

# Comparaisons et combinaisons d'extracteurs d'entités spatiales sur un corpus multilingue.

---

Zijian Wang<sup>1</sup>, Yoann Dupont<sup>2</sup>, Tian Tian<sup>2</sup>, Gaël Lejeune<sup>2</sup>

Sorbonne Université, STIH,

<sup>1</sup> wangzijian1994@hotmail.com

<sup>2</sup> prenom.nom@sorbonne-universite.fr

24 Juin 2020

Contexte : Humanités Numériques et TAL

EN en humanité numérique : données

Extracteurs d'entités nommées existants

Résultats et discussion

# Contexte : Humanités Numériques et TAL

---

- ANalyse auTOmatique et NumérisatiOn des MAZarinades<sup>1</sup>  
DIM STCN

---

1. <http://lejeunegael.fr/antonomaz.html>

- ANalyse auTOmatique et NumérisatiOn des MAZarinades<sup>1</sup>  
DIM STCN
- Le Traitement Automatique des Langues
- Les Humanités Numériques
- → leur interaction

Deux axes de réflexion :

- *Input* : Documents numérisés, quelles contraintes? (thèse de JB Tanguy 2019-...)
- *Output* : Adaptation/combinaison de systèmes de détection d'Entités Nommées (thèse de Caroline Parfait 2020-...)

---

1. <http://lejeunegael.fr/antonomaz.html>

# Contraintes sur l'Input

Gérer la variation :

- Multilinguisme
- Hétérogénéité
- Bruitage

	<b>VISQVE</b> Babillard on me nomme, le ne veux espargner nul homme, le fuis sous & remply de vin le veux parler de Tabarin De Tabarin ce Mazinique, Cest homme peruers & inique, Qui n'a ny Dieu, ny Foy, ny Loy Qui a enleué nostre Roy, Et fait assieger nostre Ville: Comme vn Meschant & Malhabille Par ce grand Prince de Condé Qu'il a enchanté sans tardé Qui a fillé, chose certaine, Les yeux de nostre bonne Reyne,	VIS QVE Babillard on me nomme InlÉfP le ne veux efpargner nul homme, WjfflsL le fuis fous & remply de vin Icvcuxrparler deTabarin DeTabarin ce Mazinique, . Cesthomrneperuers & inique, Qiii n'any Dieu, nyFoy.nyLoy Quiaenleuénoftre Roy, . Et fait affieger noftre Ville : Commervn Meschant ôC Malhabille Par ce grand Prince de Condé Qu'il a enchanté fans tarde* J Qui a fille, chose certaine. Les yeux de noftre bonne Rey ne,
---	---	---

Figure 1 – Source : Abiven et Lejeune 2019

Question des "conditions de laboratoire"

- Applicabilité des Systèmes état de l'art
- Évaluation/Interprétation
- Valeur ajoutée pour l'utilisateur

Focus aujourd'hui : Entités Nommées de lieu en Allemand  
(CLEF-HIPE 2020 P.Ortiz, Y.Dupont, G.Lejeune, T.Tian<sup>2</sup>)

---

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

Enrichir les données textuelles avec une couche "sémantique" pour permettre de "raisonner" automatiquement.

- Noms d'organisation, noms de personne
- Noms de fêtes, noms de films/chansons,
- Noms de protéines . . .

Enrichir les données textuelles avec une couche "sémantique" pour permettre de "raisonner" automatiquement.

- Noms d'organisation, noms de personne
- Noms de fêtes, noms de films/chansons,
- Noms de protéines . . .
- Noms de lieux

## Les entités nommées : définitions

À (Boncourt)*lieu* (en 1886)*date* :

La boulangerie de (M. Staempfli)*personne*, (rue du parc, 68)*lieu*.

la (confédération du (Rhin)*lieu*)*lieu??*

Adapter les systèmes d'extraction d'entités nommées standard sur :

- des corpus anciens
- non nativement numériques
- multilingues

## EN en humanité numérique : données

---

## Textes anciens : problèmes

transcription à la main -> données bruités

segmentation des textes initiaux -> tiret séparant les mots tâches,

espace en trop -> caractères en trop, OOV (out-of-vocabulary)

Corpus d'apprentissage et d'évaluation de la tâche CLEF-HIPE  
 En français et allemand, annotés en entités de type *Personne*,  
*Organisation*, *Location*, *Product* et *Time*

	tokens	nb doc	nb seg	nb d'entités nommées annotées				
				Pers	Loc	Org	Time	Prod
<i>train Fr</i>	166217	158	19183	3067	2513	833	273	198
<i>dev Fr</i>	37592	43	4423	771	677	158	69	48
<i>train De</i>	86960	104	10353	1747	1170	358	118	112
<i>dev De</i>	36175	40	4186	664	428	172	73	53

# Extrait du corpus français

```
# language = fr
# newspaper = EXP
# date = 1918-04-22
# document_id = EXP-1918-04-22-a-i0077
# segment_iiif_link = https://iiif.dhlab.epfl.ch/.../1186,1881,474,79/full/0/default.jpg
Lettre      0      ]
de          0      -
la          0      -
Su          B-loc  NoSpaceAfter
.           I-loc  -
_           I-loc  NoSpaceAfter
sss        I-loc  -
allemands  I-loc  EndOfLine
# segment_iiif_link = https://iiif.dhlab.epfl.ch/.../1190,1967,493,52/full/0/default.jpg
(           0      NoSpaceAfter
Nous       0      -
serons    0      -
heureux   0      -
de        0      -
publier   0      -
de        0      -
temps     0      -
à         0      EndOfLine
# segment_iiif_link = https://iiif.dhlab.epfl.ch/.../1165,1995,517,53/full/0/default.jpg
autre     0      NoSpaceAfter
,         0      -
sous      0      -
cette     0      -
rubrique  0      NoSpaceAfter
,         0      -
des       0      -
```

## Extrait du corpus français - Exemple de bruit

Lettre	0	-
de	0	-
la	0	-
Su	B-loc	NoSpaceAfter
.	I-loc	-
_	I-loc	NoSpaceAfter
sss	I-loc	-
allemands	I-loc	EndOfLine

Forme non-standard :

"Su. \_sss allemands" -> "Suisse allemande"

## Extracteurs d'entités nommées existants

---

## Extracteurs multilingues d'entités nommées

- Modèles appris et configurés sur une langue "standard" (la langue journalistique).
- Adaptation avec apprentissage automatique à partir de données annotées dans la langue/le domaine cible

nom	langue initiale	algorithme	corpus d'apprentissage
SpaCy	anglais	CNN	Conll2003 (journal)
Standford NLP	anglais	CRFs	Conll2003 (journal)
SEM	français	CRFs	French Treebank

Présentation :

- réseau de neurones avec CNN.
- multilingue, mais qualité variable.

Présentation :

- réseau de neurones avec CNN.
- multilingue, mais qualité variable.

Adaptation :

- entraînement depuis zéro (blank).
- adaptation depuis modèle existant (de\_core\_web\_sm).

Présentation :

- SEM : outil utilisant CRF, Wapiti.
- Enrichissements textuels, gestion de format.
- Uniquement pour le français.

## Présentation :

- SEM : outil utilisant CRF, Wapiti.
- Enrichissements textuels, gestion de format.
- Uniquement pour le français.

## Adaptation :

- Enrichissement des lexiques existants.
- Entraînement sur le corpus.
- Beaucoup de travail manuel.

Présentation :

- architecture état-de-l'art pour l'étiquetage de séquences.
- LSTM bi-directionnel (contextes gauche et droit) + CRF (cohérence des étiquettes).
- "plug and play" : utiliser des embeddings pré-entraînés et adapter reste du réseau sur données.

## Présentation :

- architecture état-de-l'art pour l'étiquetage de séquences.
- LSTM bi-directionnel (contextes gauche et droit) + CRF (cohérence des étiquettes).
- "plug and play" : utiliser des embeddings pré-entraînés et adapter reste du réseau sur données.

## Adaptation :

- Concaténation de 3 représentations : FastText original, Flair + FrELMo (ELMo entraîné sur OSCAR-fr).
- Entraînement sur le corpus (représentations figées).

## Résultats et discussion

---

ystème	P	R	F
Spacy (blank)	64,8%	53%	58,3
Spacy (pré-entraîné)	70,9%	59,1%	64,42
SEM	64,4%	43,8%	52,1
Bi-LSTM-CRF	<b>63.1%</b>	<b>66,6%</b>	<b>64,8</b>

## Quelques observations

- Réseaux de neurones meilleurs que CRF (métriques)
- RN plus facilement adaptables que CRF (pré-entraînement vs ré-entraînement)
- adapter un modèle pré-entraîné améliore les résultats (divers papiers)

Plafond de verre ?