Nachrichten **Saturday, March 7, 1891 – p.3**
In deutschen Blättern liest man folgendes „ Gedicht an Otto von Bismarck , ehemaligen Reichskanzler : Der alte Otto ! ...
- Hugo (p . n . k .) 1815 : Geburt **Otto von Bismarcks** in Schönhausen . Nachmaliger deutscher
- EIDGENOSENSCHAFT Die Tat **Tuesday, April 1, 1941 – p.6**
Graubünden Die Gründungsversammlung der lungliberalen Bewegung Graubündens fand in Thusis statt . An der öffentlichen Kundgebung referierte der Präsident der...
- von Bismarck , der jüngere Bruder des Grafen **Otto von Bismarck** und Enkel des Eisernen Kanzlers , wegen
- Einmarsch der Russen in Jugoslawien Die Tat **Monday, October 2, 1944 – p.2**
Calais erstürmt Moskau , 1. Okt . (Exchange .) Das weitaus bedeutendste Ereignis in den Kämpfen der...
- Im Schloß Friedrichsruh bei Hamburg , wo Fürst **Otto von Bismarck** begraben liegt , hat Himmler dem schwedischen
- Im letzten Stadium Die Tat **Wednesday, May 2, 1945 – p.1**
Im Schloß Friedrichsruh bei Hamburg , wo Fürst Otto von Bismarck begraben liegt , hat Himmler dem schwedischen Grafen Bernadotte...
- , daß nicht Karl Marx das Rennen macht , sondern **Otto von Bismarck** , sei es in Gestalt einer neuen Partei oder
- thy Paris , q . Nov . Die nationale Konfe... Die Tat **Monday, November 10, 1947 – p.2**
-thy . Paris , q . Nov . Die nationale Konferenz der Gewerkschaftsfaktion « Force ouvrière » , die am...
- (UP) Der Enkel des « Eisernen Kanzlers » **Otto von Bismarck** , Graf Gottfried von Bismarck , erlitt am
- ab . Der Verunfallte ist auf dem Transpor... Die Tat **Saturday, September 17, 1949 – p.5**
ab . Der Verunfallte ist auf dem Transport ins Kantonsspital Luzern gestorben . • Nyon . Der 73 jährige Landwirt...

Pagination controls at the bottom of the main content area show pages 1 through 8.

Tasks

1. NERC

Recognition and classification of entity mentions with

- subtask 1: coarse types
- subtask 2: fine-grained types.

	is_NIL loc.adm.nat Morocco	
1	3 L. e	A ££ uroo .
	is_NIL pers.ind Stephen Pichon	
	-is_NIL comp.name	
2	M. PIclaon	à ffl Madrid .
	-is_NIL comp.title	
	-is_NIL loc.adm.town Madrid	
	-is_NIL comp.function	
3	M. Pichon,	ministre des affaires
	-is_NIL comp.name	
	-is_NIL comp.title	
	-is_NIL pers.ind Stephen Pichon	
4	étrangères , esl arrive' à	Madrid . Aujourd'hui,
5	il ^ dîne chez le roi. Diverses fêles sont	
6	projetées en son honneur.	

Types	Sub-types	
	Coarse	Fine
pers	pers.ind pers.coll	pers.ind.articleauthor
org	org.ent org.adm	org.ent.pressagency
prod	prod.media prod.doctr	
date	time.date.abs	
loc	loc.adm	loc.adm.town loc.adm.reg loc.adm.nat loc.adm.sup
		loc.phys
		loc.geo loc.hydro loc.astro
		loc.oro loc.fac
		loc.add
		loc.add.phys loc.add.elec

Table 1: Entity types used for NERC tasks.

Tasks

2. Entity Linking

Towards Wikidata QID or NIL



- end-to-end EL: w/o mention boundaries
- EL-only: with mention boundaries



The screenshot shows a NER interface processing a French text. It highlights entities such as "Stephen Pichon" (with a link to Q1069292), "M. Pichon, ministre des affairesétrangères" (with a link to Q1069292), and "Madrid". The interface also shows various entity types like "loc.adm.nat", "comp.name", and "comp.function". Below the main text, a summary sentence is provided: "étrangères , esl arrive' à Madrid . Aujourd'hui, il ^ dîne chez le roi. Diverses fêles sont".

Participation bundles:

Bundle id	Associated tasks
bundle1	NERC-coarse and NERC-fine and NEL
bundle2	NERC-coarse and NEL
bundle3	NERC-coarse and NERC-fine
bundle4	NERC-coarse
bundle5	NEL

Participation guidelines: [10.5281/zenodo.3677171](https://doi.org/10.5281/zenodo.3677171)

Corpus selection

- Digitized newspaper archives (CH, LU, US)
- Diachronic: from 1738 to 2019
- Multilingual: fr, de, en
- Sampling and manual triage:
 - journalistic content
 - no feuilleton, cross-words, meteo, etc.
 - exclusion of extreme OCR noise
 - no provision of different OCR → real-life setting

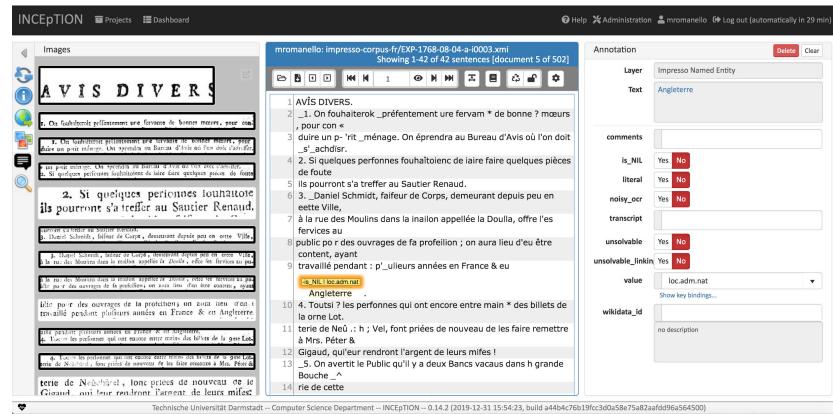


3 L. e A EE uroo.
 M. Pichon à Madrid.
M. Pichon, ministre des affaires étrangères, est arrivé à Madrid. Aujourd'hui, il ^ dîne chez le roi. Diverses fêtes sont projetées en son honneur.
 . Le correspondant du *Temps* en Espagne télégraphie qu'à Madrid on attache une grande importance au voyage de M. Pichon. Les cercles politiques, financiers et militaires croient qu'on approche du moment décisif dans les affaires du Maroc. On estime qu'en présence de l'agitation marocaine contre l'organisation de la police, une action énergique peut prévenir bien des violences, et l'on s'attend à ce que, des Conférences que va tenir M. Pichon avec le roi et ses ministres, il résulte une action combinée plus active de la France et de l'Espagne.

Corpus annotation

- Trilingual annotators, trained on a mini-ref
- INCEpTION platform

- NERC annotation difficulties:
 - NE mention boundaries
 - consideration of multiple languages
 - what is to be annotated or not
 - definition at time x
 - metonymy



The screenshot shows the INCEpTION platform's annotation interface. On the left, a document page displays a section titled "AVIS DIVERS" with several numbered items in French. On the right, a sidebar titled "Annotation" allows users to select a "Layer" (e.g., "Impresso Named Entity", "Text", "Anglais") and "Text" (e.g., "Anglais"). A "comments" field is present, along with checkboxes for "is_NIL", "literal", "noisy_ocr", "transcript", "unsolvable", "unsolvable_linkin", and "wikidata_id". A dropdown menu for "value" shows "loc.adm.net" selected. A note at the bottom indicates "Show key bindings...".

M. Curtoys d' Anduaga, doyen du corps diplojtelfsue espagnol, et ministre plenipotentiaire pendant 50 ans

Zürichputsch, Baslerpropaganda

Commission imperiale, Die französische Regierung

Is Savoie or Moldavia a region or a country?

Corpus annotation

- Trilingual annotators, trained on a mini-ref
- INCEpTION platform
- NERC annotation difficulties:
 - NE mention boundaries
 - consideration of multiple languages
 - what is to be annotated or not
 - definition at time x
 - metonymy
- EL annotation difficulties:
 - Requires historical knowledge + Sherlock Holmes skills
 - Historical statuses of entities unequally represented in KB

Germany, Q183

- 962-1813: Holy Roman Empire, Q12548
- 1806-1813: Confederation of the Rhine, Q154741
- 1815-1866: German Confederation, Q151624
- 1867-1870: North German Confederation, Q150981
- 1871-1918: German Empire, Q43287
- 1918-1933: Weimar Republic, Q41304
- 1933-1945: Nazi Germany, Q7318
- 1949-1990: West Germany, Q713750
- 1949-1990: East Germany, Q16957

Corpus characteristics

newspaper articles	563
tokens	444,596
(linked) mentions	18,962
metonymy	1252
components	6,219
noisy mentions (test set)	10%
NIL	25.72%

mentions: 10,923 (Fr), 6584 (De), 1455 (En)

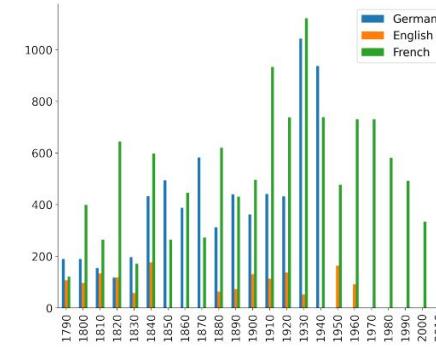


Fig. 2: Diachronic distribution of mentions across languages.

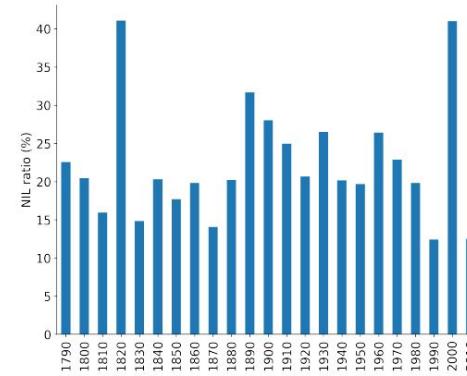


Fig. 3: Diachronic ratio of NIL entities.

Corpus release

- train/dev/test (70/15/15)
- no train set for English
- no sentence segmentation
- no sophisticated tokenization
- document metadata

TOKEN	NE-COARSE-LIT	NE-COARSE-METO	NE-FINE-LIT	NE-FINE-METO	NE-FINE-COMP	NE-NESTED	NEL-LIT	NEL-METO	MISC
# language	= fr								
# newspaper	= GDL								
# date	= 1908-01-07								
# document_id	= GDL-1908-01-07-a-i0009								
[...]									
# segment_iif_link	= https://iiif.dhlab.epfl.ch/iiif_impresso/GDL-1908-01-07-a-p0002/296,3651,514,64/full/0/default.jpg								
M	B-pers	0	B-pers.ind	0	B-comp.title	0	Q1069292	_	NoSpaceAfter
.	I-pers	0	I-pers.ind	0	I-comp.title	0	Q1069292	_	_
Piela	clon	I-pers	0	I-pers.ind	0	B-comp.name	0	Q1069292	_
à	0	0	0	0	0	_	_	_	
ffltadrid	B-loc	0	B-loc.adm.town	0	0	0	Q2807	_	NoSpaceAfter
.	0	0	0	0	0	_	_	_	EndOfLine
# segment_iif_link	= https://iiif.dhlab.epfl.ch/iiif_impresso/GDL-1908-01-07-a-p0002/213,3709,731,64/full/0/default.jpg								
M	B-pers	0	B-pers.ind	0	B-comp.title	0	Q1069292	_	NoSpaceAfter
.	I-pers	0	I-pers.ind	0	I-comp.title	0	Q1069292	_	_
Pichon	I-pers	0	I-pers.ind	0	B-comp.name	0	Q1069292	_	NoSpaceAfter
,	I-pers	0	I-pers.ind	0	0	0	Q1069292	_	_
ministre	I-pers	0	I-pers.ind	0	B-comp.function	0	Q1069292	_	_
des	I-pers	0	I-pers.ind	0	I-comp.function	0	Q1069292	_	_
affaires	I-pers	0	I-pers.ind	0	I-comp.function	0	Q1069292	_	EndOfLine
# segment_iif_link	= https://iiif.dhlab.epfl.ch/iiif_impresso/GDL-1908-01-07-a-p0002/171,3753,770,64/full/0/default.jpg								
étrangères	I-pers	0	I-pers.ind	0	I-comp.function	0	Q1069292	_	NoSpaceAfter
,	0	0	0	0	0	_	_	_	
esl	0	0	0	0	0	_	_	_	
arrive	0	0	0	0	0	_	_	_	NoSpaceAfter
'	0	0	0	0	0	_	_	_	
à	0	0	0	0	0	_	_	_	
Madrid	B-loc	0	B-loc.adm.town	0	0	0	Q2807	_	NoSpaceAfter
.	0	0	0	0	0	_	_	_	
Aujourd	0	0	0	0	0	_	_	_	NoSpaceAfter
'	0	0	0	0	0	_	_	_	NoSpaceAfter
hui	0	0	0	0	0	_	_	_	NoSpaceAfter
,	0	0	0	0	0	_	_	_	EndOfLine

CC BY-NC 4.0

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

[10.5281/zenodo.3706857](https://zenodo.10.5281/zenodo.3706857)

Auxiliary resources

In-domain Fr, De, and En embeddings:

- fastText word embeddings (with and w/o subwords)
- flair character embeddings (now integrated into the flair [framework](#))

CC BY-SA 4.0

[https://files_ifi_uzh_ch_cl_siclemat_impresso_clef-hipe-2020/
10.5281/zenodo.3706808](https://files_ifi_uzh_ch_cl_siclemat_impresso_clef-hipe-2020_10.5281_zenodo.3706808)

Evaluation

- Entities (not tokens) as the unit of reference
- Macro & **Micro** Precision, Recall and F1 measure
- Evaluation scenarios:

	NERC	EL
Strict	exact mention boundaries	consideration of the top link only, (overlapping mention boundaries)
Fuzzy	overlapping boundaries	historical mapping, cut-offs @3 and @5 (overlapping mention boundaries)

HIPE Scorer: <https://github.com/impresso/CLEF-HIPE-2020-scorer>

HIPE Eval Toolkit: <https://github.com/impresso/CLEF-HIPE-2020-eval>

Participation

40
registrations

75 runs
42% French
31% German
26% English
6 teams work on all languages

13
participating teams
All participated to NERC-Coarse
3 to NERC-Fine
5 to EL-only and end-to-end EL

11
Working Notes

Participating systems' main features

- 11 teams applied **neural approaches** for NERC;
- Most of them worked with **contextualized embeddings**, esp. **BERT**;
- Experimentation with **various input embeddings** (char, subword, word, historical or contemporary, type-level or contextualized)
- Some attempted to improve the **newspaper line-based input format** with proper sentence segmentation and tokenization;

Results overview (NERC)

- Neural system with strong embedding resource prevail;
- Performances correlates with amount of train/dev data;
- BERT-based systems > Bi-LSTM;
- Great performances diversity, but results are better than expected (6 teams > .8);
- NERC fine with 12 classes more difficult;
- NE components show reasonable performances.

F1 scores	French		German		English	
	Strict	Fuzzy	Strict	Fuzzy	Strict	Fuzzy
<i>NERC-Coarse literal</i>						
Baseline	.646	.769	.476	.585	.405	.562
Median	.677	.808	.636	.766	.463	.645
Best system	.840	.921	.797	.878	.632	.806
<i>NERC-Coarse metonymic</i>						
Best system	.783	.783	.634	.694	-	-
<i>NERC-Fine</i>						
Best system	.784	.856	.668	.771	-	-
<i>NE components</i>						
Best system	.657	.751	.642	.707	-	-

Results overview (EL)

- EL performances are lower, and as diverse;
- NERC error propagation in end-to-end setting, but EL-only not a lot better;
- Performance increase with cut-offs @3 and @5.

Overall, what helps:

- BERT;
- actively tackling the problems of OCR noise, word hyphenation and sentence segmentation;
- in-domain resources.

F1 scores	French		German		English	
	Strict	Fuzzy	Strict	Fuzzy	Strict	Fuzzy
<i>End-to-end Entity Linking (literal)</i>						
Baseline	.257	.270	.180	.195	.239	.239
Best system	.598	.617	.534	.557	.531	.531
<i>End-to-end Entity Linking (metonymic)</i>						
Best system	.297	.462	.396	.469	-	-
<i>Entity linking only (with mentions provided)</i>						
Baseline	.498	.512	.418	.437	.506	.506
Best system	.639	.659	.582	.602	.658	.658

Time-based observations

Analysis of F1 score as a function of time.

Hypothesis: the older, the more difficult.

Observation: no strong correlation between article publication date and performance.

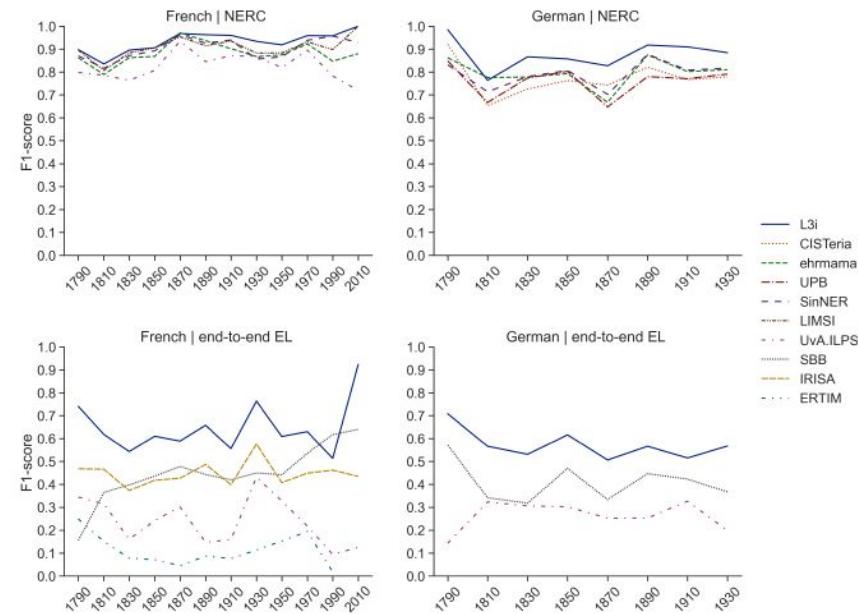
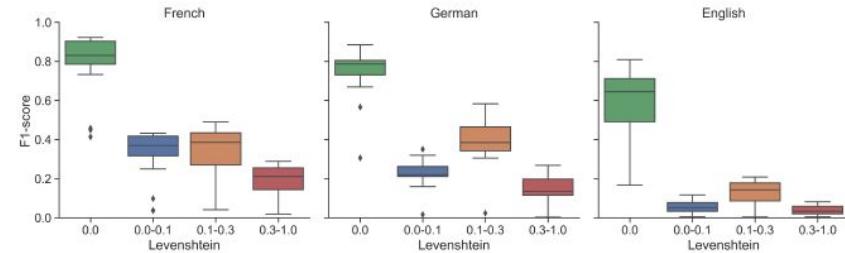


Fig. 4: F1 score as a function of time for the 5 best systems for NERC (top) and end-to-end EL (bottom) for the languages French (left) and German (right). The x-axis shows 20-years time buckets (e.g. 1790 = 1790-1809).

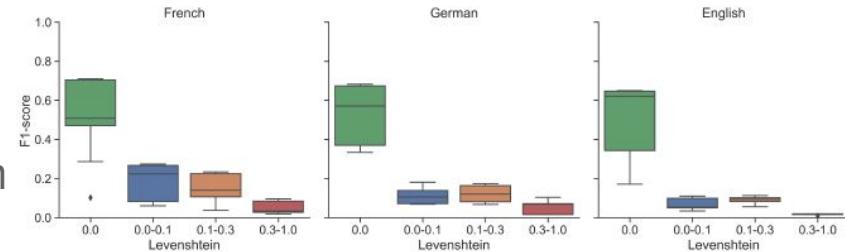
Impact of OCR noise

Evaluation on various noise levels

- noise: length-normalized Levenshtein distance between surface form and manual transcription;
- noisy vs non-noisy have remarkable differences on both NERC and EL;
- greatest performance variation at medium noise level



(a) NERC-Coarse.



(b) End-to-end EL with the relaxed evaluation regime and a cutoff @3.

Fig. 5: Impact of OCR noise: distribution of performances across systems on entities with different noise level severity for NERC (a) and end-to-end EL (b).

Conclusion (1/2)

- **Robustness test** for NERC and EL approaches on challenging historical material;
- New insights in **domain and language adaptation**;
- **Neural-based systems** with strong resources and proper segmentation are capable of dealing with historical and noisy inputs;
- EL, nested entities, entity components remain challenging;
- Performances are affected by OCR noise, but not by document publication date;

Discover more about systems this afternoon 3-6:30pm!

Conclusion (2/2)

Main Outcomes:

- Contribution to the advance of SoTA for NE processing on historical texts;
- Datasets;
- Scorer;
- Annotation guidelines;
- Step towards efficient semantic indexing of historical material.

Future Directions:

- Potential HIPE2 in 2021 with additional document types and languages.

Many thanks to

- CLEF organizers
- Participating teams (kudos!)
- *NZZ*, *Le Temps*, and the Swiss and Luxembourg national libraries
- Richard Eckart de Castillo, Clemens Neudecker, Sophie Rosset and David Smith
- INCEpTION project team
- Camille Watter, Gerold Schneider, Emmanuel Decker and Ilaria Comes
- SNSF (grant number CR-SII5 173719)

Thank you for your attention



HIPE : impresso.github.io/CLEF-HIPE-2020/



Scorer : impresso.github.io/CLEF-HIPE-2020/



Evaluation toolkit : impresso.github.io/CLEF-HIPE-2020/



Impresso project : impresso.github.io/CLEF-HIPE-2020/



@ImpressoProject

