



**HAL**  
open science

## **EVOLEX : approches psycholinguistique et computationnelle de l'accès au lexique et de la proximité sémantique entre paires de mots**

Xavier de Boissezon, Lola Danet, Cécile Fabre, Jérôme Farinas, Bruno Gaume, Nabil Hathout, Lydia-Mai Ho-Dac, Mélanie Jucla, Patrice Péran, Bénédicte Pierrejean, et al.

### ► To cite this version:

Xavier de Boissezon, Lola Danet, Cécile Fabre, Jérôme Farinas, Bruno Gaume, et al.. EVOLEX : approches psycholinguistique et computationnelle de l'accès au lexique et de la proximité sémantique entre paires de mots. Forum à la croisée des sciences : Interagissez, Imaginez, Innovez - FACS3I, Jan 2019, Toulouse, France. hal-02047651

**HAL Id: hal-02047651**

**<https://hal.science/hal-02047651v1>**

Submitted on 25 Feb 2019

**HAL** is a multi-disciplinary open access archive for the deposit and dissemination of scientific research documents, whether they are published or not. The documents may come from teaching and research institutions in France or abroad, or from public or private research centers.

L'archive ouverte pluridisciplinaire **HAL**, est destinée au dépôt et à la diffusion de documents scientifiques de niveau recherche, publiés ou non, émanant des établissements d'enseignement et de recherche français ou étrangers, des laboratoires publics ou privés.



# EVOLEX : APPROCHES PSYCHOLINGUISTIQUE ET COMPUTATIONNELLE DE L'ACCÈS AU LEXIQUE ET DE LA PROXIMITÉ SÉMANTIQUE ENTRE PAIRES DE MOTS

Xavier de Boissezon<sup>1</sup> Lola Danet<sup>1</sup> Cécile Fabre<sup>2</sup> Jérôme Farinas<sup>3</sup> Bruno Gaume<sup>2</sup> Nabil Hathout<sup>2</sup> Lydia-Mai Ho-Dac<sup>2</sup>  
Mélanie Jucla<sup>4</sup> Patrice Péran<sup>1</sup> Bénédicte Pierrejean<sup>2</sup> Julien Pinquier<sup>3</sup> Ludovic Tanguy<sup>2</sup>

<sup>1</sup>ToNIC, Inserm, UT3 <sup>2</sup>CLLE-ERSS, CNRS, UT2J <sup>3</sup>IRIT, CNRS, UT3 <sup>4</sup>Octogone-Lordat, UT2J Contact: xavier.deboissezon@inserm.fr, melanie.jucla@univ-tlse2.fr, ludovic.tanguy@univ-tlse2.fr, evolex@irit.fr

## OBJECTIFS

Développer une **méthode assistée par ordinateur** pour évaluer auprès de populations avec ou sans troubles du langage la façon dont

- ▶ nous accédons aux informations présentes dans notre lexique mental
- ▶ nous associons ces informations les unes avec les autres

Pour réaliser cet objectif, nous allions des compétences en psycholinguistique et neuropsychologie d'une part et de traitement automatique des langues naturelles d'autre part

## ORGANISATION DU PROJET EVOLEX

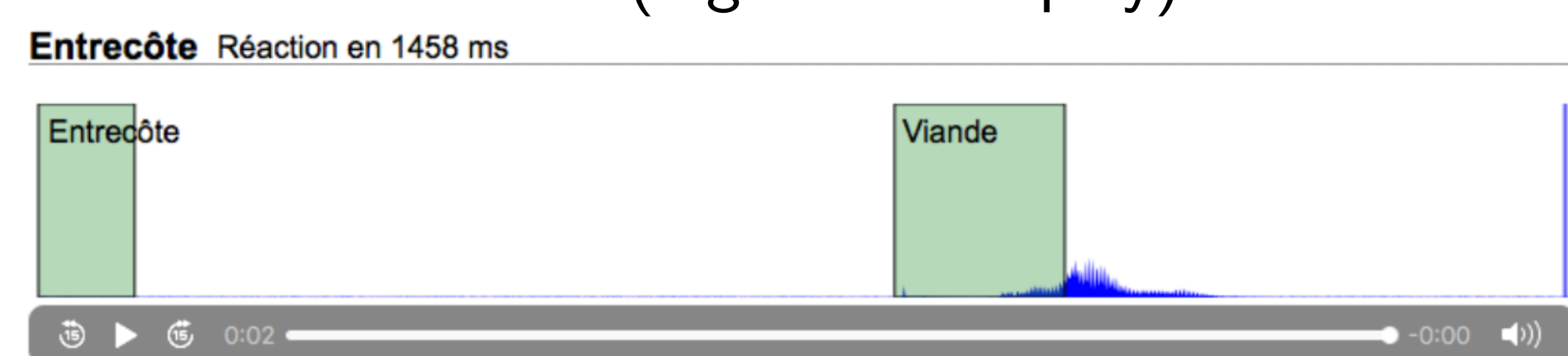
1. Développement d'outils pour la reconnaissance automatique de la parole et la correction manuelle des données récoltées (productions orales et temps de réponse) IRIT, ToNIC
2. Définition des épreuves, du matériel linguistique utilisé, collecte et prétraitement des données Octogone-Lordat, ToNIC
3. Annotation et analyses qualitatives des réponses CLLE-ERSS, Octogone-Lordat
4. Application des modèles de Sémantique Distributionnelle pour évaluer les données CLLE-ERSS
5. Mise en commun des analyses pour la caractérisation de profils de locuteurs et des réseaux sémantiques CLLE-ERSS, Octogone-Lordat, ToNIC

## PROTOCOLE EVOLEX

**3 épreuves d'accès lexical** dont une épreuve de **Génération Verbale (GV)** :

1 nom entendu ⇒ le premier mot qui vient à l'esprit (ex : *abricot* → *confiture*)

- ▶ Recueil des données (logiciel Samoplay)



- ▶ Passations auprès de locuteurs avec ou sans troubles du langage de 18 à 90 ans



Ce projet a bénéficié du soutien de l'Institut des Handicaps du CHU de

Toulouse et de la Fédération Hospitalo-Universitaire des Handicaps Cognitifs, Psychiques et Sensoriels

## RÉFÉRENCES

- [1] B. Gaume, K. Duvignau, E. Navarro, Y. Desalle, H. Cheung, S.-K. Hsieh, P. Magistry, and L. Prevot. Skilllex : a graph-based lexical score for measuring the semantic efficiency of used verbs by human subjects describing actions. *Revue TAL*, 55(3), 2016.
- [2] B. Gaume, L. Tanguy, C. Fabre, L.-M. Ho-Dac, B. Pierrejean, N. Hathout, J. Farinas, J. Pinquier, L. Danet, P. Péran, X. De Boissezon, and M. Jucla. Automatic analysis of word association data from the evolex psycholinguistic tasks using computational lexical semantic similarity measures. In 13th International Workshop on Natural Language Processing and Cognitive Science (NLPCS), Krakow, Poland, 2018.
- [3] P. Péran, J.-F. Démonet, C. Pernet, and D. Cardebat. Verb and noun generation tasks in huntington's disease. *Movement disorders : official journal of the Movement Disorder Society*, 19(5) :565-571, 2004.
- [4] L. Tanguy, F. Sajous, and N. Hathout. Évaluation sur mesure de modèles distributionnels sur un corpus spécialisé : comparaison des approches par contextes syntaxiques et par fenêtres graphiques. *Traitement Automatique des Langues*, 56(2), 2015.

## PARTICIPATION DE QUATRE LABORATOIRES TOULOUSAINS



**CLLE-ERSS** est une équipe de linguistique comportant un axe de recherche en traitement automatique des langues



**SAMoVA** de l'IRIT pour « **Structuration, Analyse et Modélisation de documents Vidéo et Audio** »



**Octogone-Lordat**, laboratoire de neuro-psycholinguistique, mène des recherches sur l'étude des (dys)fonctionnements langagiers et de la cognition



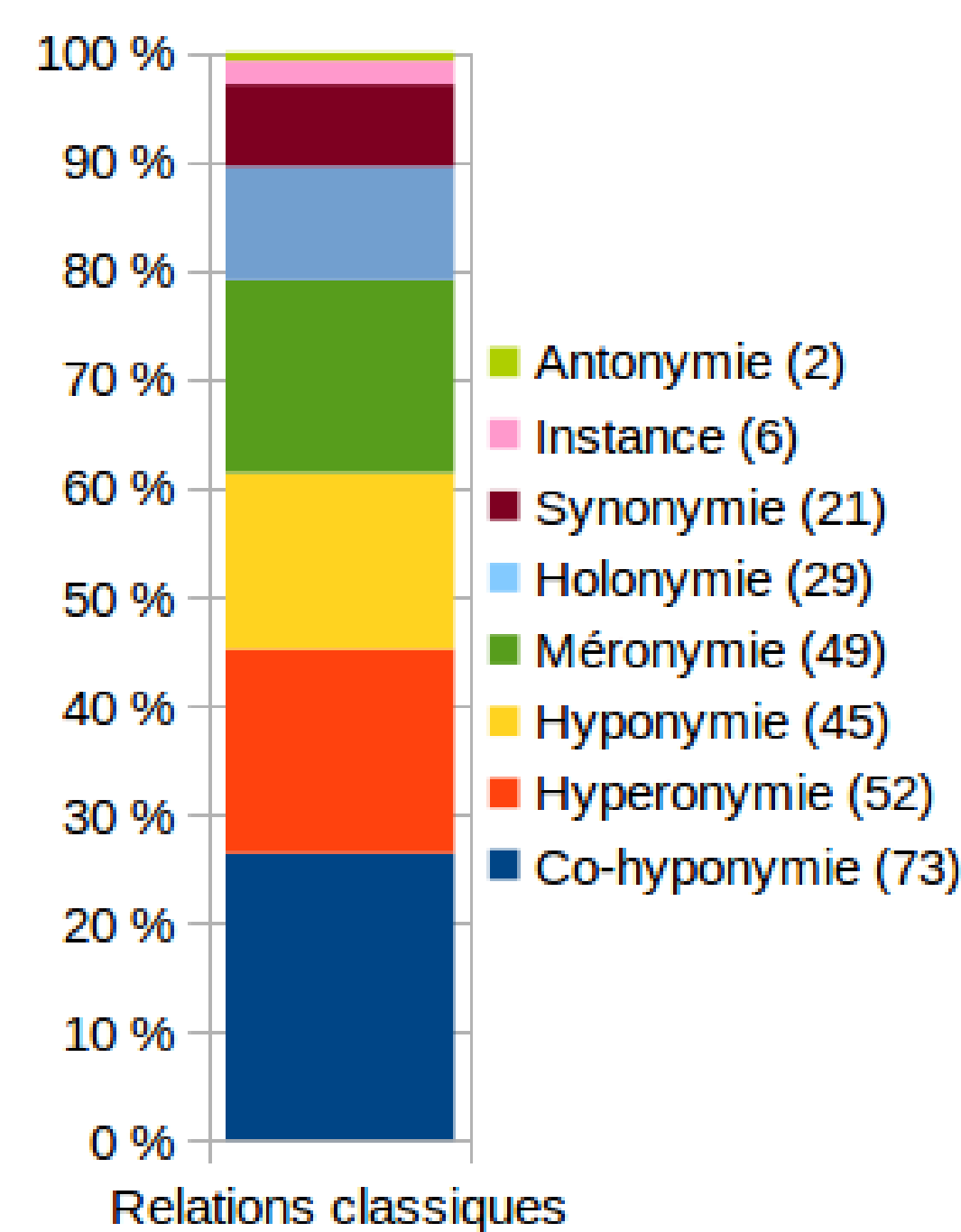
**ToNIC** pour « Toulouse Neuro Imaging Center », a pour objectif principal l'étude du cerveau humain et des principales pathologies qui l'affectent

## DONNÉES (GV, VERSION 1, 30 SUJETS)

- ▶ 1800 paires stimulus-réponse récoltées
- ▶ 1544 paires exploitables (559 différentes)

### Double-annotation des relations

⇒ des **relations "classiques"** (50%) :



⇒ des **relations "autres"** : association d'idées (36%), syntagmes (9%), pas de relation (5%), lien phonologique (1%)

## MÉTHODE

**4 mesures pour évaluer la similarité / proximité sémantique entre mots**

**2 similarités issue de corpus**

- ▶ **Corpus1st** : Plus deux mots apparaissent fréquemment ensemble en corpus, plus ces deux mots sont similaires (= collocations)

- ▶ **Corpus2nd** : Plus deux mots partagent les mêmes collocations, plus ces deux mots sont similaires (= substituabilité)

- ▶ corpus = FrWac (pages Web .fr, 2 milliards de mots)

**2 similarités issues de dictionnaires**

- ▶ **Dict1st** : Si un mot apparaît dans la définition d'un autre mot, ces deux mots sont similaires

- ▶ **Dict2nd** : Plus les définitions de deux mots partagent de mots, plus ces deux mots sont similaires

- ▶ dict. = TLF – Trésor de la Langue Française

## RÉSULTATS PRINCIPAUX

### Analyse en Composantes Principales

→ Des mesures positivement corrélées aux **relations classiques**

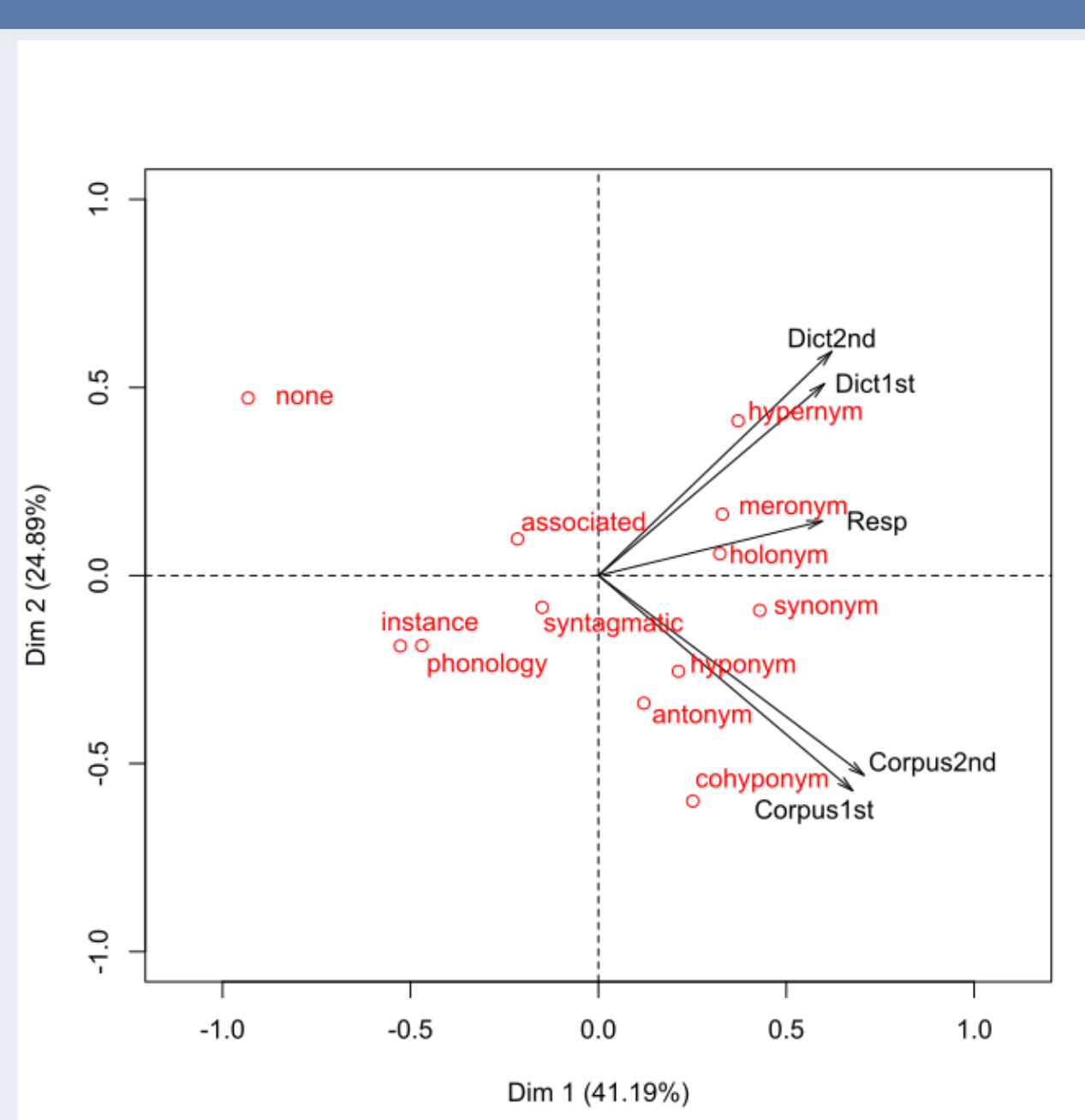
→ Des mesures négativement corrélées dans les cas de

- ▶ sens éloigné (sans relation identifiée, phonologie : *chou – caillou*)
- ▶ d'instance (*magicien – Merlin*)

→ Des aspects différents captés par les différentes méthodes

- ▶ **dictionnaire** – hyperonymie (*entrecôte – viande*)
- ▶ **corpus** – cohyponymie (*balançoire – toboggan*)

→ Association et Syntagmatique non corrélées et caractérisées



## AU DELÀ DE LA CARACTÉRISATION DES PAIRES DE MOTS...

Identification de profils de réponses

