

## LES PRINCIPAUX DOMAINES DE L'ANALYSE LINGUISTIQUE ET QUELQUES EXEMPLES D'APPROCHES D'ANALYSE AUTOMATIQUE DU LANGAGE

Dans ses différentes dimensions, orales et écrites

Domaine	unité	Problématiques	Outils / bibliothèques / ressources
Phonétique	Phone	<p>Extraire corrélats du signal (DURÉE, FRÉQUENCE, formants)</p> <p>Outils : Aligneurs texte/son</p> <p>Annotation (token : alignement, phonétisation, syllabe, Prosodie)</p>	<p>Praat : analyse du spectre <a href="http://praat.org">praat.org</a></p> <ul style="list-style-type: none"> <li>• align data</li> <li>• extract data</li> <li>• analyse data</li> <li>• train model</li> </ul> <p>EXAMPLE Adrien Méli's PhD <a href="http://mapage.noos.fr/admeli/soutenance/#/">http://mapage.noos.fr/admeli/soutenance/#/</a></p> <p>Se pilote avec python via la bibliothèque Parselmouth  <a href="https://parselmouth.readthedocs.io/en/stable/">https://parselmouth.readthedocs.io/en/stable/</a>  <a href="https://github.com/YannickJadoul/Parselmouth">https://github.com/YannickJadoul/Parselmouth</a></p> <p><a href="http://sppas.org">SPPAS</a> <a href="http://sppas.org">sppas.org</a> (représentation en SAMPA)</p> <p>(se pilote aussi en python)</p> <ul style="list-style-type: none"> <li>➔ aligneur, phonétiseur</li> <li>➔ analyse des syllabes</li> <li>➔ analyse de la prosodie (MOMEL-INTSINT)</li> </ul> <p>###ALIGNERS  <a href="https://eleanorchodroff.com/tutorial/index.html">https://eleanorchodroff.com/tutorial/index.html</a></p>

		Reconnaissance vocale (VR) TTS text to speech resynthèse vocale	<a href="#">forced-alignment-tools</a> , <a href="#">A collection of links and notes on forced alignment tools by Alberto Pettarin</a>  P2FA (aligneur de Penn)  PyTorch-NLP <ul style="list-style-type: none"> <li>• Audio I/O and Pre-Processing with torchaudio</li> <li>• Speech Command Recognition with torchaudio</li> </ul> #CMU Resources (part of the NLTK toolkit) Dedicated python packages: Pronouncing  MISC <a href="#">Phonetic conversion (and text to speech)</a>  Tacotron
Phonologie	Phonème	faire l'inventaire des systèmes phonologiques, des syllabes possibles, étudier les correspondances (sons/lettres :	PHON (Y. Rose) Algorithmes de découpages en syllabes <a href="https://www.phon.ca/phon-manual/getting_started.html">https://www.phon.ca/phon-manual/getting_started.html</a>  Algorithmes de découpages en syllabes dans SPPAS (sppas.org) API en SAMPA  API en ARBANET Ressources : le CMU <a href="http://www.speech.cs.cmu.edu/cgi-bin/cmudict">http://www.speech.cs.cmu.edu/cgi-bin/cmudict</a>

		grapho-phonématique) (OCR)	<p>(disponible dans la bibliothèque NLTK)</p> <p># le CELEX (décomposition en morphèmes, transcription, fréquence, ressource payante pour l'anglais, se méfier du découpage en syllabes)</p> <p>#French: A lexical database for contemporary French : LEXIQUE™  <a href="#">A lexical database</a>  <a href="#">Pallier et al. 2001</a></p>
Morphologie	morphème	<p>Détecter les morphèmes  <i>Singer</i> / <i>stronger</i>  Morphème dérivationnel  Morphème flexionnel</p>	<p>Pour l'anglais</p> <p>Linguistica</p> <p>Pour des analyses non-supervisées :  Chipmunk</p> <p>{-er} [sing]<sub>V</sub> / [singer]<sub>N</sub> / 2 types (25 lemmas)  {strong, stronger} [strong]<sub>A</sub> 2 tokens / 1 lemma</p> <p>Pour le français:  LEFF: Morphological and syntactic lexicon for French  <a href="http://pauillac.inria.fr/~sagot/index.html#wolf">http://pauillac.inria.fr/~sagot/index.html#wolf</a></p> <p>LEXIQUE (et fréquence)  <a href="http://www.lexique.org/">http://www.lexique.org/</a>  <a href="http://www.lexique.org/shiny/openlexicon/">http://www.lexique.org/shiny/openlexicon/</a></p>
Lexicologie	Lexème (forme de mot)	Retrouver les lemmes Tokenisation	<p>BNC (100M) CLAWS BNC2014  COCA (350M)  BLP SUBTLEX US</p>

		<p>Analyse des stopword</p> <p>Identifier les collocations</p>	<p>Ressources pour la fréquence  <a href="http://crr.ugent.be/papers/SUBTLEX-US%20frequency%20list%20with%20PoS%20and%20Zipf%20information.zip">http://crr.ugent.be/papers/SUBTLEX-US%20frequency%20list%20with%20PoS%20and%20Zipf%20information.zip</a>  l'échelle de Zipf :  <a href="http://crr.ugent.be/archives/1352">http://crr.ugent.be/archives/1352</a></p> <p>Tokeniser (tokenisation) / lemmatisation (lemma/type) TTR</p> <p>NLTK / wordnet</p> <p>Treetagger</p> <p>PARSEME (fr)  Multi-word expression / MWE / MWU</p>
Syntaxe	Syntagme Propositions syntaxiques	<p>Pos-tagging</p> <p>PARSING</p> <p>Universal Dependencies</p>	<p>Treetagger / NLTK  coreNLP</p> <p>Spacy</p> <p>Stanford corenlp (JAVA)</p> <p>Démo en ligne : <a href="https://corenlp.run/">https://corenlp.run/</a></p> <p>UDpipe  <a href="https://github.com/ufal/udpipe">https://github.com/ufal/udpipe</a>  UD models</p>

			<p><a href="https://universaldependencies.org/#language-u">https://universaldependencies.org/#language-u</a> (Sequoia) <a href="https://github.com/UniversalDependencies/UD_French-Sequoia">https://github.com/UniversalDependencies/UD_French-Sequoia</a> (Partut)</p> <p>Plus fin : <a href="https://surfacesyntacticud.github.io/conversions/">https://surfacesyntacticud.github.io/conversions/</a> SUD Surface Syntactic Universal Dependency <a href="https://surfacesyntacticud.github.io/">https://surfacesyntacticud.github.io/</a></p> <p>analyse en coréférence</p>
Sémantique	Sème/ sémantème	<p>WORDNET (anglais) WOLF (français)</p> <p>détection des synsets</p> <p>relations</p> <p>Détecter les noms propres, les références au lieu et au temps NER</p>	<p><a href="https://wordnet.princeton.edu/">https://wordnet.princeton.edu/</a> <a href="http://wordnetweb.princeton.edu/perl/webwn">http://wordnetweb.princeton.edu/perl/webwn</a></p> <p>le « wordnet français », INRIA, équipe de Benoit Sagot : <a href="http://blog.onyme.com/etude-de-lontologie-wordnet-libre-du-francais-wolf/">http://blog.onyme.com/etude-de-lontologie-wordnet-libre-du-francais-wolf/</a> <a href="http://alpage.inria.fr/~sagot/wolf.html">http://alpage.inria.fr/~sagot/wolf.html</a> représentations ontologiques : maison <a href="http://blog.onyme.com/wp-content/uploads/2013/09/wolf-0.1.5-jouet.png">http://blog.onyme.com/wp-content/uploads/2013/09/wolf-0.1.5-jouet.png</a></p> <p>Encore des difficultés : <a href="http://blog.onyme.com/wp-content/uploads/2013/09/wolf-1.0b-maison-wrong.png">http://blog.onyme.com/wp-content/uploads/2013/09/wolf-1.0b-maison-wrong.png</a></p> <p>NER ( Named Entity Recognition)    NLTK, Spacy</p>

		word embedding..	Hugging-face  (BERT, ALBERT, ROBERT .... FlaubERT) CamemBERT,BioBert  word similarities:cosine, npmi...  Glove,Word2vec
Pragmatique	Acte de langage (Austin , Grice)	Que veut-il dire quand il utilise telle forme/telle phrase ? pertinence Calcul d'implicature Maximes de la conversation Principe de coopération	L'analyse du sentiment (sentana : sentiment analysis) <a href="https://corenlp.run/">https://corenlp.run/</a>  Chatbot ...

**Pour aller plus loin:**

Emily M. Bender (2013) *Linguistic Fundamentals for Natural Language Processing: 100 Essentials from Morphology and Syntax*, Morgan & Claypool

Emily M. Bender & Alex Lascarides (2020) *Linguistic Fundamentals for Natural Language Processing II. 100 Essentials from Semantics and Pragmatics*, Morgan & Claypool

## **OUTILS POUR L'ANALYSE AUTOMATIQUE DE LA COMPLEXITE SYNTAXIQUE ET LEXICALE**

### **La page de Xiaofei Lu**

<http://www.personal.psu.edu/xxl13/download.html>

(voir aussi son analyseur [D-Level Analyzer](#) )

Pour la complexité lexicale

[Lexical Complexity Analyzer](#)

L2SCA : L2 Syntactic Complexity Analyzer

<http://www.personal.psu.edu/xxl13/downloads/l2sca.html>

### **La famille d'outils de Kyle & Crossley:**

Les outils d'analyse de la sophistication lexicale (TAALES), de la complexité (TAASC), de la cohésion (TAACO) mais aussi pour l'analyse automatique du sentiment (SEANCE)

L'interface est très intuitive et multi-OS (pour des raisons de d

pour l'analyse automatique du sentiment : SEANCE

<https://www.linguisticanalysistools.org/seance.html>

article de synthèse : <https://hal.archives-ouvertes.fr/hal-02768504v3/document>

## DES RESSOURCES LINGUISTIQUES A L'ASSISTANCE A LA REDACTION EN ANGLAIS (*COMPUTER-AIDED WRITING*)

Votre moteur de recherche préféré peut être utile (utilisez les guillemets), mais vous aurez aussi des surprises

Pour vérifier que le mot, la collocation existe et pour vérifier sa construction : les corpus de références comme le COCA (350M mots, américain) et le BNC (100M mots, britannique)

<https://www.english-corpora.org/>

<https://www.english-corpora.org/coca/>

<https://www.english-corpora.org/bnc>

A partir de votre texte rédigé en français, utiliser un moteur de traduction, <https://www.deepl.com/fr/translator> puis vérifier vos doutes éventuels avec votre moteur de recherche préféré

### **!! Attention au plagiat par traduction**

Ne soumettez pas de travaux où vous auriez juste copié- collé des phrases des articles (ou des blogs scientifiques). Il s'agirait de plagiat, ce qui est répréhensible

Pour bien comprendre le degré de reformulation autorisé qui distingue le plagiat d'un véritable travail de rédaction scientifique, étudiez les exemples proposés sur les sites des universités britanniques et américaines.

Voici par exemple le site d'Oxford.

<https://www.ox.ac.uk/students/academic/guidance/skills/plagiarism>

pour la traduction (corpus alignés)

<https://www.linguee.fr/anglais-francais/traduction/go+online.html>

(C'est la qualité initiale de ces ressources qui ont assuré le succès foudroyant de DeepL)

La détection automatique d'erreurs grammaticales et orthographiques (songez à votre traitement de texte)

Grammarly : <https://www.grammarly.com/>

La suggestion collocations plus fréquentes (COLLOCAID)

Pour les collocations : <http://www.collocaid.uk/>



