

Rencontres Minutes de l'équipe de Linguistique Computationnelle aka Speed Dating

10 septembre 2020

STIH, Sorbonne University

Le tiercé dans l'ordre

1. Karën Fort
2. Carlos González
3. Lichao Zhu
4. Vincent Lully
5. Imed Laaridh
6. Yoann Dupont
7. Zijian Wang
8. Julien Bezancon
9. Caroline Parfait
10. Gaël Lejeune
11. Fouad Aounti
12. Jean Baptiste Tanguy

Un séminaire pour quoi ?

- Des rencontres régulières (1 fois/mois)
- Pas trop formelles

Un séminaire pour quoi ?

- Des rencontres régulières (1 fois/mois)
- Pas trop formelles
- Comprendre l'éco-système
- Connaître les collègues
- Partager des idées
- Collaborer (en enseignement et en recherche)

Les rencontres minutes d'aujourd'hui : faire un tour d'horizon rapide et boire des cafés

Le Wiki <http://stih-sorbonne-universite.fr/dokuwiki/doku.php?id=seminaire>

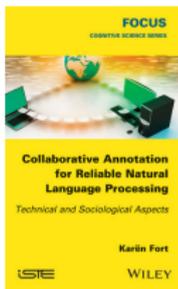
La "FAQ" <https://gitlab.com/YoannDupont/accueil-su/-/blob/master/guide-ac>

Karën Fort

D'où je parle

Voir <https://www.schplaf.org/kf/>

- ▶ Création de ressources langagières pour le TAL



- ▶ Ethique et TAL



Production participative (*crowdsourcing*)



ZOMBILINGO

RIGORMORTIS

BISAME

KRIK!

AYO!



Portail de jeux pour les langues et atelier récurrent :

L I N G O B O I N G O

Games4NLP

Projets locaux :

- ▶ MEMES : A. Gautier, G. Lejeune et G. Sioufi
- ▶ CAMARADERIE : M. Avanzi et A. Thibault

Conseils

- ▶ Conseil national des universités (CNU) 27 (informatique) :
<https://cnu27.univ-lille.fr/>
 - ▶ Qualifications
 - ▶ Promotions, CRCT, etc
- ▶ Comité d'éthique de la recherche (CER) de SU
- ▶ GDR LIFT (linguistique informatique, formelle et de terrain)
- ▶ actions européennes COST :
 - ▶ European Network for Combining Language Learning with Crowdsourcing Techniques (enetCollect)
 - ▶ Language In The Human-Machine Era (LITHME)

Carlos González



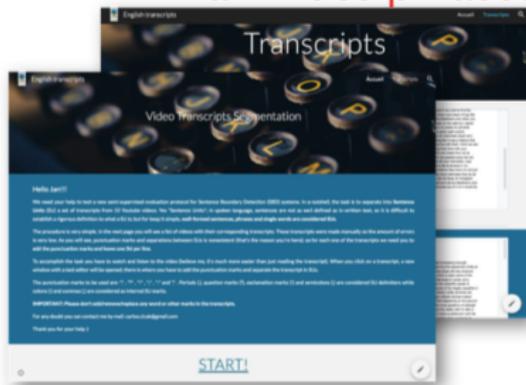
Séminaire Linguistique Computationnelle

Carlos González

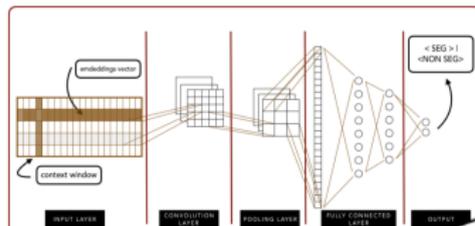
10/09/2020

Axes de recherche

Y a-t-il des phrases dans le langage oral ?



	Mots	Segmentation manuelle			Moyenne
		ref ₁	ref ₂	ref ₃	
Total	8 080	502	485	281	422±100



Information textuelle des transcriptions

- **Sac-de-mots, n -grammes de mots/caractères, plongements de mots**
- ✓ Information textuelle, ressources linguistiques disponibles
- ✗ erreurs de transcription

Signal audio et informations extraites

- **Pauses, durée des mots, niveau du signal, énergie**
- ✓ Indépendante de l'étape de transcription
- ✗ découpage des phrases incomplètes

WiSEBE: Window-Based Sentence
Boundary Evaluation

Carlos-Emiliano González-Gallardo^{1,2(✉)} and Juan-Manuel Torres-Moreo
LIA - Université d'Avignon et des Pays de Vaucluse, 339 chemin des Meiries,
84140 Avignon, France
carlos-emiliano.gonzalez-gallardo@univ-avignon.fr
juan-manuel.torres@univ-avignon.fr
LIA - GIGL, École Polytechnique de Montréal,
Montréal, Québec H3C 3A7, Canada

Axes de recherche

Résumer pour informer

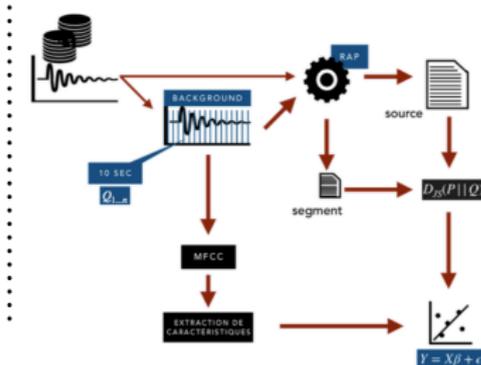
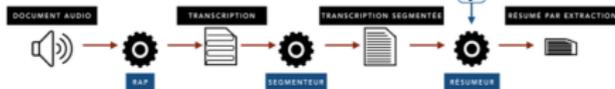
ACTUALITÉS / TÉLÉREPORTAGES



multimédia & multilingue

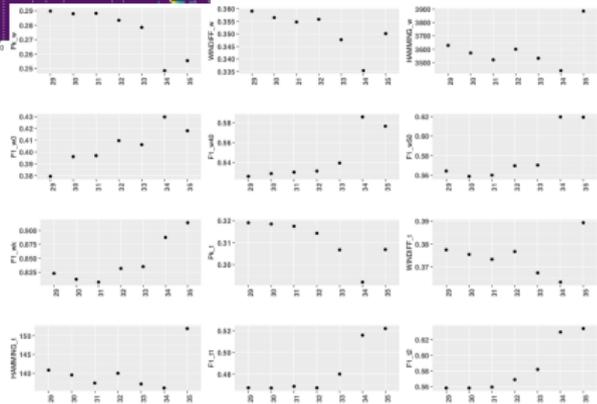
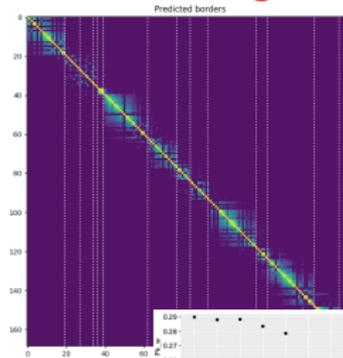
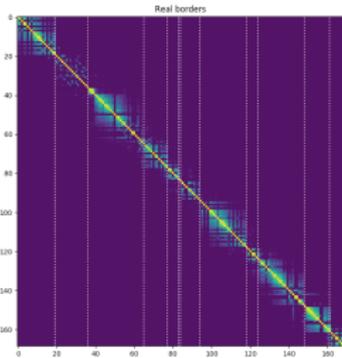


→ résumé



Axes de recherche

Diviser et régner



 Global News Podcast

carlos-emiliano.gonzalez-gallardo@paris-sorbonne.fr

Lichao Zhu

MÈMES pas mèmes

Lichao ZHU

10 septembre 2020



Qu'est-ce ?

- ▶ Mèmes : Extraction automatique et analyse par Myriadisation d'Expressions Semi-figées (Karën, Antoine, Gaël et bibi)
- ▶ Annotations et modélisations des candidats mèmes

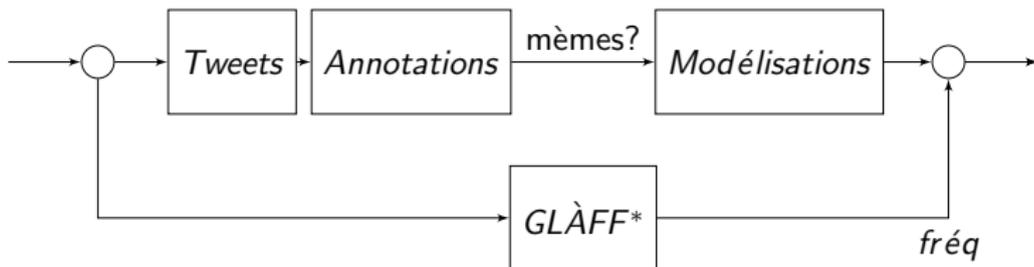
Exemples

- Roubaix, ton univers (politique) impitoyable ... !
- @PiotrSzut URSSAF, ton univers impitoyable

Questions

- ▶ Quels sont des traits formels des mèmes (et leur origine) ?
- ▶ Sur quels critères formels peut-on les repérer ?

Comment ?

**Méthode**

- ▶ annotations des candidats mèmes ;
- ▶ typologies phono-morpho-syntaxiques des cas avérés et non avérés ;
- ▶ modélisations formelles des mèmes avérés ;
- ▶ comparaisons avec les données du GLÀFF.

* : Gros Lexique À tout Faire du Français

Observations

Dans un mème,

- ▶ la structure syntaxique est souvent banale ;
- ▶ le nombre de mots outils est souvent supérieur à 2 ;
- ▶ au moins un mot est de moyenne ou de forte fréquence lexicale ;
- ▶ la ponctuation joue un rôle important.

Vincent Lully

- Vincent LULLY
- Master ILGII -> Thèse CIFRE -> ATER -> MCF
- Web sémantique et ses applications
 - Système de recommandation !!!
 - Moteur de recherche !!
 - Patrimoine culturel !
- Comparaisons intrinsèques et extrinsèques des plongements sémantiques et lexicaux
- Auto-complétion dans des systèmes de questions-réponses

Imed Laaridh

Imed Laaridh, PhD – Traitement automatique de la parole

Formation universitaire



- | | |
|-----------|--|
| 2013-2017 | Doctorat en Informatique : Laboratoire LIA, Bourse Brain and Language Research Institut (BLRI/ILCB). CERI, Université d'Avignon, Avignon, France. |
| 2010-2013 | Ingénieur en Informatique
École Nationale des Sciences de l'Informatique (ENSI), Tunisie. |
| 2008-2010 | Études préparatoires en Mathématiques et Physique
Institut Préparatoire aux Études d'Ingénieurs de Tunis (IPEIT), Tunis, Tunisie |

Positions académiques



- | | |
|-----------|--|
| 2017-2020 | Post-doc : Laboratoire IRIT, Financement projet région & ANR.
IRIT, Université de Toulouse/CNRS |
| 2017-2017 | Ingénieur de recherche : Laboratoire Informatique d'Avignon (LIA)
LIA-CERI, Université d'Avignon, Avignon. |

Thèmes de recherche

- Traitement automatique de la parole pathologique
- Évaluation de la qualité de la prononciation au niveau phonème (projet TypALoc)
 - Alignement de la parole lue des patients au niveau phonème
 - Décision au niveau phonème de la normalité ou non de la production

Thèmes de recherche

- Traitement automatique de la parole pathologique
- Évaluation de la qualité de la prononciation au niveau phonème (projet TypALoc)
 - Alignement de la parole lue des patients au niveau phonème
 - Décision au niveau phonème de la normalité ou non de la production

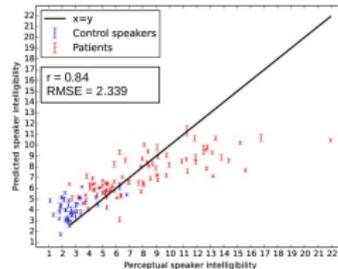
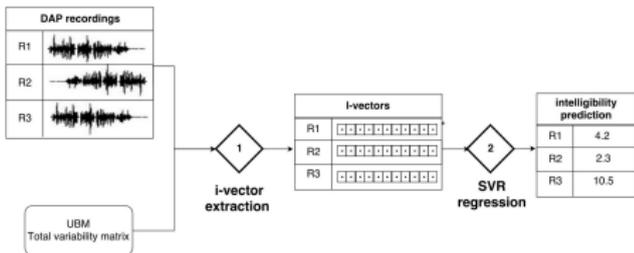
①	②	③	④	⑤	⑥	⑦																																																																																																																																																																		
<table border="0"> <tr><td>■</td><td>dans</td></tr> <tr><td>□</td><td>un</td></tr> <tr><td>■</td><td>petit</td></tr> <tr><td>■</td><td>village</td></tr> <tr><td>■</td><td>de</td></tr> <tr><td>■</td><td>la</td></tr> <tr><td>■</td><td>montagne</td></tr> </table>	■	dans	□	un	■	petit	■	village	■	de	■	la	■	montagne	<table border="0"> <tr><td>□</td><td>il</td></tr> <tr><td>□</td><td>y</td></tr> <tr><td>□</td><td>a</td></tr> <tr><td>■</td><td>un</td></tr> <tr><td>■</td><td>pauvre</td></tr> <tr><td>■</td><td>cordonni</td></tr> <tr><td>□</td><td>tout</td></tr> <tr><td>■</td><td>vieux</td></tr> <tr><td>□</td><td>et</td></tr> <tr><td>■</td><td>tout</td></tr> <tr><td>■</td><td>casse1</td></tr> </table>	□	il	□	y	□	a	■	un	■	pauvre	■	cordonni	□	tout	■	vieux	□	et	■	tout	■	casse1	<table border="0"> <tr><td>■</td><td>les</td></tr> <tr><td>■</td><td>villageois</td></tr> <tr><td>■</td><td>lui</td></tr> <tr><td>■</td><td>apporten</td></tr> <tr><td>■</td><td>dés</td></tr> <tr><td>■</td><td>chaussur</td></tr> <tr><td>■</td><td>a2</td></tr> <tr><td>■</td><td>re1parer</td></tr> <tr><td>■</td><td>mafs</td></tr> <tr><td>□</td><td>il</td></tr> <tr><td>■</td><td>ne</td></tr> <tr><td>■</td><td>travaille</td></tr> <tr><td>□</td><td>pas</td></tr> <tr><td>□</td><td>vite</td></tr> </table>	■	les	■	villageois	■	lui	■	apporten	■	dés	■	chaussur	■	a2	■	re1parer	■	mafs	□	il	■	ne	■	travaille	□	pas	□	vite	<table border="0"> <tr><td>□</td><td>tous</td></tr> <tr><td>□</td><td>les</td></tr> <tr><td>■</td><td>soirs</td></tr> <tr><td>□</td><td>il</td></tr> <tr><td>■</td><td>mange</td></tr> <tr><td>□</td><td>tout</td></tr> <tr><td>■</td><td>seul</td></tr> <tr><td>■</td><td>bien</td></tr> <tr><td>■</td><td>tristemen</td></tr> </table>	□	tous	□	les	■	soirs	□	il	■	mange	□	tout	■	seul	■	bien	■	tristemen	<table border="0"> <tr><td>■</td><td>ce</td></tr> <tr><td>■</td><td>soir</td></tr> <tr><td>□</td><td>il</td></tr> <tr><td>□</td><td>y</td></tr> <tr><td>□</td><td>a</td></tr> <tr><td>□</td><td>devant</td></tr> <tr><td>□</td><td>lui</td></tr> <tr><td>□</td><td>un</td></tr> <tr><td>□</td><td>gros</td></tr> <tr><td>□</td><td>tas</td></tr> <tr><td>□</td><td>de</td></tr> <tr><td>■</td><td>souliers</td></tr> <tr><td>□</td><td>et</td></tr> <tr><td>□</td><td>de</td></tr> <tr><td>■</td><td>cue3tres</td></tr> <tr><td>■</td><td>a2</td></tr> <tr><td>■</td><td>recoudre</td></tr> </table>	■	ce	■	soir	□	il	□	y	□	a	□	devant	□	lui	□	un	□	gros	□	tas	□	de	■	souliers	□	et	□	de	■	cue3tres	■	a2	■	recoudre	<table border="0"> <tr><td>■</td><td>jamais</td></tr> <tr><td>■</td><td>le</td></tr> <tr><td>■</td><td>ne</td></tr> <tr><td>■</td><td>pourrai</td></tr> <tr><td>■</td><td>les</td></tr> <tr><td>■</td><td>re1parer</td></tr> <tr><td>■</td><td>le</td></tr> <tr><td>■</td><td>suis</td></tr> <tr><td>■</td><td>a3ge1</td></tr> <tr><td>□</td><td>et</td></tr> <tr><td>■</td><td>trop</td></tr> <tr><td>■</td><td>malade</td></tr> </table>	■	jamais	■	le	■	ne	■	pourrai	■	les	■	re1parer	■	le	■	suis	■	a3ge1	□	et	■	trop	■	malade	<table border="0"> <tr><td>■</td><td>pre2s</td></tr> <tr><td>□</td><td>de</td></tr> <tr><td>□</td><td>lui</td></tr> <tr><td>□</td><td>la</td></tr> <tr><td>■</td><td>grosse</td></tr> <tr><td>□</td><td>horloge</td></tr> <tr><td>□</td><td>rait</td></tr> <tr><td>□</td><td>tic</td></tr> <tr><td>□</td><td>tac</td></tr> <tr><td>□</td><td>tic</td></tr> <tr><td>□</td><td>tac</td></tr> </table>	■	pre2s	□	de	□	lui	□	la	■	grosse	□	horloge	□	rait	□	tic	□	tac	□	tic	□	tac
■	dans																																																																																																																																																																							
□	un																																																																																																																																																																							
■	petit																																																																																																																																																																							
■	village																																																																																																																																																																							
■	de																																																																																																																																																																							
■	la																																																																																																																																																																							
■	montagne																																																																																																																																																																							
□	il																																																																																																																																																																							
□	y																																																																																																																																																																							
□	a																																																																																																																																																																							
■	un																																																																																																																																																																							
■	pauvre																																																																																																																																																																							
■	cordonni																																																																																																																																																																							
□	tout																																																																																																																																																																							
■	vieux																																																																																																																																																																							
□	et																																																																																																																																																																							
■	tout																																																																																																																																																																							
■	casse1																																																																																																																																																																							
■	les																																																																																																																																																																							
■	villageois																																																																																																																																																																							
■	lui																																																																																																																																																																							
■	apporten																																																																																																																																																																							
■	dés																																																																																																																																																																							
■	chaussur																																																																																																																																																																							
■	a2																																																																																																																																																																							
■	re1parer																																																																																																																																																																							
■	mafs																																																																																																																																																																							
□	il																																																																																																																																																																							
■	ne																																																																																																																																																																							
■	travaille																																																																																																																																																																							
□	pas																																																																																																																																																																							
□	vite																																																																																																																																																																							
□	tous																																																																																																																																																																							
□	les																																																																																																																																																																							
■	soirs																																																																																																																																																																							
□	il																																																																																																																																																																							
■	mange																																																																																																																																																																							
□	tout																																																																																																																																																																							
■	seul																																																																																																																																																																							
■	bien																																																																																																																																																																							
■	tristemen																																																																																																																																																																							
■	ce																																																																																																																																																																							
■	soir																																																																																																																																																																							
□	il																																																																																																																																																																							
□	y																																																																																																																																																																							
□	a																																																																																																																																																																							
□	devant																																																																																																																																																																							
□	lui																																																																																																																																																																							
□	un																																																																																																																																																																							
□	gros																																																																																																																																																																							
□	tas																																																																																																																																																																							
□	de																																																																																																																																																																							
■	souliers																																																																																																																																																																							
□	et																																																																																																																																																																							
□	de																																																																																																																																																																							
■	cue3tres																																																																																																																																																																							
■	a2																																																																																																																																																																							
■	recoudre																																																																																																																																																																							
■	jamais																																																																																																																																																																							
■	le																																																																																																																																																																							
■	ne																																																																																																																																																																							
■	pourrai																																																																																																																																																																							
■	les																																																																																																																																																																							
■	re1parer																																																																																																																																																																							
■	le																																																																																																																																																																							
■	suis																																																																																																																																																																							
■	a3ge1																																																																																																																																																																							
□	et																																																																																																																																																																							
■	trop																																																																																																																																																																							
■	malade																																																																																																																																																																							
■	pre2s																																																																																																																																																																							
□	de																																																																																																																																																																							
□	lui																																																																																																																																																																							
□	la																																																																																																																																																																							
■	grosse																																																																																																																																																																							
□	horloge																																																																																																																																																																							
□	rait																																																																																																																																																																							
□	tic																																																																																																																																																																							
□	tac																																																																																																																																																																							
□	tic																																																																																																																																																																							
□	tac																																																																																																																																																																							

Thèmes de recherche

- Prédiction automatique de l'intelligibilité de patients atteints de la dysarthrie et de cancer de la tête et du cou (projet C2SI)
 - Utilisation des i-vecteur adaptés de la reconnaissance du locuteur pour prédire intelligibilité des patients et la sévérité de la maladie
 - Utilisation de différentes annotations automatiques (reconnaissance au niveau phonème, score de vraisemblance, confusion phonémique) pour la prédiction de l'intelligibilité

Thèmes de recherche

- Prédiction automatique de l'intelligibilité de patients atteints de la dysarthrie et de cancer de la tête et du cou (projet C2SI)
 - Utilisation des i-vecteur adaptés de la reconnaissance du locuteur pour prédire intelligibilité des patients et la sévérité de la maladie
 - Utilisation de différentes annotations automatiques (reconnaissance au niveau phonème, score de vraisemblance, confusion phonémique) pour la prédiction de l'intelligibilité



Thèmes de recherche

- Utilisation de la reconnaissance automatique de la parole pour le réglage des prothèses auditives (projet Phonics)
 - Simulation des effets des traumatismes sonores sur la parole
- Recherche de paramètres distinctifs automatique entre les patients atteints de la maladie de Parkinson et l'Atrophie Multi-Systématisée AMS (projet Voice4PD-MSA)
 - Étude du rythme de la parole chez les patients
 - Étude de la confusion phonémique et la qualité de prononciation comme marqueur de la pathologie

Yoann Dupont

Entités nommées : systèmes et données

Yoann Dupont

Les entités nommées (EN)

Japon : la droite au pouvoir remporte les élections

22 octobre 2017. – Le **Parti libéral-démocrate** de **Shinzō Abe** a remporté largement les élections législatives organisées ce dimanche. Le Premier ministre pourrait ainsi théoriquement occuper ce poste jusqu'en 2021, alors qu'il est au pouvoir depuis 2012. S'il y parvient, il battra alors le record de longévité en tant que chef du gouvernement, dans un pays habitué à l'instabilité politique. **Shinzō Abe** a réussi son pari d'orienter la campagne électorale vers les menaces étrangères, à savoir la **Chine** et la **Corée du Nord**, dans un contexte de tensions militaires très fortes. Le Premier ministre entend réformer la constitution japonaise pour supprimer son caractère pacifique, et développer l'armée nationale.

Son parti de la droite conservatrice réussit à décrocher 311 sièges sur les 465 que compte la **Chambre des représentants**. S'il perd quelques sièges, il conserve une large majorité absolue. Deux coalitions s'opposaient à lui : l'une de gauche et pacifiste, qui obtient 67 sièges soit presque deux fois plus que dans la législature sortante. La coalition libérale « **Koike** » obtient elle 58 sièges, treize de moins qu'auparavant.



Figure 1: source: <https://fr.wikinews.org>

Les entités nommées (EN)

Japon : la droite au pouvoir remporte les élections

22 octobre 2017. – Le **Parti libéral-démocrate** de **Shinzō Abe** a remporté largement les élections législatives organisées ce dimanche. Le Premier ministre pourrait ainsi théoriquement occuper ce poste jusqu'en 2021, alors qu'il est au pouvoir depuis 2012. S'il y parvient, il battra alors le record de longévité en tant que chef du gouvernement, dans un pays habitué à l'instabilité politique. **Shinzō Abe** a réussi son pari d'orienter la campagne électorale vers les menaces étrangères, à savoir la **Chine** et la **Corée du Nord**, dans un contexte de tensions militaires très fortes. Le Premier ministre entend réformer la constitution japonaise pour supprimer son caractère pacifique, et développer l'armée nationale.

Son parti de la droite conservatrice réussit à décrocher 311 sièges sur les 465 que compte la **Chambre des représentants**. S'il perd quelques sièges, il conserve une large majorité absolue. Deux coalitions s'opposaient à lui : l'une de gauche et pacifiste, qui obtient 67 sièges soit presque deux fois plus que dans la législature sortante. La coalition libérale « **Koike** » obtient elle 58 sièges, treize de moins qu'auparavant.



Figure 1: source: <https://fr.wikinews.org>

Apprentissage : système entraîné à reproduire une tâche à partir d'exemples.

Les entités nommées (EN)

Japon : la droite au pouvoir remporte les élections

22 octobre 2017. – Le Parti libéral-démocrate de Shinzō Abe a remporté largement les élections législatives organisées ce dimanche. Le Premier ministre pourrait ainsi théoriquement occuper ce poste jusqu'en 2021, alors qu'il est au pouvoir depuis 2012. S'il y parvient, il battra alors le record de longévité en tant que chef du gouvernement, dans un pays habitué à l'instabilité politique. Shinzō Abe a réussi son pari d'orienter la campagne électorale vers les menaces étrangères, à savoir la Chine et la Corée du Nord, dans un contexte de tensions militaires très fortes. Le Premier ministre entend réformer la constitution japonaise pour supprimer son caractère pacifique, et développer l'armée nationale.

Son parti de la droite conservatrice réussit à décrocher 311 sièges sur les 465 que compte la Chambre des représentants. S'il perd quelques sièges, il conserve une large majorité absolue. Deux coalitions s'opposaient à lui : l'une de gauche et pacifiste, qui obtient 67 sièges soit presque deux fois plus que dans la législature sortante. La coalition libérale « Koike » obtient elle 58 sièges, treize de moins qu'auparavant.



Figure 1: source: <https://fr.wikinews.org>

Apprentissage : système entraîné à reproduire une tâche à partir d'exemples.

Exemples d'applications :

- extraction d'informations (qui ? quoi ? quand ? où etc...)
- systèmes de question-réponse
- anonymisation

Travaux en cours / prévus : structuration

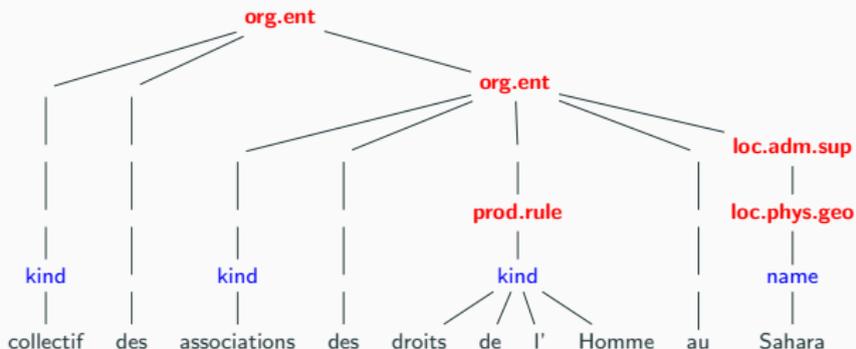
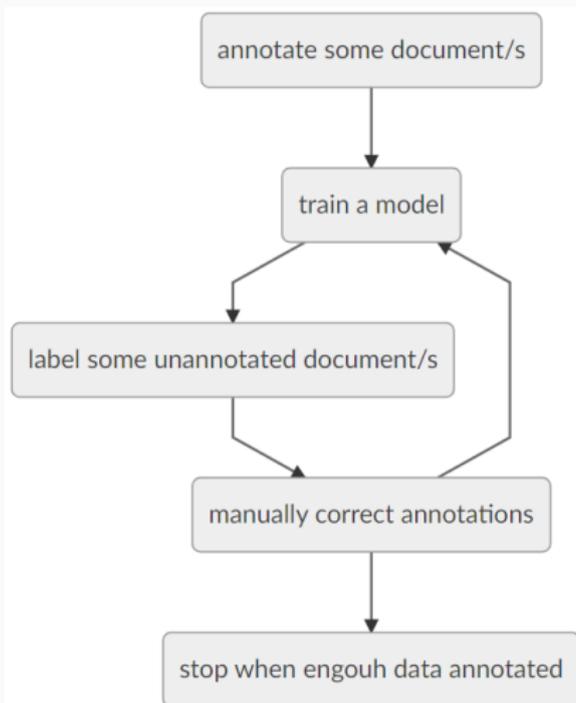


Figure 2: une entité nommée structurée. Source : corpus Quaero Grouin *et al.* (2011)

Travaux en cours / prévus : apprentissage actif



Travaux en cours / prévus : BERT Devlin et al. (2018)

Use the output of the masked word's position to predict the masked word

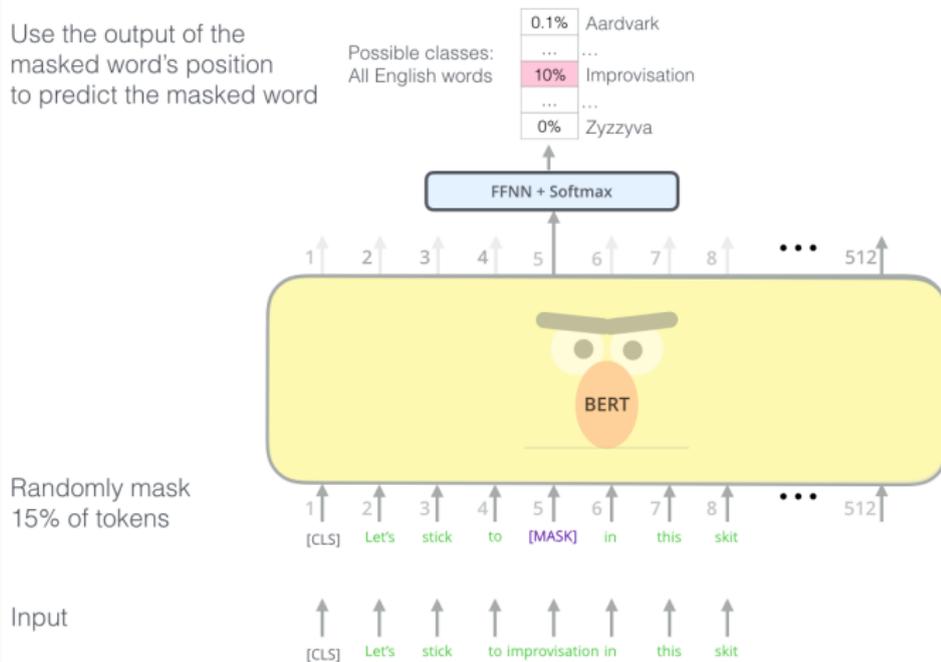


Figure 3: source : <https://jalamar.github.io/illustrated-bert/>

References

DEVLIN J., CHANG M.-W., LEE K. & TOUTANOVA K. (2018). **Bert: Pre-training of deep bidirectional transformers for language understanding.** *arXiv preprint arXiv:1810.04805*.

GROUIN C., ROSSET S., ZWEIGENBAUM P., FORT K., GALIBERT O. & QUINTARD L. (2011). **Proposal for an extension of traditional named entities: From guidelines to evaluation, an overview.** In *Proceedings of the 5th Linguistic Annotation Workshop*, p. 92–100: Association for Computational Linguistics.

Zijian Wang



Comparaison des résultats de différents outils pour la détection des entités nommées

Zijian WANG Master 2 Langue et Informatique

September 9, 2020

Évaluation des résultats

- *Détection des entités nommées par des outils existants*
- *Ré-entraînement des modèles*
- *Évaluation par diverses façons*
- *Visualisation des résultats pour ces évaluations*

Quelle est votre destination?

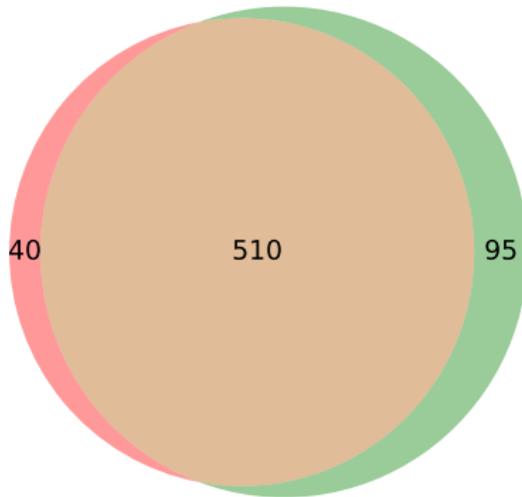
Paris LOC

. Quels sites voulez-vous visiter? La

Tour Eiffel LOC

.

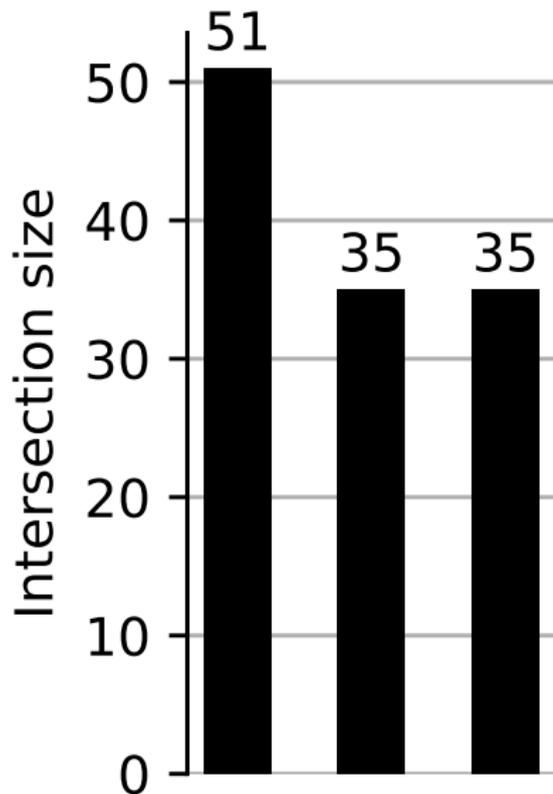
- Ce sont des syntagmes ne portant qu'un sens dans la langue.
- Ce sont des attributs pertinents pour l'apprentissage d'une phrase.



Spacy sm

Spacy md

Which is the best for ORG?

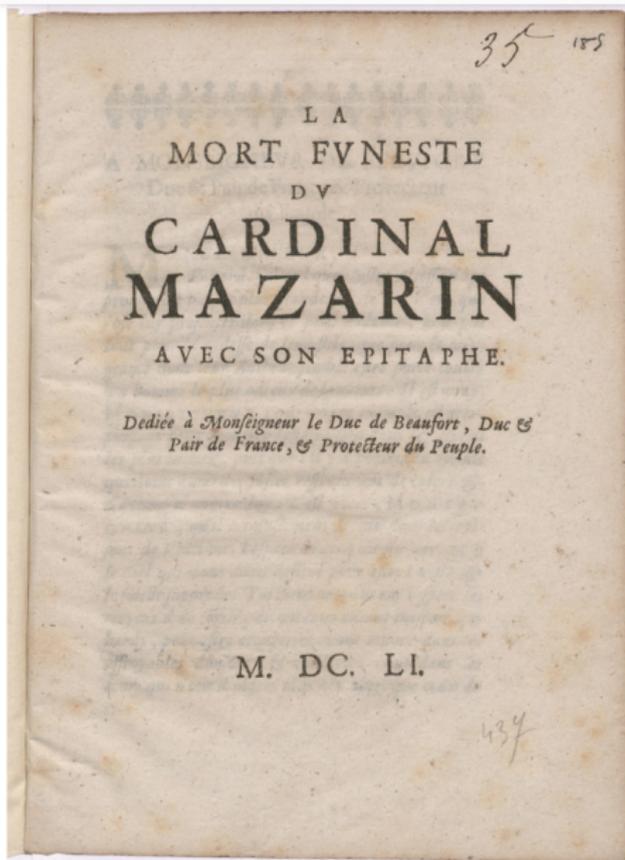


25

Julien Bezancon

stage.pdf

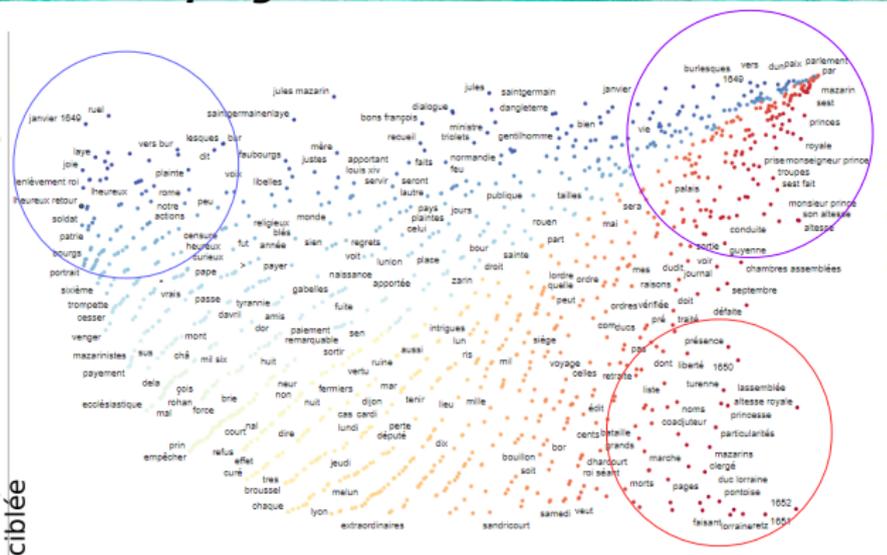
Julien
Bezançon



Compte-rendu
de stage de L3

Graphique obtenu avec le programme créé

Fréquence des mots de la catégorie
ciblée



Top 1649
janvier 1649
ruei
jules mazarin
jules
saintgermainlaye
vers burlesques
conférence
burlesques
dialogue
pour bien
saintgermain
1649
laye
bons français

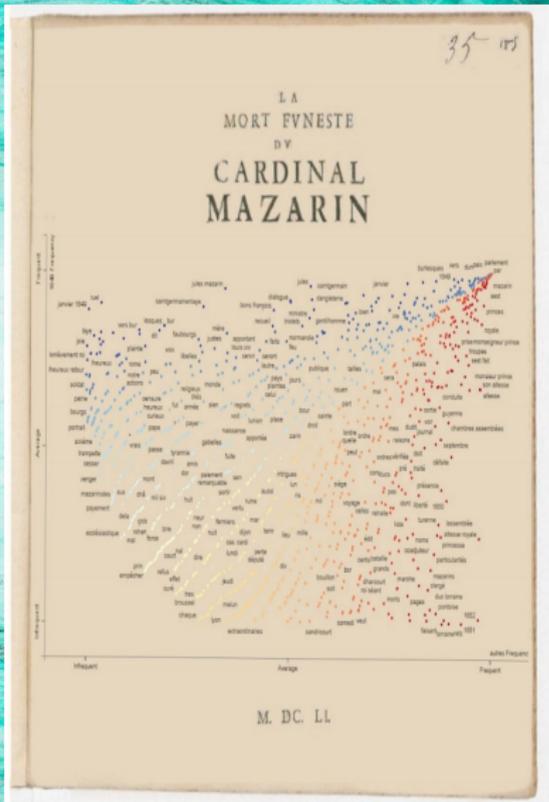
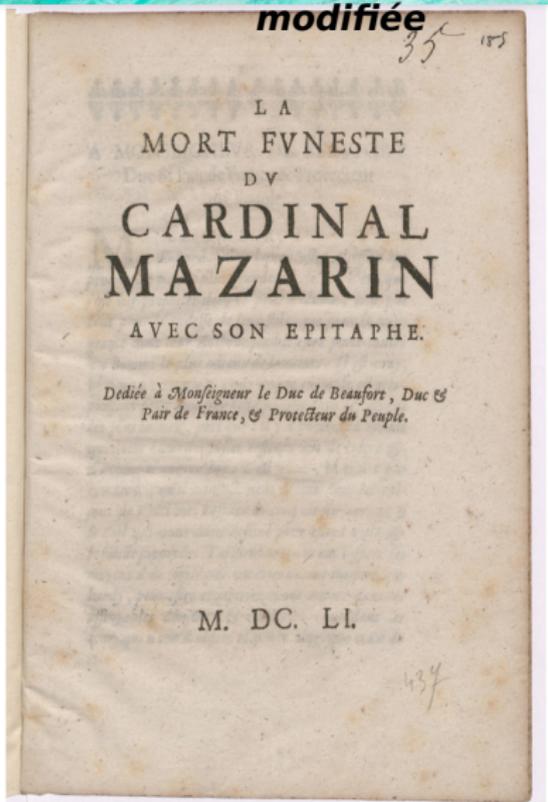
Characteristic
mazarin
parlement
messieurs
sest
altesse
monseigneur
faite
burlesques
touchant
ladite
sieur
remembrance
lettre
quil
contenant
conty
quils
jusques
seigneurs
sieurs
harangue
paix
officiers
manifeste
reine
saintgermain
dangleterre
guyenne
mazarins
parisiens

Top autres
1652
1651
retz
pontoise
lorraine
altesse royale
cardinal retz
duc lorraine
altesse
son altesse
faisant
jun 1652
clergé
particularités

Fréquence des mots dans les autres catégories

1649 document count: 1,577; word count: 20,580
autres document count: 2,789; word count: 47,225

**Couverture originale / couverture
modifiée**



Caroline Parfait

Speed dating recherche, Septembre 2020

contrat doctoral SCAI



Analyse de l'espace littéraire : apprentissage automatique et évaluation des systèmes de reconnaissance des entités nommées

Caroline Parfait, Humanités Numériques



Objectifs

- Dépasser les problèmes de **variabilités**
- Utiliser le potentiel des outils de **machine learning pour la REN** dans la Littérature

Livrables :

- Un **modèle adaptatif** d'outils numériques pour la REN spatiales
- Un **corpus standard** annoté
- un **guide d'annotation** pour les Humanités Numériques littéraires

Méthodologie proposée : “conception centrée utilisateur”

- **Identification des besoins** et des contextes d'utilisation → penser l'interdisciplinarité
- **Évaluer** les modèles **existants** pour la REN spatiale :
 - Corpus : Récit de voyages du XIXème et XXème siècle
 - **ANR Chapitre**, Paris 3 Sorbonne Nouvelle → 3000 documents
 - Corpus de la **Très Grande Bibliothèque** → 35000 documents
- **Modéliser** un nouvel outil
- **Évaluation** auprès des **utilisateurs**

Perspectives

- **Applicabilité transversale** -> extraction de concepts médicaux
- Entériner des **analyses littéraires à partir des résultats** (stylométrie, analyses comparatives, analyses des sentiments et des émotions...)
- Exploiter les **technologies du web sémantique pour la désambiguïsation** des EN (data linking)

Merci de votre attention

Gaël Lejeune

Qui suis-je ?

- Gaël Lejeune MCF en informatique
- Thèse en 2013 : Veille Epidémiologique Multilingue

Et depuis ?

Différentes tâches :

- Extraction de Contenu/de structure (PDF, PNG, HTML...)
- Classification (polarité, émotion, dialectes, datation ...)
- Extraction d'Information et Reconnaissance d'Entités Nommées
- Indexation/Terminologie
- Evaluation Intrinsèque VS Extrinsèque (plongements...)

Qui suis-je ?

- Gaël Lejeune MCF en informatique
- Thèse en 2013 : Veille Epidémiologique Multilingue

Et depuis ?

Différentes tâches :

- Extraction de Contenu/de structure (PDF, PNG, HTML...)
- Classification (polarité, émotion, dialectes, datation ...)
- Extraction d'Information et Reconnaissance d'Entités Nommées
- Indexation/Terminologie
- Evaluation Intrinsèque VS Extrinsèque (plongements...)

Dans un contexte de variation dans les données

- Multilinguisme : comment analyser n langues ?
- Hétérogénéité : comment traiter n états de textes ?
- Massification : comment travailler sur n To de textes ?

Qui suis-je ?

- Gaël Lejeune MCF en informatique
- Thèse en 2013 : Veille Epidémiologique Multilingue

Et depuis ?

Différentes tâches :

- Extraction de Contenu/de structure (PDF, PNG, HTML...)
- Classification (polarité, émotion, dialectes, datation ...)
- Extraction d'Information et Reconnaissance d'Entités Nommées
- Indexation/Terminologie
- Evaluation Intrinsèque VS Extrinsèque (plongements...)

Dans un contexte de variation dans les données

- Multilinguisme : comment analyser n langues ?
- Hétérogénéité : comment traiter n états de textes ?
- Massification : comment travailler sur n To de textes ?

Qu'est-ce qui m'intéresse ?

Il n'existe pas de telle chose qu'une donnée parfaite

Tout pré-traitement amène son lot de désagréments

- Enlever les ponctuations, pourquoi ?
- Découper en mots, pourquoi ?
- Découper en phrases, pourquoi ?
- → Tendance de l'informatique à javelliser/uniformiser les données et les approches

Sur quoi je travaille ?

Projets en cours

- MEMES (2019-2021) Memes : Extraction automatique et analyse par Myriadisation d'Expressions Semi-figées (K.Fort, A.Gautier, L.Zhu)
- ANTONOMAZ (2018-2022) ANalyse auTOMatique et NumérisatiOn des MAZarinades (K.Abiven, G.Roe , JB.Tanguyet *al.*)
- OBVIL-NER (2020-2023) Apprentissage automatique et évaluation des systèmes de reconnaissance des entités nommées (C. Parfait, M. Alrahabi, G. Roe)
- WADDLE et DANIEL (2018-...) Exploitation de Données Textuelles hétérogènes (A. Barbaresi, E. Giguet)
- SINNER (2020-...) Extraction d'Entités Nommées en contexte Multilingue et Bruité (Y.Dupont, P. Ortiz, T. Tian)

Fouad Aounti

Détection d'objets dans des images médiévales

Fouad AOUINTI

Maison de la Recherche
Sorbonne-Université, Paris

10 septembre 2020

Préparation de données

- Dataset original, environ de 1000 images :
 - train set : 80% d'images ;
 - test set : le reste 20%.



- Techniques d'augmentation : scale, rotation, zoom, flip, etc.
- 2 classes : *graphic*, *no-graphic*.
- Deep Learning en Python : Keras avec Tensorflow.

Classification d'images

- Modèles CNN^a pré-entraînés : VGG16, ResNet, etc.
- Stratégie du Transfer Learning : entraîner tout le modèle.

a. Convolutional Neural Network

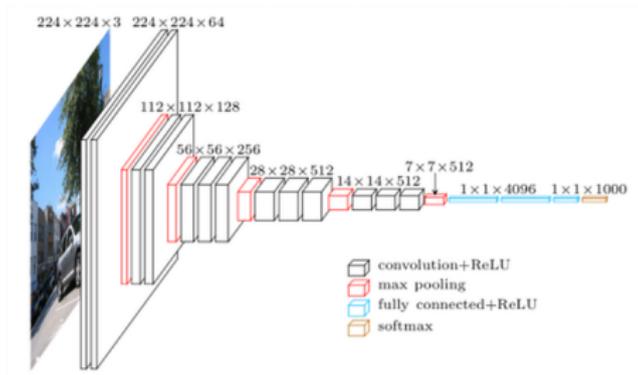
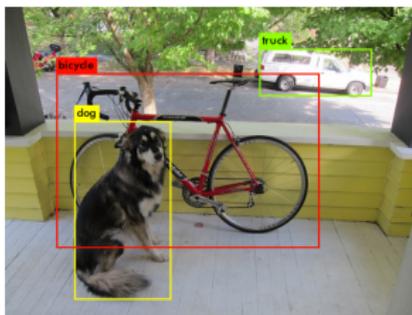


FIGURE 1 – Architecture du VGG16.

Détection d'objets



(a)



(b)

FIGURE 2 – (a) Image réelle, (b) Images médiévales.

Jean Baptiste Tanguy

Une thèse en humanités numériques

Encrages linguistiques de la politisation de l'imprimé pendant la Fronde (1648-1653) : la masse polémique des Mazarinades

Directeur : Glenn Roe

Encadrants : Karine Abiven et Gaël Lejeune

Rattachements : LabEx OBVIL et équipe STIH

Partenaire : Bibliothèque Mazarine

Financement : Région Île-de-France



Objet d'étude : les Mazarinades

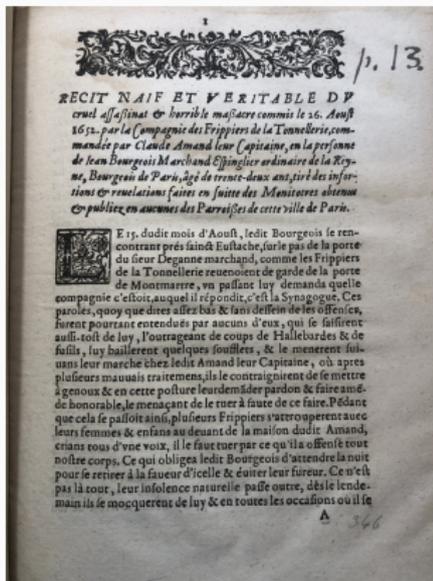


Figure 1 – *Récit naïf et véritable...*, s.l.n.d., 7 pages, Bibl. Maz. M. 12232

- Entre 5000-6000 Mazarinades : plusieurs bibliographies existent mais aucune n'est/ne peut être exhaustive
- Elles nous sont mystérieuses : date, lieu, auteur, conditions de parution, de tirage, de lecture, contexte...
- Pamphlets, actes de la cours, récits de guerre ou de fête, lettres... parus en France pendant la Fonde (1648-1653) : politisation du petit imprimé (C. Jouhaud)
- Une multitude de formes qui se **voient** et **s'entendent** dans la ville à l'actualité affolée
- Français pré-classique (ex. : u=v, i=j, présence du s long, etc.) mais aussi latin, italien, patois
- Exemplaires datant du XVIIème siècle, donc : tâches, restaurations, etc.

Axes de recherche

- Identifier des petits corpus cohérents :
 - Des "faisceaux" (C. Jouhaud), des réseaux, etc. : travail des historiens
 - Réaliser du *clustering* pour rapprocher des pièces jusqu'alors restées isolées
- Construire des corpus numériques :
 - Acquisition des données (sur documents historiques) : numérisation, OCR, évaluation (transcription)
 - Mise à disposition pérenne des données : consortium CAHIER et FAIRisation des données
- Étude des phénomènes linguistiques mis en jeu dans la politisation de l'imprimé :
 - Néologie, isotopie, phraséologie, rhétorique...
 - Linguistique sur corpus, à partir de données textuelles océrisées, donc bruitées

Figure 1 – Source : Jeff Stahler 2014

