

Séminaire de rentrée : Recontres Minute Science

09 septembre 2021

STIH, Sorbonne University

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

Alexandre Bartz

Alexandre Bartz

Sorbonne Université

9 septembre 2021

Projet Antonomaz

- Master TNAH (*Technologies numériques appliquées à l'histoire*) de l'École nationale des chartes (Paris)
- Depuis avril 2021, ingénieur du projet Antonomaz

Antonomaz

ANalyse au**TO**matique et **N**umérisasi**ON** des **MAZ**arinades

- Exploiter un corpus d'environ 5.000 libelles rédigés au milieu du XVIIe siècle, sous la Fronde
- But : analyser automatiquement ce corpus et le mettre à la disposition des chercheuses et chercheurs

Mes missions

- Encodage des documents en XML-TEI, contrôle qualité grâce à des schémas d'encodage (ODD)
- Travail sur les métadonnées
- Mettre à disposition les données : publication sur Github, aussi bien de la chaîne de production que des textes encodés

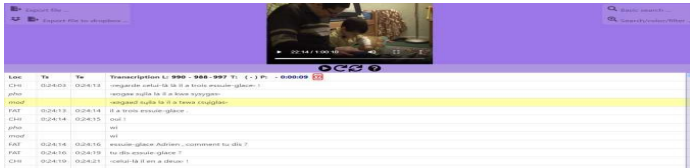
- Rendre les données textuelles produites exploitables et interrogeables par tout.e.s : lemmatisation, annotation morphosyntaxique, normalisations et *NER* (en collaboration avec l'INRIA, Simon Gabay et Philippe Gambette)
- Proposer une application web pour mettre les données en valeur
- Projet secondaire : reconnaissance automatique de l'imprimeur grâce à la casse utilisée (voir <https://github.com/Imprinters> en collaboration avec Simon Gabay)

Andrea Briglia

Andrea Briglia, PhD en Linguistique, 10/2017- 3/2021 (cotutelle Messina & Montpellier)

« Approches statistiques et computationnelles à l'acquisition du français L1 chez l'enfant »

Données CoLaE (ANR, Morgenstern & Parisse, 2012). 6 enfants enregistrés une heure par mois, tous les mois, d'un an jusqu'à cinq ans. Transcriptions en trois types différents.



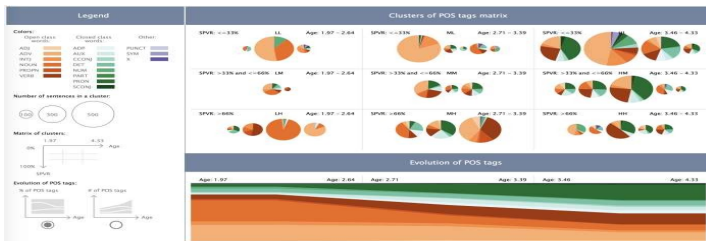
The screenshot shows a video player interface with a purple header. The video content displays a child's face. Below the video, there is a transcription table with columns for speaker, start time, end time, and the transcription text. The table contains several rows of data, with some rows highlighted in yellow.

Loc.	Ts	Te	Transcription L: 990 - 998 - 997 F: (-) Ph. - 0:00:09
Chi	0:2:4:03	0:2:4:12	regarde celui-là il a trois ans et demi
John			regarde celui-là il a trois ans et demi
moof			regarde celui-là il a trois ans et demi
FAT	0:2:4:12	0:2:4:14	il a trois ans et demi
Chi	0:2:4:14	0:2:4:15	oui
John			oui
moof			oui
FAT	0:2:4:14	0:2:4:16	regarde celui-là il a trois ans et demi
FAT	0:2:4:16	0:2:4:19	regarde celui-là il a trois ans et demi
Chi	0:2:4:19	0:2:4:21	regarde celui-là il a trois ans et demi

En total, environ 8'000 énoncés pour chaque enfant. Test de représentativité statistique (capture rate) 1%

HP Suite des variation phonético/phonologiques, étapes acquisition du fr L1, tendances et contraintes, comparaison intra-enfant et inter-enfants. Deux niveaux d'analyse : phonétique et syntaxique.

Méthodes Chi2, POS tagging (stanza, UD), clusterisation avec EM, visualisation avec Multistreamgraph, modélisation et reproduction avec un CNN



<http://advanse.lirmm.fr/EMClustering/>

Réseau neuronal de type convolutionnel : https://colab.research.google.com/drive/1fa0ak1kyWfmsCxtfZpl6vYtEY_Ppw5

Limites données bruitées, prosodie et non-verbal, langage adulte, distinction attaque/coda et coarticulations, difficulté de représenter la phonologie en utilisant des séquences des caractères... et en plus, on constate un grande variabilité entre les enfant (exceptions, retour en arrière...)

Conclusions « n'importe quelle variation ne varie pas en n'importe quelle autre » (Sauvage, 2015), cadre interprétatif ; théorie des traits phonologique de N. Clements.

∅ /gR/ → /dR/

∅ /dR/ → /gR/

∅ /tR/ → /kR/

∅ /kR/ → /tR/

∅ « 52 Q : ben attends, on essaie de/ de **l'trouver**, si on le **trouve** pas ze/ ben c'est pas grave hein, ça c'est un gros euh c'est bien... si on le **krouve** pas, alors c'est pas grave...donc, [...] »

∅ Jérémi « c'est [dra] ça ? » prononcé exprès à la place de « c'est grand ça ? » dans le but d'étudier la réaction de l'enfant, qui répond ainsi « on dit pas [dra], on dit [dra] !! ». (Sauvage, 2015)

Là, on trouve toute la difficulté d'étudier la conscience phonologique en cours de développement

Caroline Parfait

Reconnaissance des Entités Nommées (REN) spatiales dans un corpus littéraire : robustesse des systèmes existants ?

Caroline Parfait (1,2) Gaël Lejeune (2)
Motasem Alrahabi (1) Glenn Roe (1)

Speed Dating (Vos travaux de recherche en 3mn), 2021

9 Septembre 2021



Sens Texte
Informatique
Histoire



caroline.parfait@sorbonne-universite.fr
gael.lejeune@sorbonne-universite.fr



REN : Vocabulaire & application concrète

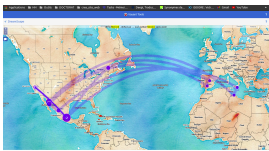


Figure: Voyant Tools, "Les trappeurs de l'Arkansas (4e édition)", Gustave Aimard, 1858.

- Humanités Numériques (HN)
- Entités nommées spatiales (ENS)= Toponymes
- Reconnaissance Optique de Caractères (OCR)
- Bruit(s) :
 - OCR = Ajout, transformations de caractères, autres (?)
 - REN = NON ENS
- Silence(s?) :
 - OCR = Suppression de caractères (?)
 - REN = NON-Reconnaissance ENS (VP, FP, FN, etc.)

Problématique(s) et Méthode

Comment rendre les **systèmes de NER plus robustes** face aux **variations** dans les données qui leur sont soumises par les **utilisateurs** ?

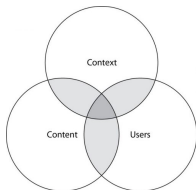


Figure: Peter Morville, "Three Circles of Information Architecture", 2004

- **Utilisateurs** : chercheurs SHS, HN, TAL ... Interdisciplinarité
- **Données** : Corpus ELTeC (Romans Français 19/20ème s.)
- **Variabilités** : diachronie, diatopie, **qualité de la transcription OCR ?**

- Enquêtes **Utilisateurs** quantitative et qualitative avec J-B. Tanguy
- **Évaluation des systèmes** SPACY et STANZA
- **Angle Expé.** : "Impact du bruit des transcriptions OCR sur la REN ?"

Exemples :

Version	Context	Output
Ref.	<i>prendre la diligence de Châlons</i>	Châlons
Kraken	<i>prendre la diligence de Ch_lons</i>	Ch_lons
Tess	<i>prendre la diligence de Chalons</i>	Chalons
Tess fr	<i>prendre la diligence de Châlons</i>	Châlons

Eva Lacroix



Eva Schaeffer-Lacroix

Maitresse de conférences HDR (sections 7 et 12)

Inspé de Paris (Sorbonne Université)



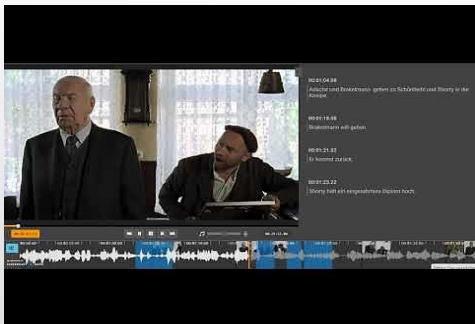
elacroix@inspe-paris.fr
<http://didaktik.hautefort.com/>

Mon parcours

- Magister artium Germanistik & Romanistik, Universität Konstanz (1992)
- Professeure agrégée d'allemand (1994-2011)
- Master en didactique des langues et cultures, Paris 3 (2004-2005)
- Thèse en sciences du langage, sections 7 et 12, Paris 3 (2006-2009) :
Corpus numériques et production écrite en allemand langue étrangère. Une recherche avec des apprenants d'allemand.
- Maitresse de conférences à l'Inspé de Paris depuis 2011
- Habilitation en sciences du langage, sections 7 et 12, Université de Strasbourg (2020) :
Recherches et pratiques outillées en allemand et français - Didactique, linguistique de corpus et traduction.

Mes recherches

- [Mymap](#)
- Recherche à dominante engagée
- Expérience "usager" : observer les actions des personnes et ce qu'elles en disent
- Outillage des apprentissages et analyses
- Collaboratrice scientifique au département d'allemand de l'université de Genève : linguiste de corpus & projet FNS (Fonds national suisse) sur la politique de la langue en Suisse (mars 2015 – octobre 2016)
- Initiatrice d'une collaboration franco-allemande ayant mené à un dépôt de projet ANR-DFG sur la traductibilité de scripts d'audiodescription en mars 2021



Projet TADS (Traduction de scripts d'audiodescription)

Annotation XML-TEI

element	tag XML-TEI	exemple Buettewarder (épisode 12)
indications de temps	<time>	10: 06: 50
parties à enregistrer (<i>speaker</i>)	<sp>	Brakelmann fährt durchs Dorf. [Brakelmann traverse le village.]
instructions concernant les actions à faire ou à éviter pendant l'enregistrement	<stage>	(Rest übersprechen) [Chevaucher le reste]
débit	<stage type="delivery">	s, ss, n [rapide, très rapide, normal]
légendes qui s'affichent à l'écran (titres, textes de pancartes, etc.)	<caption>	Schild am Straßenrand. "Goethe Eier 5 Euro 10" [Panneau au bord de la route. "Oeufs à la manière de Goethe 5 euros 10]
(fins de) répliques	<prompt>	"Goethe gehört mir!" [Goethe est à moi !]

EXMARaLDA

- [1]
- | | | | |
|---------|--------------------------|------------------------|-------------------------------|
| | 0 [00:00.0] 1 [00:15.3] | 2 [00:16.62] [00:28.8] | |
| SP [v] | Drei verkleidete Männer. | | Zwei Männer mit Händen in den |
| Ads [v] | | | Ja, ich finde das ja |
- [2]
- | | | | |
|-----------|-----------------------------------|-------------|------------------|
| | 4 [00:45.7] | 5 [00:42.5] | 6 [00:43.3] |
| SP [v] | Taschen. | | |
| Ads [v] | immer schön, wenn einer arbeitet. | | Ja , da kann ich |
| Bra [v] | | | Tatsache ? |
| sound [v] | HAMMERN | | |
- [3]
- | | | | |
|-----------|--------------------------------------|-------------------------|----------------------|
| | 7 [00:45.6] | 8 [00:46.639] [00:47.2] | |
| Ads [v] | stundenlang zukucken , stundenlang ! | | Ja, wenn er dan n so |
| Bra [v] | | | Aha. |
| sound [v] | | | |
- [4]
- | | | | |
|-----------|--------------------------------|--------------|--|
| | 10 [00:49.4] | 11 [00:51.1] | |
| Ads [v] | Handgriffe macht und nachdenkt | | |
| sound [v] | HAMMERN | | |

Nous recherchons une personne prête à nous aider à créer une feuille de style XSL afin de pouvoir associer les vidéos et les fichiers XML, annotés (entre autres) avec des indications de temps.

Ibtihel Ben Ltaifa

Speed Dating Recherche

Équipe de Linguistique Computationnelle

Ibtihel BEN LTAIFA

Doctor en Informatique

Ibtihel.BEN_Ltaifa@paris-sorbonne.fr

Domaine de recherche

- Fouille de Données
- Réseaux Sociaux
- Intelligence artificielle

Mots clés : Traitement du langage naturel, Extraction de connaissances à partir de données, Indexation sémantique, Annotation sémantique, Classification, Clustering, Recherche d'Information, Détection de communautés, Apprentissage supervisé/non supervisé

Cybernetics and Systems 2020 (Impact Factor = 1.880)

[J1] Ben Ltaifa Ibtihel, Lobna Hlaoua, and Lotfi Ben Romdhane. **Hybrid deep neural network-based text representation model to improve microblog retrieval.** *Cybernetics and Systems*, 51(2):115–139, 2020.

KES 2019

[C2] Ben Ltaifa Ibtihel, Hlaoua Lobna, and Ben Romdhane Lotfi. **A deep learning-based ranking approach for microblog retrieval.** *Procedia Computer Science*, 159:352–362, 2019.

KES 2018

[C1] Ben Ltaifa Ibtihel, Hlaoua Lobna, and Ben Jemaa Maher. **A semantic approach for tweet categorization.** *Procedia Computer Science*, 126:335–344, 2018.

-Approche linguistique pour l'extraction des connaissances dans un texte Arabe en utilisant l'apprentissage automatique profond.

Jean Baptiste Tanguy

Une thèse en humanités numériques

Du bruit dans les données textuelles océrisées : mesure, impact et évaluation. Le cas des mazarinades (1648-1653).

Directeur : Glenn Roe

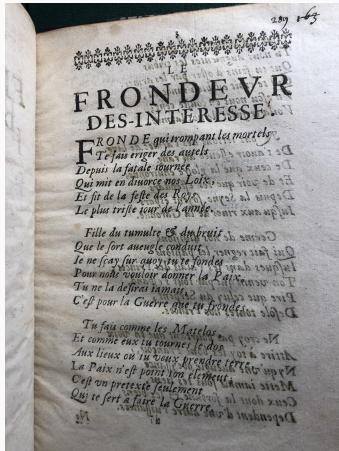
Encadrants : Karine Abiven et Gaël Lejeune

Partenaire : Bibliothèque Mazarine

Financement : Région Île-de-France



Un premier objet d'étude : le "corpus" des mazarinades

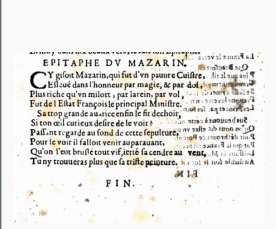


- Pamphlets, poésies burlesques, actes légaux, etc., imprimés entre 1648 et 1653 (La Fronde)
- Ensemble d'environ 6000 pièces (+ existence d'éditions et d'états)
- Français pré-classique (ex. : u=v, i=j, présence du s long, etc.)
- Exemplaires datant du XVIIème siècle, donc : tâches, restaurations, etc.

Environ 500 pièces sélectionnées pour la première vague de numérisation

Acquisition automatique des données textuelles (OCR) sans correction

OCR = reconnaissance optique de caractères \neq HTR

Image	ABBYY (payant)	Kraken (gratuit)
	<p>EPITAPHE DV MAZARIN.</p> <p>CY glist Mazarin, qui fut d'un pauvre Cuffre, di il sus au^h Eileu dans l'honneur par magic, OZ par doi Plus riche qu'un mlort, par larcin, par voi à; 5 7 Fut de l'Etat François le principal Ministre. Sa trop grande auarice enfin le fit déchoir. Si ton œil curieux desire de le voir? -j in she-j émuouûdu? • Paffint regarde au fond de cette leputure, nv' U7-Ln" Pour le voir il falloit venir asparaunt, et, n Qu'on l'eut bruOé tout vif, jette la cendre au vent, d.ii.i indEia Tu ny trouveras plus que la triste peinture, h < ol si nul jdsbuA</p> <p>FIN. MH.</p>	<p>EPITAPHE DV MAZARIN.</p> <p>CY glist Mazarin, qui fut d'un pauvre Cuffre, Eileu dans l'honneur par magic, & par doi, Plus riche qu'un mlort, par larcin, par voi, Fut de l'Etat François le principal Ministre.</p> <p>Sa trop grande auarice enfin le fit déchoir,</p> <p>Si ton œil curieux desire de le voir?</p> <p>Paffant regarde au fond de cette leputure,</p> <p>Pour le voir il falloit venir asparaunt,</p> <p>Qu'on l'eut brué tout vif, jette la cendre au vent, Tu ny trouveras plus que la triste peinture.</p> <p>FIN.</p>

- Observer le bruit → mesurer le bruit
- Comprendre ce sur quoi les logiciels d'OCR buttent
- Erreurs d'OCR et évaluations extrinsèques : dépendance à la tâche, au modèle, aux données ?
- De telles erreurs excluent-elles toute étude en HN ? (ressenti vs. réel)

Lejeune Gael

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
- Humanités Numériques

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
- Humanités Numériques

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
- Humanités Numériques

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 ?

De manière transversale : comment tirer avantage de la variation ?

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) → Thèse à Venir
- ANTONOMAZ (2018-2022) ANalyse auTOMatique et NumérisatiOn des MAZarinades (K.Abiven, G.Roe , JB.Tanguy, A. Bartz)
- OBVIL-NER (2020-2023) Humanités Numériques et Entités Nommées (C. Parfait, M. Alrahabi, G. Roe)
- WADDLE et DANIEL (2018-...) Données Textuelles hétérogènes (A. Barbaresi, E. Giguet) et multilingues (S. Mutuvi, E.Boros, A.Doucet)
- CERES Centre d'expérimentation en Méthodes Numériques pour els SHS (V. Julliard, C. Marti, C.Guillotet, T. Bottini, E. Papinot, F. Allié) + collaboration GEMASS

Manuela Yapomo

Clustering de textes multilingues pour l'extraction de néologismes

Manuela Yapomo

STIH - Sens Texte Informatique Histoire

09 septembre 2021

Contexte

- Création semi-manuelle d'un petit corpus de textes multilingues et multithématiques structuré ;
- Élargissement de ce corpus au moyen du clustering ;
- Extraction de néologismes (candidats) et leurs traductions (candidates) à partir des clusters formés.

Extraits d'articles comparables (biogaz – Wikipédia)

fr

Le **biogaz** est le **gaz** produit par la **fermentation** de **matières organiques animales** ou **végétales** en l'absence d'**oxygène**. Cette **fermentation** appelée aussi **méthanisation** se produit naturellement (dans les marais) ou spontanément dans les **décharges** contenant des **déchets organiques**, [...].

en

Biogas typically refers to a mixture of **gases** produced by the **breakdown** of **organic matter** in the absence of **oxygen**. **Biogas** can be produced from regionally available **raw materials** such as **recycled waste**.

de

Biogas ist ein **brennbares Gas**, das durch **Vergärung** von **Biomasse** jeder Art entsteht. Es wird in **Biogasanlagen** hergestellt, wozu sowohl **Abfälle** als auch **nachwachsende Rohstoffe** vergoren werden.

Motivation scientifique de la recherche

- Vérification de la plus-value du clustering pour l'extraction de néologismes bilingues dans un domaine ;
- Contribution à la science terminologique : documentation des phénomènes néologiques de ce domaine de spécialité ; contraste de langues.

Nicolas Hiebel

Stagiaire M2 et futur doctorant

Nicolas Hiebel

9 septembre 2021

- Master Langue et Informatique en Sorbonne
- Stage au Laboratoire Interdisciplinaire des Sciences du Numérique (LISN)
- Encadrement
 - Aurélie Névéol (LISN)
 - Olivier Ferret (CEA)
 - Karën Fort (LORIA / STIH)
- Thèse prévue à partir d'octobre

Stage M2 : Similarité phrastique dans des corpus médicaux

Problématiques

- Similarité
- Domaine médical

Thèse en octobre : Création éthique de données textuelles artificielles : application au domaine biomédical

Nikola Lakovic

Thèse de Nikola Lackovic

« De la voix à la génération de métadonnées sur la relation client :
Application à l'amélioration de la satisfaction client »

(Thèse CIFRE Malakoff Humanis, Dir. C. Montacié)

- Actes de Dialogue appliqués à la reconnaissance de la parole
- Combiner les informations du texte (retranscription) et de la parole (paralinguistique)
- Objectif 1 : Modèle amélioré de reconnaissance pour le français
- Objectif 2 : Meilleure "collaboration" TAL/TAP



Solveig Poder



Séminaire de rentrée : linguistique computationnelle

Clustering de corpus et topic modeling



Objectifs de la méthode

- Alternative aux outils de textométrie
- Pistes d'analyses pour corpus diachroniques
- Méthode non supervisée et généraliste
- Visualisation claire des résultats



Données

- Deux corpus d'articles de presse (Europresse)
 - Affaire Sackler (crise des opioïdes)
 - PMA
- Un corpus littéraire : 681 Mazarinades
- Formatage des données :
 - Conversion en JSON (scrapping de HTML)
 - Suppression des doublons (Levenshtein)



Méthode

- **Clustering** : propagation d'affinité
- **Topic modeling** : Latent Dirichlet Allocation
- Représentation des données :
 - Plusieurs types de vecteurs (fréquences, TF-IDF)
 - Plusieurs tailles de n-gram (1 à 3 mots)
- Clustering par périodes de 1 à 3 ans
- 1 cluster --> 5 descripteurs (n-grams)

Visualisation (API Flask)

<https://clustertool.herokuapp.com/>

Formulaire :

CLUSTERING DE DONNÉES

Sur quel corpus souhaitez-vous travailler ?

PMA (articles de presse) ▾

Comment souhaitez-vous représenter les données ?

En vecteurs de TF-IDF ▾

Taille minimale des n-grams 1 ▾ Taille maximale des n-grams 2 ▾

Taille de la fenêtre temporelle (en nombre d'année) 1 ▾

Valider

Visualisation (API Flask)

Résultats:

Lien vers motionchart



pma

Lien vers nuage de mots

Lien vers l'article



1994

Cluster 1 "enfant, embryon, don, très, organes" : 2 document(s)

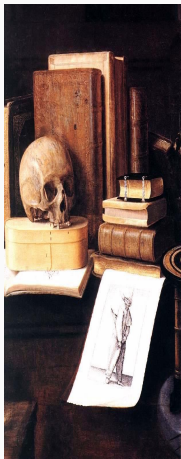
Ethique biomédicale: la France va disposer d'une législation

Libérale Espagne, restrictive Allemagne

Cluster Divers "enfant, embryon, don, très, organes" : 1 document(s)

Des garde-fous pour limiter les dérives et protéger les enfants à naître

Corina Chutaux



Corina CHUTAUX MILA

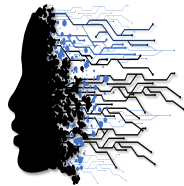
Formation
Double Master

Marché de l'art (ICART) & Littérature Générale et comparée (SORBONNE-NOUVELLE)

Thèse de doctorat (Sorbonne-Nouvelle)

Titre : *Dématérialisation de l'art et de la littérature à l'aube de la digitalisation*

Sujet : Analyse des œuvres littéraires et artistiques générées par des Intelligences artificielles et les conséquences qui en découlent dans ce que j'ai appelé, « le siècle de la dématérialisation »



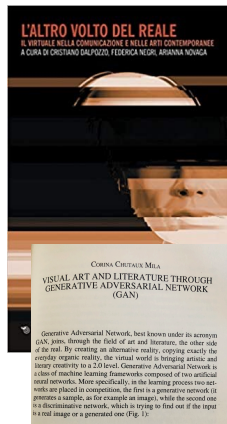
Publications



Ouvrage théorisant pour la première l'art invisible : un art sans œuvre d'art, sans spectateur traditionnel et parfois sans auteur, À paraître en septembre

Ouvrage collectif, en collaboration avec l'Université de Vérone, paru en 2020

Chapitre intitulé *Visual Art and literature trough Generative Adversarial Networks,*



CORINA CHUTAUX MILA
VISUAL ART AND LITERATURE THROUGH
GENERATIVE ADVERSARIAL NETWORK
(GAN)

Generative Adversarial Network, best known under its acronym GAN, joins, through the field of art and literature, the other side of the real. By creating an alternative reality, copying exactly the everyday organic reality, the virtual world is bringing artistic and literary creativity to a 2.0 level. Generative Adversarial Network is a class of machine learning frameworks composed of two artificial neural networks. More specifically, in the learning process two networks are placed in competition, the first is a generative network (it generates a sample, as for example an image), while the second one is a discriminative network, which is trying to find out if the input is a real image or a generated one (Fig. 1).



Expériences professionnelles
antérieures dans l'enseignement

ENDA

Cours :

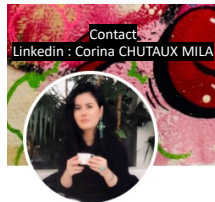
- Histoire de l'art
- Art Invisuel vs Art transhumaniste



UNIVERSITÉ
**SORBONNE
NOUVELLE**
PARIS 3

Cours :

- Méthodologie de la Recherche documentaire
- Communication digitale



SciencesPo.

Cours :

- Art, Littérature et Politique du 18e au 21e siècle
- Enseignante-Référente

Heesoo Choi

Séminaire STIH - Linguistique computationnelle

Hee-Soo Choi

9 septembre 2021

Qui je suis



- Précédemment en Licence Langue Française et Techniques Informatiques (LFTI)
- Actuellement en Master 2 Langue et Informatique et en stage de recherche au LORIA sous la direction de Karèn Fort et Bruno Guillaume
- Prochainement en Thèse...

Vérification d'universaux linguistiques sur des corpus multilingues annotés en syntaxe de dépendances

- 4 universaux de Greenberg testés
- 141 corpus, 74 langues d'Universal Dependencies
- *Investigating Dominant Word Order on Universal Dependencies With Graph Rewriting*, Hee-Soo Choi, Bruno Guillaume, Karën Fort, Guy Perrier (RANLP2021)



Lier des ressources lexicales du français en vue d'une interopérabilité entre niveaux linguistiques

- Encadré par Mathieu Constant, Karën Fort et Bruno Guillaume
- Laboratoires : ATILF et LORIA
- Financement : École doctorale SLTC, Université de Lorraine



Lier des ressources lexicales du français en vue d'une interopérabilité entre niveaux linguistiques

- Lier automatiquement des ressources hétérogènes mais complémentaires
- Développer des méthodes de liage cohérentes
- Utiliser une représentation par graphes
- Obtenir un ensemble lié de ressources pour les outils de TAL

Karen Fort



Création de ressources langagières et éthique pour le TAL

Karèn Fort

karen.fort@sorbonne-universite.fr / <https://members.loria.fr/KFort/>

Séminaire d'équipe, Speed dating, 9 septembre 2021



Production participative (*crowdsourcing*)



ZOMBILINGO

RIGORMORTIS

BISAME

KRIK!

AYO!



Portails de science participative et de jeux pour les langues :

SCIENCE ENSEMBLE

L I N G O B O I N G O

Atelier récurrent : Games4NLP

Éthique et TAL

- ▶ Création de données pour un TAL plus éthique :
 - ▶ corpus du français pour l'évaluation des biais stéréotypés des modèles de DL (A. Névéol, Y. Dupont et J. Besançon) : <https://languagearc.com/projects/19> (Participez !)
 - ▶ corpus de dossiers médicaux synthétiques en français (A. Névéol, O. Ferret, N. Hiebel)
- ▶ Analyse des sections éthiques de la conférence NAACL 2021 (E. Bender, M. Mitchell et E. van Miltenburg)
- ▶ Grille d'analyse déontologique pour le TAL (M. Amblard)

Projets en cours

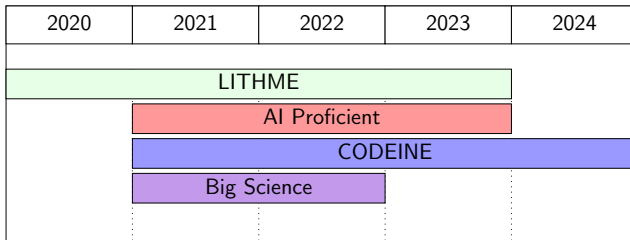


FIGURE – Projets en cours (en vert les actions COST, en rouge les projets européens, en bleu les projets ANR, en violet les projets avec une entreprise)

Encadrements de thèses

- ▶ **Lucas Lambrey** (avec Gilles Siouffi et Antoine Gautier) :
Détection et production de défigements linguistiques dans les réseaux sociaux assistés par les sciences participatives : fertilisation croisée entre traitement informatique et analyse linguistique. (CERES)
- ▶ **Heesoo Choi** (avec Mathieu Constant et Bruno Guillaume) :
Lier des ressources lexicales du français en vue d'une interopérabilité entre niveaux linguistiques. (Univ. de Lorraine).
- ▶ **Nicolas Hiebel** (avec Aurélie Névéol et Olivier Ferret) :
Création éthique de données textuelles artificielles : application au domaine biomédical. (CODEINE).

Responsabilités

- ▶ Chair du comité d'éthique d'ACL (Association for Computational Linguistics)
- ▶ Comité d'éthique de la recherche (CER) de SU
- ▶ Conseil national des universités (CNU) 27 (informatique) : <https://cnu27.univ-lille.fr/>
 - ▶ Qualifications
 - ▶ Promotions, CRCT, etc
- ▶ GDR LIFT (linguistique informatique, formelle et de terrain)
 - ▶ École d'été sur l'annotation en 2022 (à confirmer)
- ▶ Actions européennes COST :
 - ▶ Language In The Human-Machine Era (LITHME)