



**HAL**  
open science

# Construction d'une base de données lexicale pour les mots français abstraits et concrets

Daria Goriachun

► **To cite this version:**

Daria Goriachun. Construction d'une base de données lexicale pour les mots français abstraits et concrets. Sciences de l'Homme et Société. 2020. dumas-03024192

**HAL Id: dumas-03024192**

**<https://dumas.ccsd.cnrs.fr/dumas-03024192>**

Submitted on 18 Jan 2021

**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.

Aix-Marseille Université  
Faculté des Arts, Lettres, Langues, Sciences Humaines



Mémoire Master 2 Sciences du Langage

**Construction d'une base de données lexicale  
pour les mots français abstraits et concrets  
Contribution au projet ANR ALECTOR**

Daria GORIACHUN  
Sous la direction de Núria GALA

Mai 2020

## Remerciements

Je tiens à exprimer ma reconnaissance à Nuria Gala pour avoir accepté de m'encadrer dans cette étude. Je la remercie pour soutien et ses encouragements tout au long de ce travail et ses conseils méthodologiques indispensables.

Je remercie l'Agence Nationale pour la Recherche pour avoir financé ce travail, à travers le projet ALECTOR (ANR-16-CE28-0005).

Je remercie profondément M. Frank Sajous d'avoir fourni l'accès aux données complètes de la base de données *Les Voisins De Le Monde*, ainsi que Christelle Zielinski ingénieure d'étude de l'ILCB pour nous avoir aidé avec la mise en place de l'expérience en ligne.

## Résumé

Les expériences en linguistique expérimentale ont montré que les mots abstraits et concrets sont perçus différemment par notre système cognitif. En tenant compte de ce constat, nous faisons l'hypothèse que les mots peuvent avoir les traits inhérents à une catégorie particulière : abstraite ou concrète. Cependant, les lexiques disponibles pour le français ne contiennent pas cette information, d'autant plus qu'elle est difficile à cerner. Au cours de notre travail nous avons essayé d'identifier ces traits particuliers en partant des variables disponibles dans les ressources existantes, comme la fréquence et la structure morphologique. Le seul trait fiable qui permet de distinguer les mots concrets des mots abstraits, cependant, était l'iconicité. Pour pouvoir cibler ce trait nous avons ensuite développé une base de données lexicale de façon semi-automatique en partant d'une liste de 369 noms abstraits et concrets et en utilisant deux méthodes distributionnelles. Notre objectif était d'observer le comportement des voisins distributionnels et des cooccurrences syntaxiques dans les deux classes des noms et de construire une liste par propagation à partir des listes initiales. Les résultats obtenus, 7.898 noms annotés automatiquement, ont permis de valider notre approche, ce qui nous permettra de poursuivre la recherche avec l'utilisation dans d'autres travaux expérimentaux. Notre objectif à moyen terme étant d'explorer l'impact du trait 'iconicité' dans la compréhension et la lecture de mots en contexte,

## Table des matières

1. Introduction .....	6
2. Caractérisation des mots abstraits et concrets .....	8
2.1. Aspects neuropsycholinguistiques .....	8
2.2. Sur les notions d'abstrait et de concret .....	10
2.3. Méthodes d'annotation de mots abstraits et concrets .....	12
3. Expérience 1. Impact de la fréquence sur l'abstractivité du lexique.....	15
3.1. Stimuli .....	15
3.2. Participants et tâche .....	16
3.3. Résultats et discussion .....	16
3.4. Biais d'étude .....	17
4. Expérience 2. Impact de la morphologie sur l'abstractivité du lexique.....	18
4.1. Stimuli.....	18
4.2. Questionnaire et tâche.....	19
4.3. Participants.....	20
4.4. Résultats.....	20
4.5. Discussion.....	22
5. Expérience 3. Annotation automatique du lexique en traits abstrait et concret.....	23
5.1. Données .....	24
5.2. Méthodologie .....	25
5.3. Résultats.....	26
5.4. Évaluation .....	28
5.4.1. Stimuli et procédure .....	28
5.4.2. Participants .....	30
5.4.3. Analyse statistique et résultats .....	30
A. Analyse de la précision des méthodes distributionnels .....	30
B. Analyse de l'impact de temps de réaction .....	32
C. Analyse de l'impact de la langue maternelle et des troubles du langage.....	32
D. Analyse d'effet de l'ambiguïté abstrait/concret.....	33
E. Analyse de l'impact de la fréquence des stimuli .....	33
5.4.4. Discussion .....	33
6. Conclusion.....	34
Bibliographie .....	37

Annexes .....	44
Annexe 1. Questionnaire avec des stimuli. Expérience 1.....	44
Annexe 2. Données des participants. Expérience 1.....	47
Annexe 4. Le rapport des mots fréquents et rares à l'accord entre annotateurs. Expérience 1 .....	49
Annexe 5. Le rapport des mots concrets et abstraits à l'accord entre annotateurs. Expérience 1 .....	50
Annexe 6. Stimuli. Expérience 2 .....	51
Annexe 7. Questionnaire. Expérience 2 .....	52
Annexe 8. Graphiques : Mots concrets avec des suffixes et mots concrets sans suffixes. Expérience 2 .....	53
Annexe 9. Graphiques : Mots abstraits avec des suffixes et mots abstraits sans suffixes. Expérience 2 .....	54
Annexe 10. Graphique : Mots très concrets. Expérience 2 .....	55
Annexe 11. Les moyennes et les écarts types. Expérience 2.....	56
Annexe 12. Listes initiales.....	58
Annexe 13. Exemples des données filtrées. ....	59
Annexe 14. Stimuli pour l'expérience en ligne. ....	62

## 1. Introduction

L'utilisation de variables cognitives pour augmenter la précision des algorithmes dans certaines tâches de TALN est important en linguistique informatique et en neurosciences. Des études récentes incluent, par exemple, l'annotation automatique de la polarité des sentiments au niveau des textes (Mohammad, 2016), des phrases (Wilson et al., 2009) et des mots (Gala & Brun, 2012). Pourtant si certains éléments peuvent être obtenus à partir d'informations contextuelles, il n'y a pas de réel consensus sur ce qu'est le niveau d'iconicité des mots et comment peut-il être obtenu automatiquement.

La différence entre les deux catégories, abstrait et concret, a été montrée dans des expériences psycholinguistiques il y a longtemps. La base de ces études est la théorie à double codage, décrite par Paivio (1965, 1991) qui décrit deux sous-systèmes cognitifs distincts : deux façons, verbale et non verbale, de décoder l'information. Leur activation dépend du degré d'iconicité du mot. Si les mots concrets utilisent d'une manière égale ces deux systèmes parce qu'ils ont une image comme support dans la mémoire du locuteur (notion d'iconicité), les mots abstraits peuvent uniquement être décodés verbalement. Plus tard, cette théorie a également été prouvée par des tests du potentiel évoqué (PE) et par l'imagerie par résonance magnétique fonctionnelle (IRMf) (Just et al., 2004) qui ont montré la distinction détaillée entre l'activation de différentes zones cérébrales lors du traitement de mots abstraits et concrets.

Le terme « effet de concrétude » (*concreteness effect*) fait référence aux temps de réaction plus rapides pour les mots concrets dans différents types de tâches cognitives (Jessen et al., 2000). Un certain nombre de théories expliquant l'effet de concrétude sur les lecteurs normaux et les personnes ayant des troubles de lecture ont été proposées. Plaut et Shallice (1993), dans leur modèle connexionniste, expliquent que les mots concrets sont lus plus facilement, grâce à la facilité de leur caractérisation (iconicité). Une étude récente a montré un impact de l'iconicité et de la régularité orthographique des mots sur la précision de lecture des mots et sur l'efficacité d'apprentissage des mots varié sur la régularité orthographique et l'imageabilité (Steady & Compton, 2019).

En s'appuyant sur ces théories, on peut supposer que les mots concrets sont lus plus facilement que les mots abstraits, car ils bénéficient de deux façons de décodage. L'hypothèse a été confirmée par Kroll & Merves (1986) et James (1975) dans une étude où ils se réfèrent à la facilité avec laquelle un mot évoque une image mentale. Également, Shallice (1988) et

Schwanenflugel (1991) affirment que les mots hautement imaginables ont une représentation sémantique plus riche ou plus facilement accessible.

De nombreuses études ont été menées pour prouver que des facteurs tels que la longueur des mots, ainsi que certaines caractéristiques morphologiques et phonétiques affectaient la perception et la lecture des mots, mais peu de recherches ont été menées pour prouver la complexité des mots abstraits par rapport aux mots concrets.

Bien qu'il y ait des recherches qui ont montré que les mots abstraits présentent des difficultés de lecture pour les personnes dyslexiques en raison de leur faible caractère iconique, le problème est double à l'heure actuelle :

- 1) Comment identifier les mots complexes en contexte, en tenant compte de leur trait abstrait/concret, car cette notion n'est pas modélisée ;
- 2) Il n'y a pas de ressource lexicale avec des informations explicites, exhaustives et accessibles d'où on peut extraire cette information pour modaliser la notion.

Ainsi dans le but d'enrichir un système d'identification de mots complexes à lire et à comprendre, nous souhaitons constituer une ressource qui renseigne le niveau d'iconicité des noms en français. Dans ce travail, nous visons à identifier automatiquement les noms abstraits et concrets en démarrant à partir d'une liste initiale annotée manuellement. Dans les sections suivantes, nous abordons d'abord la question des différences entre les mots abstraits et concrets en termes de leur traitement cognitif (section 2). Dans la section 3, nous présentons une caractérisation des mots abstraits et concrets et des méthodes de TAL qui ont été utilisées pour l'annotation du lexique abstrait et concret.

Les trois sections suivantes présentent des stratégies et mises en place du travail expérimental. Dans la section 4, nous décrivons la méthodologie et les résultats de notre première expérience destinée à déterminer le rôle de la fréquence dans la reconnaissance d'un mot polysémique comme abstrait ou concret. Dans la section 5, nous présentons la méthodologie et les résultats de la deuxième expérience qui a au pour le but de déterminer l'impact de la morphologie (mots construits) sur la reconnaissance d'un mot comme abstrait ou concret. La section 6 comprend la troisième expérience qui vise à créer la base de données lexicale des noms abstraits et concrets à partir de la liste de mots initiale et l'évaluation d'échantillon de 120 mots extrait de la base pour déterminer si l'annotation automatique correspond aux jugements humains. Nous concluons enfin par une discussion générale sur la disponibilité des données et sur les usages possibles de la ressource.

Notre étude a été réalisée dans le cadre du projet ANR ALECTOR<sup>1</sup> avec le but d'étudier le facteur d'abstractivité, afin de pouvoir l'intégrer dans un outil de simplification de textes destiné à des enfants dyslexiques et faibles lecteurs.

Les objectifs principaux de notre travail étaient :

- Identifier s'il existait des traits formels pouvant nous indiquer le degré de l'iconicité d'un mot
- Etablir une typologie abstrait/concret
- Explorer la possibilité d'annotation automatique du lexique selon le degré d'abstractivité
- Construire une base de données des noms abstraits et concrets à partir de méthodes distributionnelles
- Evaluer les méthodes de construction et la base de données annotées obtenues (en comparant l'échantillon avec des annotations faites automatiquement et annotations humaines)

## 2. Caractérisation des mots abstraits et concrets

### 2.1. Aspects neuropsycholinguistiques

Il est considéré dans la littérature qu'il existe deux systèmes du codage de l'information verbale et visuelle. Ces codages sont réalisés dans des régions du cerveau différentes, ce qui est prouvé par des expériences avec l'utilisation de l'imagerie par résonance magnétique fonctionnelle (IRMf) et l'électroencéphalographie (EEG).

M. Just et ses collègues (Just et al., 2004) ont observé que les mots abstraits sont souvent associés aux régions du cerveau touchées chez les enfants dyslexiques ; cependant, il est nécessaire de distinguer les types de dyslexie. Les erreurs sémantiques ne sont pas uniquement présentes dans les cas de dyslexie profonde, acquise et non observée chez les enfants. Une expérience avec des phrases de haute et basse iconicité montre que cela prend plus de temps de répondre 'vrai' ou 'faux' pour les phrases avec un haut degré d'iconicité.

Cependant, Paivio a constaté que l'effet de concrétude, caractérisé par le temps de réponse plus rapide à des mots concrets, ne se produit pas seulement face à des données provenant

---

<sup>1</sup> ANR ALECTOR. Consulté le 20 mai 2020, à l'adresse <https://alectorsite.wordpress.com/>.

d'individus souffrant de dyslexie profonde, mais également chez les normo-lecteurs (Paivio, 1991). Les principales théories expliquant l'effet de concrétude chez les normo-lecteurs incluent la théorie du double codage (Paivio, 1990) qui soutient que les mots concrets ont un avantage en termes de traitement car ils activent le système verbal (linguistique) et le système non-verbal (d'imagerie), tandis que les mots abstraits activent seulement le système verbal (Figure 1).

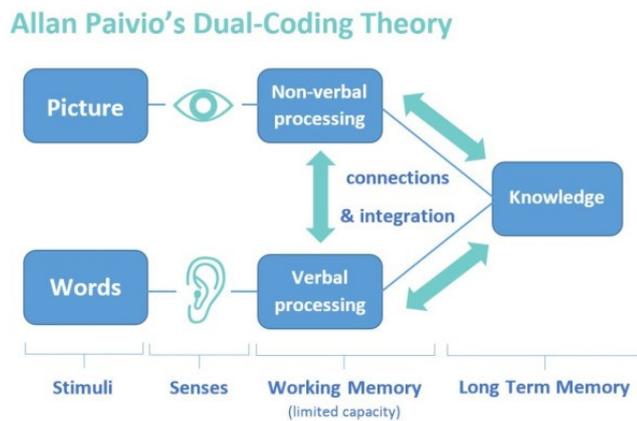


Figure 1. Système de double codage

Une autre explication de l'effet de concrétude est la théorie de la disponibilité du contexte (Schwanenflugel et al., 1988; Schwanenflugel & Shoben, 1983; Schwanenflugel & Stowe, 1989) qui soutient que les mots concrets sont fortement associés à quelques contextes, tandis que les mots abstraits sont faiblement associés à de nombreux contextes.

L'effet de concrétude est toutefois exagéré chez les personnes souffrant de dyslexie profonde, de sorte qu'il peut être impossible de lire des mots abstraits à cause d'un déficit sémantique de ces mots. Certaines preuves suggèrent que cet effet de concrétude exagéré se reflète également dans les différences d'activation neuronale chez les normo-lecteurs et les personnes aphasiques (Sandberg & Kiran, 2014). Diverses théories expliquant l'effet de concrétude dans la dyslexie profonde ont été proposées. Selon l'hypothèse de Coltheart (Coltheart et al., 1988) l'hémisphère gauche permet une lecture abstraite des mots. Les lecteurs souffrant de dyslexie profonde ayant des dommages à l'hémisphère gauche utilisent fortement l'hémisphère droit, ce qui entraîne des difficultés avec les mots abstraits. Morton et Patterson (1980) proposent un modèle à double voie dans lequel la dyslexie profonde résulte de lésions multiples. Dans ce modèle, la lecture s'effectue via la voie sémantique ; cependant, la sémantique des mots abstraits est altérée. De même, Plaut et Shallice (1993), dans leur modèle connexionniste, proposent un avantage pour la lecture de mots concrets, car les mots concrets sont plus simples à caractériser que les mots

abstrait. En outre, le modèle des différents cadres de représentation de Crutch et Warrington (2005) propose que les mots concrets sont représentés dans un cadre catégorique (basé sur la similarité sémantique) et les mots abstraits sont principalement représentés par une association sémantique (contextes linguistiques). Cette théorie soutient que les mots concrets partagent davantage de représentations avec d'autres mots similaires (par exemple, *vache – mouton*) qu'avec d'autres mots associés (par exemple, *vache – étable*), tandis que les mots abstraits partagent davantage de représentations avec d'autres mots associés (par exemple, *vol-punition*) qu'avec d'autres mots similaires (par exemple, *vol - crime*). En conséquence, les lecteurs souffrant de dyslexie profonde produisent plus d'erreurs associatives, comme *vol – punition*, que d'erreurs sémantiquement similaires, comme *vol - crime* en réponse à des mots cibles abstraits et plus d'erreurs sémantiquement similaires que des erreurs associatives en réponse à des mots cibles concrets.

## 2.2. Sur les notions d'abstrait et de concret

La caractérisation des mots en concrets et abstraits reste une tâche difficile. Premièrement, par des mots concrets on comprend des mots qui ont un degré élevé d'iconicité. Selon Tellier et al., (2018), les mots concrets sont associés à une grande iconicité, notamment en termes de représentation mentale, tandis que les mots abstraits sont plutôt encodés verbalement (Paivio, 1986). Les mots concrets sont davantage associés aux informations contextuelles et aux expériences sensorimotrices que les mots abstraits, dans la mesure où les mots concrets sont liés à une haute iconicité et les mots abstraits à une faible iconicité.

La notion de nom 'concret' fait référence aux objets, matériaux, sources de sensations relativement directes (Gorman, 1961); la notion de nom 'abstrait' fait référence à des objets, des matériaux et des sources de sensations indirectes. Un mot peut être générique (nommer un groupe ou une catégorie) ou spécifique (nommer une idée ou un objet spécifique) et abstrait, ainsi que générique ou spécifique et concret. On classe comme 'abstrait' tous les noms de mesures, processus, types d'humains, avec un trait sensoriel. Les noms des créatures mystiques sont classés comme concrets. Les états, périodes et qualités, phénomènes et événements sont classés comme abstraits.

La notion de concrétude concerne aussi les mots qui peuvent être ressentis par l'un des cinq sens (Dove, 2016). Les mots concrets se réfèrent généralement à des concepts qui sont

spatialement et physiquement perceptibles, alors que les mots abstraits se réfèrent souvent à des concepts composés d'information sociale ou introspectif (Danguécan & Buchanan, 2016) (cf. Table 1).

<b>Mots abstraits</b>		<b>Mots concrets</b>	
Processus, états et périodes	confinement, espoir, mois	Perceptibles spatialement	table, arbre
Mesures et qualités	degré, gentillesse	Physiquement perceptible par l'un des cinq sens	musique, arc-en-ciel, amertume
Phénomènes et événements	conseil, soirée	Tous les êtres vivants	femme, chat
Traits d'humains	menteur, génie	Creatures mythologiques	troll, dragon

Table 1 : Typologie des noms abstraits et concrets.

Une classification binaire des mots en abstraits et concrets, cependant, reste assez subjective, premièrement, parce que chaque personne a une expérience linguistique différente, et deuxièmement, parce que dans le vocabulaire de chaque langue, il y a beaucoup de mots polysémiques qui souvent ont des significations liées à différentes catégories sur l'échelle de l'iconicité.

Même si la nature binaire d'une telle division peut sembler être un obstacle à la classification, dans cette étude, nous adhérons à une telle binarité. On suppose que si des études précédentes ont pu prouver la différence dans la perception des mots abstraits et concrets par le cerveau humain, la ligne entre l'abstrait et le concret existe dans le lexique et peut se refléter dans des caractéristiques spécifiques inhérentes au vocabulaire. En revanche, cette binarité n'est pas absolue : à la lumière des résultats de notre évaluation par des humains (cf. section 6.3) il y a une certaine gradation dans la perception de l'iconicité. Par exemple, 'gare' sera perçu comme très concret, 'signe' ou 'nation' au milieu de l'échelle, et 'manie' comme plutôt abstrait.

Il existe cependant quelques bases de données contenant des informations sur les mots abstraits. Elles reflètent généralement les résultats d'annotations humaines, contiennent moins d'un millier de mots, peu de traits sémantiques ou lexicales. (Brysbaert et al., 2014; Bonin et al., 2003). Par exemple, Ferrand et Alario ont utilisé une base de données contenant 260 mots abstraits

(Ferrand, 2001) et 366 mots concrets (Ferrand & Alario, 1998) afin de mener une expérience d'associations de mots. Ces listes de mots hors contexte avec le niveau d'iconicité ont été compilées sur la base des corpus américains et canadiens traduits et approuvés par des francophones. Une ressource comme JeuxDeMots (Lafourcade, 2007), réseau lexical de référence pour le français, ne contient pas, à ce jour, des informations de ce type.

### 2.3. Méthodes d'annotation de mots abstraits et concrets

Différentes tentatives de construction de listes annotées de mots abstraits et concrets sont décrites dans la littérature. Rabinovich et al. (2018) utilisent une approche faiblement supervisée pour prédire l'abstractivité des mots et des expressions en l'absence totale de données étiquetées. Ils exploitent uniquement les indices morphologiques en tant que suffixes et préfixes et l'environnement contextuel d'un mot tel qu'il apparaît dans le texte. Leurs résultats montrent que les indices proposés sont suffisamment puissants pour obtenir une forte corrélation avec les marqueurs humains. Les résultats démontrent également qu'un indice morphologique minimum et un corpus textuel sont suffisants pour fournir quelques prédictions. Les auteurs ont utilisé l'ensemble des « indicateurs d'abstractivité » en anglais, comme les suffixes *-ness*, *-ence*, *-ety*, *-ship* etc.

D'autres recherches en anglais montrent différents degrés de concrétude pour les formes de mots construits dans la représentation mentale. Les mots à structure opaque (*'departement'*) peuvent être plus difficiles à catégoriser que les mots qui peuvent être facilement décomposés en une racine avec une forte signification sémantique et un morphème qui forme le dérivé (*'happiness'*) (Marslen-Wilson et al., 2013).

Avec l'essor récent des techniques de plongement de mots (ou *word embeddings*), les méthodes de construction ont évolué permettant d'étendre automatiquement les réseaux de distribution en utilisant les informations de proximité sémantique comme vecteurs. Des études impliquant l'utilisation des algorithmes de *word embedding* pour prédire le caractère concret des mots dans une langue et entre les langues ont été proposées par Ljubešić et al. (2018). Pierrejean & Tanguy (2019) ont également étudié le problème de la stabilité du plongement des mots en fonction de l'affectation à la catégorie concreté ou abstraite. Les résultats de cette étude ont montré que la propagation de mots concrets est plus performante que la propagation de mots abstraits. Enfin, Abnar et collaborateurs (2018) ont mené des expériences en utilisant plusieurs algorithmes

pour comparer leurs performances aux résultats de l'activité cérébrale dans le but de trouver une meilleure solution pour arriver à la classification des noms en abstraits et concrets.

L'approche par plongement de mots est très puissante en TALN. Cependant, elle a des inconvénients, de même que de nombreux autres mécanismes d'apprentissage automatique, à savoir, le fait qu'il représente souvent une 'boîte noire' pour le chercheur : ce qui se passe à l'intérieur de l'opérateur de l'algorithme reste vague et limité à l'interprétation des résultats (Chen et al., 2018). Dans notre étude, nous nous intéressons non seulement à ce qui se passe après l'application d'un algorithme de TALN, mais aussi quelle est la différence entre les résultats de l'annotation automatique et du jugement humain, et pour quelle catégorie, abstrait ou concret, on peut obtenir moins de différence dans les résultats. Notre objectif est de rendre possible une propagation à partir d'une liste de mots abstraits et concrets annotée manuellement et de savoir si cette propagation fonctionne mieux pour les noms abstraits ou pour les noms concrets. Notre hypothèse est que les noms abstraits sont sémantiquement liés à d'autres noms abstraits et que les noms concrets sont sémantiquement liés à des noms concrets. On évite d'utiliser le terme 'synonymes' car les méthodes qu'on utilise dans cette étude en plus des synonymes incluent d'autres relations lexicales telles que les analogies, les antonymes et les associations de mots.

Le voisinage sémantique des mots peut être utilisé dans des algorithmes d'apprentissage automatique qui se concentrent sur la récupération de différents types d'informations sémantiques et lexicales afin d'améliorer la désambiguïsation des mots abstraits et concrets. Ces études sont souvent placées à la frontière de différents domaines scientifiques. Une étude de Hessel et al. (2018) a prouvé que les concepts concrets sont plus facilement reconnus par les algorithmes du TALN, et la méthode distributionnelle de *k* plus proches voisins (*k-nearest neighbors*) fonctionne mieux et est plus applicable pour les mots abstraits que pour les mots concrets. Cela peut être expliqué par l'existence de milliers d'images de mots concrets sur Internet, qui peuvent être facilement associés aux mots. Pour les mots abstraits, ces images auront une représentation moins homogène. Une autre étude (Reilly & Desai, 2017) soutient cependant que la densité de voisinage sémantique est plus élevée pour les mots concrets.

Comme on peut le voir, la variable sémantique de l'iconicité et son impact sur le comportement des mots associés restent non étudiés. Les recherches en sciences cognitives et en linguistique informatique tentent de faire un parallèle entre les schémas de traitement en cerveau et les algorithmes artificiels et statistiques. Hultén et al. (2018) ont réussi à montrer que le décodage neuronal du sens abstrait ou concret des mots est fondé sur le système cognitif verbal à

travers les régularités de l'usage et peut être recréé en utilisant des algorithmes de calcul mesurant ces régularités.

Dans la suite du mémoire, nous allons présenter différentes expériences visant à :

- 1) caractériser le lexique abstrait et concret ;
- 2) observer la dépendance de la structure morphologique des mots et leur niveau de l'abstractivité ;
- 3) créer une base de données des noms abstraits et concrets ou ce trait soit explicite.

Notre première expérience vise à observer l'impact de la fréquence des noms polysémiques sur la reconnaissance de ces noms comme abstraits ou concrets (Goriachun, 2019a). 36 stimuli abstraits de haute et basse fréquence et 36 stimuli concrets de haute et basse fréquence ont été placés dans les contextes cohérents et présentés aux participants sous la forme de questionnaire. La tâche des participants était de classer les noms en gras (stimuli) comme abstraits ou concrets. La deuxième expérience vise à tester l'hypothèse de la dépendance de la structure morphologique des noms sur leurs niveau d'abstractivité. Nous avons établi le questionnaire consistant en 10 mots abstraits construits avec 10 synonymes simples et 10 mots concrets construits avec 10 synonymes simples et 10 mots très concrets simples. Le questionnaire avec les stimuli sans contexte a été présenté à des juges humains sous la forme de Google Form et comprenait 4 choix pour le classement : abstrait, plutôt abstrait, plutôt concret et concret (Goriachun, 2019b). La troisième expérience vise à construire automatiquement la base des noms abstrait et concrets à partir de la liste initiale de 61 noms. L'expérience comprend deux étapes : la propagation des informations d'une liste initiale annotée manuellement à l'aide des méthodes distributionnels (voisins distributionnels et cooccurrences syntaxiques) et l'évaluation du résultat par comparaison avec des jugements humains. Utilisation de deux méthodes distributionnelles nous a permis de tester l'hypothèse de la différence dans l'organisation des réseaux sémantiques des noms abstraits et concrets. Nous avons pu identifier la pertinence des méthodes dans la tâche d'annotation automatique selon la catégorie abstraite/concrète. La deuxième étape consistait d'évaluer l'échantillon de 120 noms (60 noms concrets (30 voisins distributionnels et 30 cooccurrences syntaxiques) et 60 abstraits (30 voisins distributionnels et 30 cooccurrences syntaxiques). L'évaluation a été menée en ligne avec 1083 participants au total. Tous les participants ont eu pour a tache de classer les stimuli sans contexte selon l'échelle glissière de -100 (très abstrait) a 100 (très concret) (Goriachun et Gala, 2020.).

### 3. Expérience 1. Impact de la fréquence sur l'abstractivité du lexique

Les informations théoriques présentées dans la première partie du rapport nous permettent de constater que d'après la littérature, des mots concrets et abstraits ont un impact sur la complexité de la perception d'un texte par les personnes ayant des difficultés de lecture.

Étant donné qu'une expérience avec un groupe contrôle et un groupe de personnes ayant des troubles du langage était impossible, nous avons décidé de commencer par mener une expérience avec des normo lecteurs. Le but de cette première étude était de savoir si les mots polysémiques abstraits ou concrets étaient facilement repérables par des normo lecteurs.

Nous avons établi l'hypothèse que la reconnaissance d'un mot comme concret ou abstrait est affectée par la fréquence d'un mot donné, des mots fréquents sont plus concrets par rapport à des mots moins fréquents. Nous avons ajouté à cette hypothèse la possibilité de corrélation avec la longueur d'une phrase (la quantité de contexte) et la position du mot cible dans la phrase.

#### 3.1. Stimuli

Les mots fréquents et les mots rares ont été sélectionnés à l'aide d'une base de données Lexique (New et al., 2001). Le critère de fréquence de cette ressource a été obtenu à partir des livres (plus de 50 pour les mots fréquents, moins de 30 pour les rares). La longueur des mots était contrôlée - tous les stimuli se composent de sept lettres et appartiennent à la catégorie grammaticale de nom. Parmi les mots ont été sélectionnés ceux qui, selon la ressource ReSyf (Billami et al., 2018), base lexicale de synonymes classés par le degré de difficulté et désambiguïsés sémantiquement, ont plusieurs sens abstraits et concrets (par exemple, 'branche', 'rapport', 'contact', 'courant').

18 mots fréquents et 18 mots rares ont été sélectionnés, chacun avec deux sens. Au total, 72 stimuli (cf. Annexe 1).

Pour l'annotation il était nécessaire de trouver un contexte approprié pour chaque mot : un contexte aussi proche que possible du sens abstrait ou concret du concept (*La dernière fois qu'il est allé en prison, il a explosé le **plafond** de ses cartes de crédit.* – sens abstrait ; *À cause du trou*

*dans mon **plafond**, une chouette est entrée.* – sens concret.) Les contextes ont été extraits d'un corpus parallèle.<sup>2</sup>

### 3.2. Participants et tâche

L'expérience a été réalisée en mars 2019. 27 personnes l'ont participé gracieusement, dont 25 francophones natifs et 2 avec une bonne maîtrise de la langue (niveau C). L'âge moyen des participants est 26 ans avec un niveau universitaire entre Bac+1 et Bac+8 (Annexe 2).

72 stimuli avec des contextes ont été randomisés et mélangés pour éviter de répéter le même mot un par un. Le questionnaire comprenait trois colonnes, une colonne avec des phrases et deux colonnes pour la réponse concrète et abstraite (cf. Annexe 1).

La tâche des annotateurs était de lire la phrase et classer le mot cible en contexte comme abstrait ou concret. Notre étude a consisté à analyser la perception des mots concrets et abstraits par les lecteurs, afin de déterminer la dépendance du niveau d'iconicité principalement par rapport à la fréquence du mot (Annexe 3).

Les paramètres de comparaison et d'analyse, obtenus à l'aide de la méthode de calcul Kappa de Fleiss (1981) et la formule simplifiée de calcul de gamma de Gwet (2008) pour obtenir l'accord entre les annotateurs pour chaque stimulus. Aussi pour le niveau d'iconicité a été pris le pourcentage de réponses 'concret' par rapport à toutes les réponses pour chaque stimulus. Nous avons utilisé le logiciel R pour tous les opérations statistiques.

### 3.3. Résultats et discussion

Le niveau d'accord de Kappa de Fleiss était 0,309, ce qui correspond à un accord faible (attendu pour la tâche de décision sémantique). Les résultats montrent qu'il n'y a pas de rapport entre la fréquence et l'accord ni pour les stimuli avec l'accord important ni pour ceux avec l'accord moyen et faible. En général, il y a même l'effet inverse : parmi les stimuli avec la plus grande valeur d'accord, il y a plus des mots rares et parmi les stimuli avec un accord faible il y a plus de mots fréquents. Pour les stimuli avec l'accord important les participants ont annoté 16 stimuli fréquents et 20 stimuli rares, pour l'accord faible : 20 stimuli fréquents et 16 rares (cf. Annexe 4).

---

<sup>2</sup> Linguee. Consulté le 20 mai 2020 à l'adresse <https://www.linguee.fr/>.

Une telle distribution peut être expliquée par le fait que les mots fréquents ont une application plus large et ont éventuellement plus de synonymes, ce qui rend difficile leur définition dans une catégorie particulière.

Si on compare ces résultats avec ce que nous avons obtenu pour les stimuli avec le plus grand nombre de réponses ‘concret’, comparé aux stimuli avec le plus grand nombre de réponses ‘abstrait’, dans le groupe avec un accord important entre les participants, il y a plus des mots concrets (22 contre 14), et dans le groupe avec l’accord faible, plus des mots plus abstraits (19 contre 17). Les stimuli qui ont obtenu plus de 50% des réponses ‘concret’ ont été considérées comme concrets et moins de 50% des réponses ‘concret’ ont été considérés comme abstraits (cf. Annexe 5). Cela s'explique par le fait que des mots concrets, possédant une forte iconicité, ont une réflexion matérielle dans le monde et sont plus faciles à imaginer que des mots plus abstraits.

La longueur de la phrase n’a eu aucun impact sur la définition du mot cible dans le stimulus comme étant plus ou moins concret, ni sur le degré d’accord entre les participants. La phrase la plus courte comportait 47 caractères, la plus longue 127, la longueur moyenne des phrases étant de 78 caractères.

Enfin, la position du mot n’a pas eu d’impact sur les choix des participants. Les mots cibles étaient placés au début, au milieu et à la fin des phrases, les stimuli les plus abstraits et les plus concrets étant à la fois des phrases avec le mot cible dans la position initiale, et avec le mot cible à la fin et au milieu de la phrase.

### 3.4. Biais d’étude

Au cours de cette expérience nous avons utilisé les stimuli avec des contextes. Cependant, les études montrent (Swaab et al., 2002) que l’effet du contexte dans l’expérience du jugement sémantique avec de paires des mots possédant un haut et bas degré d’iconicité peut annuler l’analyse de stimuli par rapport à l’iconicité. Ceci prouverait que le contexte n’est pas utile pour tester l’abstractivité.

Dans notre expérience nous avons pris comme les stimuli les mots polysémiques, ayant deux ou plus significations, l’une étant est nécessairement concrète et l’autre abstraite, par exemple, ‘branche’, ‘fortune’, ‘passage’. En revanche, nous a été impossible de vérifier si les fréquences d’usage de ces deux notions sont identiques ou très différentes.

Une autre expérience intéressante (Jager & Cleland, 2016) nous incite à penser que la présence de la polysémie dans les mots-stimuli a pu avoir un impact sur les résultats aussi. Les chercheurs ont observé que dans les stimuli concrets il n’y avait pas d’effet de polysémie dans la tâche de décision lexicale, cependant dans les stimuli abstraits la polysémie a été liée aux temps de réaction moins importants dans la même tâche.

#### 4. Expérience 2. Impact de la morphologie sur l’abstractivité du lexique

Dans notre deuxième expérience nous explorons le lien entre l’abstractivité et la structure morphologique des mots français par le biais d’une étude qui consistait à analyser la différence entre la perception des mots concrets et abstraits avec la structure morphologique (mots construits ou simples) Notre objectif était de déterminer la dépendance du niveau d'iconicité par rapport à la structure morphématique des mots. L’expérience a été menée dans le cadre d’un stage d’été entre mai et juillet 2019.

Notre hypothèse était que les noms abstraits construits des certains suffixes (ex. **fiction**, **jugement**) sont perçus comme plus abstraits par rapport à leurs synonymes morphologiquement simples (ex. conte, critique).

##### 4.1. Stimuli

Pour cette tâche de décision lexicale nous avons choisi 20 mots avec des suffixes dont on fait l'hypothèse qu'ils augmentent la valeur d’abstractivité (cf. Table 1). Les suffixes choisis sont les plus fréquents dans Manulex<sup>3</sup> la base de données lexicales de la langue française créée à partir de 54 manuels scolaires (Lété, 2004). Pour chaque stimuli abstrait et concret avec les suffixes nous avons choisi les mots non construits (cf. Table 2).

Fonction	Suffixes
action, résultat de l'action	-ation, -ition, -(s)sion, -xion, -isation
	-(e)ment
qualité, propriété, fonction	-ance, -ence, -escence
	-eur

<sup>3</sup> Manulex. Consulté le 2 février 2020, à l’adresse <http://www.manulex.org/>

état	-age
------	------

Table 2. Suffixes nominaux et ses fonctions.

Stimuli			
Abstrait construits	Abstrait simples	Concret construits	Concret simples
mission fiction	rôle conte	Station édition	arrêt revue
traitement jugement	soin critique	mouvement logement	geste studio
patience puissance	calm pouvoir	licence audience	diplôme public
terreur longueur	distance crainte	erreur secteur	faute zone
usage avantage	emploi succès	baggage garage	valise abri

Table 3. Stimuli.

Nous avons constitué une liste de 40 mots, leur fréquence ne dépassait pas 170 dans Lexique 3, (la valeur plus basse était 6, les synonymes (ou possible) avaient les fréquences proches). La fréquence moyenne des mots abstraits était 36,5, celle des mots concrets 40,2.

10 mots très concrets (au début appelés '*fillers*') ont été ajoutés pour comparer les noms construits et les noms sémantiquement proches des mots construits avec la structure lexicale simple, et les noms simples qui ont les synonymes relativement simples (forêt, barbe, piscine, coffre, sable, jambe, chocolat, guitare, tasse, singe).

Les stimuli construits ont un seul suffixe et n'ont pas de préfixes, ce qui permet d'observer l'impact sur la perception du mot du morphème particulier (Annexe 6).

## 4.2. Questionnaire et tâche

La tâche pour les participants était de décider pour 50 stimuli si le mot est abstrait ou concret. Dans cette expérience nous avons décidé de changer l'échelle binaire pour une échelle à 4 choix pour deux raisons. Premièrement, lors de la première expérience, de nombreux participants ont indiqué dans les commentaires qu'ils ne pouvaient pas choisir entre deux extrêmes et préféreraient avoir un choix plus flexible. Deuxièmement, une échelle plus large a été utilisée dans les expériences visant à mesurer l'iconicité et la concrétude des mots (Bonin et al., 2003; Brysbaert et al., 2014; Paivio, 1965). Nous avons proposé l'échelle plus large, de 1 à 4 où :

- 1 – abstrait (par exemple, amour, explication) ;
- 2 – plutôt abstrait, mais possibles à visualiser ou ressentir (par exemple, course, chaleur, faim) ;
- 3 – plutôt concret, facile à identifier avec le contexte (par exemple, homme, acteur) ;
- 4 – concret (par exemple, abeille, papier).

Le questionnaire a été présenté aux participants en ligne sous la forme de Google Form (Annexe 7).

### 4.3. Participants

50 participants ont passé l'expérience, parmi eux 49 francophones natifs et 1 participant avec la langue maternelle catalan. L'âge des participants varie de 19 à 60 ans, âge moyenne = 28,7. Nombre de femmes = 42, nombre d'hommes = 7, autres = 1. Les niveaux d'études des participants sont très divers de Bac à Bac+8 et BEP. 46 participants ont étudié d'autres langues étrangères : anglais, espagnol, italien, suédois, russe, chinois, LSF, danois, néerlandais, allemand, norvégien.

### 4.4. Résultats

Les résultats (analyses et visualisation) ont été obtenus à l'aide d'analyses statistiques réalisées avec le logiciel R. Nous avons gardé les stimuli très concrets (*fillers*) pour laisser aux participant la possibilité d'attribuer à certains noms la valeur 4. Cependant, après les analyses statistiques, nous sommes arrivés à la conclusion que nous devons les considérer comme des stimuli au même niveau que les autres. En effet, ils peuvent être un point d'appui fort pour confirmer qu'il y a une différence entre les mots concrets avec des suffixes et ses synonymes et les mots concrets simple et ses synonymes.

Le premier résultat est le taux d'accord entre les participants obtenu avec le calcul de Kappa de Fleiss (1981). Nous obtenons ici une concordance faible (0,226) ce qui est courant dans ce type de tâche de décision lexicale (forte subjectivité).

Pour aller plus dans le détail, nous avons trouvé qu'il y a une grande diversité parmi les réponses des participants pour :

- les mots concrets construits et les mots concrets simples qui sont proches sémantiquement (Annexe 8).

- les mots abstraits construits et les mots abstraits simples qui sont proches sémantiquement (Annexe 9).

Les mots choisis comme les stimuli très concrets ont reçu un accord entre les annotateurs important sauf ‘forêt’ qui a reçu des réponses différentes (Annexe 10). Pour vérifier que les participants étaient en accord pour les stimuli très concrets nous avons calculé le Kappa pour ces 10 mots. Le résultat obtenu est un accord moyen (0,528) et ce qui est assez grand pour la tâche de décision lexicale, mais attendu puisque les mots sont très concrets.

Les moyennes des stimuli nous montrent que les réponses les plus fréquentes sont 2 et 3 - 29 stimuli. Les participants ont eu des difficultés à classer les noms comme très abstraits ou très concrets, en choisissant plus souvent les deux variantes au milieu. Tous les stimuli concrets ont des moyennes supérieures à 2. Les écarts types pour les mots concrets simples et construits et pour les mots abstraits simples et construits sont très proches. Seulement les stimuli très concrets montrent les valeurs des écarts types très différents, plus basses que pour les autres stimuli (cf. Table 3).

Stimuli	Moyenne	ET
Concrets construits	2,86	0,84
Concrets simples	2,96	0,79
Abstrait construits	2	0,84
Abstrait simples	2,14	0,87

Table 4. Les moyens des accords et des écarts types des stimuli. Pour la version détaillée voir l’Annexe 11.

Si on compare les moyennes et les écarts types des mots concrets choisis comme les stimuli essentiels et les mêmes valeurs des stimuli très concrets (*‘fillers’*), nous observons une grande différence entre ces deux types des mots : les moyennes des stimuli concrets construits et simples sont principalement entre 2 et 3, les moyennes des fillers sont tous proches de 4. Cela confirme l’existence de mots perçus comme très concrets indépendamment du contexte, et que ces mots n’ont pas potentiellement des relations sémantiques proches avec des mots construits des suffixes *-tion, -age, -ment, -ence, -eur*.

Nous considérons que ce fait est lié à l'iconicité élevée des certains mots. Cela prouve aussi le fait que les mots construits et les mots sémantiquement liés à ces mots sont moins iconiques que les mots simples qui n'ont pas des synonymes construits des suffixes indiqués ci-dessus.

Le facteur de la fréquence, on suppose, joue un rôle assez important. Des stimuli plus concrets/construits peuvent avoir des synonymes très iconiques mais ce sont des mots avec les fréquences élevées : par exemple, le mot 'logement' a une fréquence = 11 dans Lexique 3. Son synonyme est 'maison' (facile à imaginer) avec une fréquence très élevée – 605 dans Lexique 3. Le mot 'bagage' a une fréquence = 29 dans Lexique 3, son synonyme (ici hyponyme) plus iconique est 'sac' avec une fréquence 124 dans Lexique 3.

Les moyennes des écarts types des groupes des stimuli concrets construits et simples et abstraits construits et simples sont presque pareils. Les mots concrets simples ont des écarts types en peu moins élevés que les autres groupes. Comme c'était prévu, les mots abstraits construits ont des écarts types un peu plus élevés. Les écarts-types montrent dans ce cas le niveau de l'accord entre les participants pour chaque stimulus. Les écarts types bas montrent qu'il n'y a pas trop de variation de valeurs autour de la moyenne – il n'y a pas trop de variation dans le choix de valeur entre 1 et 4 pour le stimulus. Les écarts types hauts montrent qu'il y a beaucoup de variation de valeurs autour de la moyenne – il y a trop de variation dans le choix de valeur entre 1 et 4 pour le stimulus.

Les moyennes totales des mots concrets sont plus grandes, un peu plus élevées pour les mots sans suffixes. Les moyennes des stimuli abstraits et abstraits avec des suffixes sont les plus basses.

#### 4.5. Discussion

Dans cette étude, nous avons voulu analyser la différence entre la perception des mots concrets et abstraits en tenant compte de la structure morphologique des mots du français, afin de déterminer la dépendance du niveau d'iconicité par rapport à la structure morphématique des mots. Nous avons mené une expérience avec 50 stimuli, sans contexte, dont 20 sont des mots simples et 20 des mots construits des suffixes marquant l'abstractivité. 10 mots dans chaque catégorie étaient des mots plutôt abstraits et 10 mots leurs synonymes concrets. 10 mots simples sans relations sémantiques avec les mots construits, ont été ajoutés au début comme des stimuli-fillers. Après on

a regardé les données obtenues pour ces stimuli plus précisément et ce groupe des stimuli est devenu le groupe de référence « très concrets ».

Comme il était attendu les résultats moyennes des mots choisis comme les stimuli concrets sont plus hautes que pour les stimuli choisis comme les abstraits (le plus haut est la moyenne, le plus haut est le niveau de l'iconicité). Cependant les valeurs des écarts types sont relativement proches pour ces deux groupes des mots. Il n'y a pas de grande différence entre les valeurs des stimuli construits et stimuli simples. Les écarts types sont assez grands et les moyennes sont entre 2 et 3. Cela nous permet de dire que les mots construits et leurs synonymes simples avec les fréquences proches appartiennent souvent à la même catégorie sur l'échelle de l'abstractivité. Ce point est important car il nous permet d'envisager l'annotation automatique du lexique avec les traits abstraits/concrets. Cette hypothèse prouve le fait que les mots-stimuli très concrets (fillers) sont les seuls qui ont des valeurs moyennes entre 3 et 4.

Il faut noter qu'au début du projet, au moment de choisir les stimuli, il a été très difficile de trouver les noms plutôt concrets parmi des mots construits des suffixes, les mots 'logement', 'bagage' et 'garage' étaient les seuls qui nous semblaient concrets parmi les mots possédant les suffixes d'abstractivité et cela est confirmé par les résultats d'annotation.

## 5. Expérience 3. Annotation automatique du lexique en traits abstrait et concret

Après avoir mené les deux expériences précédentes, nous sommes arrivés à la conclusion que des variables telles que la fréquence, la longueur du mot, la position du mot dans la phrase, ainsi que la structure morphologique ne corrélaient pas avec le trait abstrait/concret (Goriachun, 2019a). Ces variables formelles n'étaient pas significatives, nous nous sommes alors tournés vers la variable « iconicité ». Ainsi, en fondant nos hypothèses sur des théories liées à la différence dans l'organisation des cadres sémantiques des mots abstraits et concrets (Crutch & Warrington, 2005), nous avons décidé de mener des recherches sur la possibilité d'annoter automatiquement des noms en fonction de leur niveau d'iconicité.

Cette notion n'était pas explicitement encodée dans les ressources existantes. Nous avons décidé d'annoter une liste restreinte initiale avec un petit nombre de mots annotés manuellement, dont on connaît le niveau d'iconicité (Ferrand, 2001; Ferrand et Alario, 1998) et de propager cette annotation à des mots liés par des relations syntaxiques (cooccurrences syntaxiques) et sémantiques (voisins distributionnels) Ainsi nous nous sommes fixés pour objectif de créer une

base de données contenant des noms français annotés automatiquement sur la base de deux méthodes distributionnelles en faisant la propagation du trait abstrait/concret d'un mot vers ses voisins distributionnels et cooccurrences syntaxiques.. Nous visons à proposer une méthode fiable validée par une évaluation humaine importante de plus de 1000 participants pour une annotation automatique qui permet de créer une base de mots annotés de taille importante.

Les cooccurrences syntaxiques sont les mots qui apparaissent fréquemment à côté du mot cible ; les voisins distributionnels sont les mots qui partagent les mêmes contextes (van der Plas, 2009). Nous avons décidé de fonder nos recherches sur la base de ces deux méthodes car elles montrent deux relations distinctes : des liens sémantiques dans le cas des voisins distributionnels et une proximité syntaxique dans le contexte dans le cas de cooccurrences syntaxiques. Ceci peut être crucial pour la détermination et la distinction des mots abstraits et concrets. Par exemple, les mots 'plante' et 'fleur' sont les voisins distributionnels du mot concret 'arbre', tandis que 'branche' et 'ombre' sont ses cooccurrents syntaxiques, c'est-à-dire ils peuvent apparaître dans les mêmes contextes. 'Inquiétude' et 'peur' sont les voisins distributionnels du mot abstrait 'crainte', tandis que 'dissipation' et 'reflet' sont ses cooccurrents syntaxiques.

Par évaluation des résultats de l'annotation automatique nous avons voulu établir quel est le pourcentage d'unités de l'échantillon obtenues à partir de la méthode des voisins distributionnels et des cooccurrences syntaxiques correspond mieux aux résultats de l'annotation humaine.

Nous nous sommes également intéressés aux différences dans la précision de la prédiction entre les mots abstraits et concrets et aux différences de taille des réseaux sémantiques des mots abstraits et concrets, s'il y en avait. Crutch & Warrington (2005) suggèrent que les mots concrets sont organisés selon un principe de similarité sémantique, tandis que les mots abstraits sont organisés par leur association avec d'autres mots. Si on trouve une prévalence quantitative et plus grande homogénéité dans les résultats pour les mots abstraits ou concrets lors de l'extension de la liste principale ou lors de l'évaluation humaine, il sera possible d'indiquer quel modèle est plus pertinent de recréer à l'aide des outils du TAL.

## 5.1. Données

Pour créer notre liste de mots annotés automatiquement en utilisant les voisins distributionnels et les cooccurrences syntaxiques, nous avons d'abord créé une liste initiale de mots contenant des noms abstraits et concrets. Pour cette tâche, nous avons utilisé les mots issus de deux

études (Ferrand, 2001; Ferrand et Alario, 1998) qui contiennent les données annotées pour la langue française selon des échelles d'iconicité et d'abstractivité. Les listes des mots ont été réduites aux noms de haute fréquence et monosémiques.

Notre liste initiale contient 19 noms abstraits (Ferrand, 2001) et 42 noms concrets (Ferrand et Alario, 1998) avec un score de fréquence  $\geq 40$  (selon la base de données lexicale Lexique 3<sup>4</sup>). Les noms abstraits sont souvent monosémiques selon la ressource lexicale ReSyf (Billami et al., 2018). Les noms annotés comme très concrets et avec un indicateur haute fréquence dans Ferrand et Alario (1998) étaient souvent polysémiques. Pour cette raison, nous avons décidé d'utiliser des mots concrets qui n'ont pas de signification abstraite.

## 5.2. Méthodologie

Par la suite, nous avons extrait manuellement de la base de données distributionnelle *Les Voisins de Le Monde*<sup>5</sup> 50 unités lexicales pour chaque nom. L'ensemble de données expérimental initial a été réduit à seulement 50 voisins les plus proches car nous avons identifié qu'après ces 50 premiers voisins les relations sont devenues trop distancées (selon le score de distance fourni par *Les Voisins De Le Monde*). Le type de relation distributionnelle était limitée aux voisins distributionnels car, à cette étape, cette méthode montrait une meilleure cohérence pour notre tâche. La sortie de la méthode des cooccurrences syntaxiques comprenait beaucoup des non noms.

Après les deux premières étapes, nous avons obtenu une liste complète composée de 369 unités (180 noms concrets et 189 abstraits) (cf. Table 4). Les deux listes ont été nettoyées des répétitions et des non-noms. Nous avons aussi supprimé de la liste des mots abstraits les noms avec un niveau élevé d'iconicité (par exemple, le mot 'chat'). De même, si dans la liste des mots concrets figurait le mot 'bonheur', il a été supprimé en raison de son haut niveau d'abstractivité. Cette expérience nous a également permis d'identifier la première différence entre les listes de noms abstraits et concrets : un nombre réduit de voisins des mots concrets, ainsi qu'un grand nombre de non-noms et de noms abstraits dans la liste des mots concrets (c'est pour cette raison que nous avons dû étendre la liste initiale à 42 noms concrets afin que le nombre de voisins pour

---

<sup>4</sup> Lexique 3. Consulté le 19 février 2020, à l'adresse <http://www.lexique.org/>

<sup>5</sup> Les Voisins De Le Monde. Consulté le 4 février 2020, à l'adresse <http://redac.univ-tlse2.fr/voisinsdelemonde/>

la liste expérimentale de noms concrets (180) soit approximativement égale à la longueur de la liste abstraite (189 noms)).

Les mots abstraits ont montré plus de cohérence en termes de catégorie grammaticale et de catégorie abstrait/concret : les unités supprimées étaient principalement des répétitions.

Catégorie	Abstrait	Concret
Listes initiales (Ferrand and Alario, 1998; Ferrand, 2001)	19	42
Avant le traitement manuel	910	1638
Après traitement (données expérimentales)	189	180

Table 5. Listes initiales obtenues manuellement à partir de la ressource *Les Voisins de Le Monde*.

L'étape suivante a été d'extraire automatiquement, pour chacune des 369 unités, leurs voisins distributionnels et cooccurrences syntaxiques à partir de la ressource *Les Voisins de Le Monde*. Nous avons par la suite comparé les résultats obtenus avec chaque approche.

### 5.3. Résultats

À partir des listes initiales presque égales en nombre (180 mots abstraits et 189 mots concrets), nous avons obtenu une sortie quantitativement différente : 62.174 mots abstraits et 31.333 mots concrets, ce qui signifie que le nombre des voisins distributionnels et les cooccurrences syntaxiques des mots concrets comprend la moitié du nombre des voisins distributionnels et cooccurrences syntaxiques des mots abstraits. Après avoir éliminé toutes les répétitions, nous avons obtenu 4.222 mots concrets uniques et 3.676 mots abstraits uniques (cf. Table 6), ce qui nous permet de vérifier la propriété selon laquelle les mots abstraits ont tendance à apparaître dans des contextes plus similaires. En d'autres termes, un mot abstrait aléatoire X est plus susceptible d'avoir un autre mot abstrait aléatoire Y comme voisin distributionnel (ou comme cooccurrent syntaxique) que deux mots concrets aléatoires Z et W, qui auront la possibilité d'apparaître dans le même réseau sémantique.

Catégorie	Abstrait	Concret
Données brutes	62.174	31.333
Données traitées	3.676	4.222

Table 6. Nombre de mots annotés Abstrait / Concrets.

Une différence entre les voisins distributionnels et les cooccurrences syntaxiques a également été observée. Pour les noms concrets, la sortie des voisins et des cooccurrences est presque égale, mais pour les noms abstraits, ces nombres sont très différents (cf. Table 7). Ce résultat montre que les voisins distributionnels (VD) s'avère être la méthode qui fonctionne le mieux pour les mots abstraits, et les cooccurrences syntaxiques (CS) est la méthode qui fonctionne mieux pour les mots concrets. Ce résultat confirme une hypothèse de différences de représentations sémantiques entre les mots concrets et abstraits.

Catégorie	Abstrait	Concret
Données brutes VD	45.340	16.223
Données brutes CS	16.834	15.110
Données traitées VD	2.129	1.631
Données traitées CS	1.546	2.592

Table 7. Nombre de mots Abstrait / Concrets obtenus après propagation des annotations à partir de la liste expérimentale (VD – voisins distributionnels, CS – cooccurrences syntaxiques).

Il est important de signaler que la polysémie était très forte dans les données collectées : beaucoup de mots dans la liste des mots abstraits existaient dans la liste des mots concrets et vice versa. Le phénomène de la polysémie est inévitable dans le lexique et on pouvait s'attendre à ce que nous l'ayons aussi dans la distinction des mots abstraits et concrets. Les mots qui apparaissent uniquement dans une liste abstraite ou concrète après l'application de la méthode distributionnelle, sont indiqués ci-dessous (cf. Table 8).

Catégorie	Abstrait	Concret
VD	551	400
CS	654	1.352
Total	1.205	1.753

Table 8. Nombre des mots Abstrait et Concrets uniques.

Comme on peut le constater, le nombre des mots concrets uniques est plus important que le nombre des mots abstraits uniques. Le nombre des cooccurrences syntaxiques est beaucoup plus grand que le nombre des voisins distributionnels pour les mots concrets. Dans les mots abstraits cette différence vers la plus grande quantité des cooccurrences syntaxiques est moins importante, ce qui confirme les résultats précédents sur les différences de représentations sémantiques entre

mots concrets et abstraits selon les méthodes distributionnelles appliquées. Cependant, on suppose que les mots uniques peuvent avoir une fréquence plus basse que les mots dans les données traitées.

## 5.4. Évaluation

### 5.4.1. Stimuli et procédure

Nous avons voulu évaluer les annotations obtenues automatiquement à partir de *Les Voisins De Le Monde* par rapport à des annotations humaines. Nous avons donc créé un échantillon de 120 noms de la liste initiale : 60 noms concrets (30 voisins distributionnels et 30 cooccurrences syntaxiques des mots de la liste initiale de 180 noms concrets) et 60 abstraits (30 voisins distributionnels et 30 cooccurrences syntaxiques des mots de la liste initiale de 189 noms abstraits), que nous avons soumis à des juges (Annexe 14). L'échantillon a été choisi au hasard dans les données globales pour éviter les biais d'échantillonnage. Nous avons alors créé, à l'aide de C. Zelinski du Centre de Ressources Expérimentales de l'ILCB (CREX<sup>6</sup>) un questionnaire en ligne qui a été distribué parmi les étudiants et le personnel de l'université d'Aix-Marseille via le réseau universitaire (cf. Figure 2). Pour cette expérience nous avons utilisé une échelle de plus grande amplitude, modélisée sous la forme d'une glissière (slider), au lieu de l'échelle de Likert utilisée dans les deux expériences précédentes (cf. Figure 3). Ce choix s'explique par le fait que les attitudes des participants pour un stimulus particulier ont la forme d'un continuum multidimensionnel. L'échelle de Likert est unidimensionnelle et donne seulement quelques options de choix réduites. De même, souvent, les participants évitent de choisir les options « extrêmes » sur l'échelle, en essayant de rester au milieu ce qui est encore plus probable dans la tâche de décision sémantique ou les stimuli causent beaucoup d'hésitation. Dans cette expérience les participants ont dû annoter chaque mot selon une échelle graduée entre -100 (très abstrait) et 100 (très concret). Le test a été lancé le 26/03 à 17h. A 21h plus de mille participants y avaient répondu. Nous avons donc été capables de récolter 1083 réponses en quatre heures.

Les résultats de l'annotation ont été soumis à une analyse statistique avec le logiciel R. Tous les choix inférieurs à 0 ont été considérés comme des choix vers la notion d'abstrait et tous

---

<sup>6</sup> ILCB | Institute of Language, Communication and the Brain. Consulté le 28 mars 2020, à l'adresse <https://www.ilcb.fr/about/crex/>

les choix supérieurs à 0 sont considérés comme des choix vers le concret. Les mots au milieu de l'échelle de -10 à 10 ont été considérés comme très ambigus. Cette échelle nous permet d'observer le degré d'abstractivité d'un mot selon le jugement humain et de le comparer avec le résultat de l'annotation automatique. Dans les consignes pour les participants au début de l'épreuve il était conseillé de ne pas utiliser trop souvent le choix au milieu de l'échelle. Nous avons fait l'hypothèse que l'obtention d'une valeur moyenne indiquera un fort degré de polysémie. Par ailleurs, les participants qui ont trop souvent validé une réponse avec le *slider* au milieu, ont été éliminés de l'étude (nous avons considéré que les réponses étaient invalides).

Figure 2. Page d'accueil d'expérience en ligne.

Figure 3. Présentation du stimulus.

#### 5.4.2. Participants

1.083 participants ont passé l'expérience. Leur âge varie de 18 à 75 ans, âge moyenne = 26 ans. Nombre de femmes = 687, nombre d'hommes = 396. Les niveaux d'études des participants sont très divers de Bac à Bac+8 et BEP. 599 participants ont étudié d'autres langues étrangères et 74 participants n'étaient pas des francophones natifs ce qui n'a pas eu l'impact sur les résultats (voir partie C de la section suivante).

#### 5.4.3. Analyse statistique et résultats

Le but de notre analyse statistique était de comparer la sortie obtenue par l'annotation automatique et le jugement humain. Premièrement, nous voulions savoir si les noms annotés par l'une des deux méthodes distributionnelles comme abstraits ou concrets appartiennent ou non à ces catégories selon l'annotation des participants. Deuxièmement, notre but était d'identifier quelle méthode (voisins distributionnels ou cooccurrences syntaxiques) donne les résultats plus précis dans la tâche de l'annotation du lexique abstrait et concret.

Nos stimuli ont été choisis pour représenter chaque catégorie et chaque méthode. Au début c'étaient 30 voisins distributionnels abstraits, 30 cooccurrences syntaxiques abstraites, 30 voisins distributionnels concrets, 30 cooccurrences syntaxiques concrètes.

Étant donné que certains mots de la catégorie des noms abstraits se trouvent également dans la catégorie concrète, il existe pour eux deux options. Par exemple, le nom 'mairie' est le voisin distributionnel du nom concret 'banque' mais en même temps c'est le voisin distributionnel du nom abstrait 'institution' ; 'milieu' est le cooccurrent syntaxique du nom concret 'salle' mais aussi le cooccurrent syntaxique du nom abstrait 'origine'. Certains stimuli ne sont pas ambigus, ils apparaissent dans une seule catégorie. Par exemple, le nom 'orge' a été annoté dans la catégorie des voisins distributionnels du nom concret 'betterave'. On a décidé de ne pas se débarrasser de ce type d'ambiguïté et de voir si les participants hésiteraient aussi pendant l'annotation des stimuli qui apparaissent aux deux catégories.

##### A. Analyse de la précision des méthodes distributionnelles

La première étape était de comparer les résultats obtenus par l'annotation par le jugement humain avec la sortie obtenue après l'annotation automatique. Comme c'était mentionné ci-dessus, les stimuli ont été également distribués dans chaque catégorie.

Pour convertir les données reçues de l'annotation nous avons calculé les moyennes pour chaque stimulus en se basant sur les résultats de toutes les 1.083 passations. Tous les stimuli avec le niveau moyen supérieur à zéro ont été considéré comme concrets et tous les stimuli avec le moyen inférieur à zéro comme abstraits (cf. Figure 4).

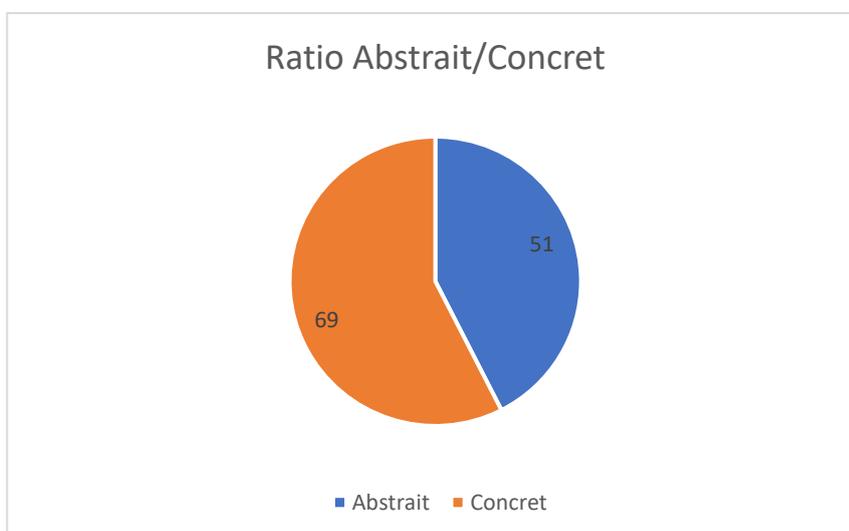


Figure 4. Ratio des noms abstrait et concrets dans les résultats d'annotation.

Catégorie	Abstrait	Concret
VD	21 (70%)	25 (83,3%)
CCS	15 (50%)	19 (63,3%)

Table 9. Nombre de réponses correspondant au sorti et le niveau de précision

En utilisant la formule Kappa de Fleiss (1981), nous avons calculé l'accord entre les annotateurs qui est égal à 0.256 (un accord faible mais attendu dans une expérience de type décision sémantique lexicale avec une grande échelle de -100 à 100). La moyenne des écarts types de stimuli était égale à 56 ce qui confirme une fois de plus la difficulté de la tâche. Les stimuli avec les écarts types moins grands (<40) étaient tous les noms concrets, parmi les stimuli avec les écarts types plus grands (>65) il y avait des noms abstraits et des noms concrets polysémiques. La dépendance linéaire entre les valeurs moyennes des stimuli (leurs niveaux de concrétude) et les écarts types correspond à une corrélation forte ( $r = -0.6210$ ) ce qui signifie que plus grand est le degré de concrétude, plus bas est la valeur des écarts types. En général, plus un mot est concret, moins les participants hésitent lors de l'annotation.

Partant du fait que la tâche de décision sémantique est complexe et ambiguë même pour un humain, nous avons obtenu de meilleurs résultats en utilisant la méthode des voisins

distributionnels. Pour les deux méthodes, les mots concrets ont donné de meilleurs résultats que les mots abstraits, comme le montre la Table 10.

Catégorie	Abstrait	Concret	Total
Voisins Distrib.	21 de 30	25 de 30	46
Précision	<b>70%</b>	<b>83,3%</b>	<b>77 %</b>
Cooccs. Syntaxiques	15 de 30	19 de 30	34
Précision	50%	63,3%	57 %

Table 10. Nombre de correspondances jugement humain / annotation automatique et score de précision.

#### B. Analyse de l'impact de temps de réaction

La plateforme utilisée pour le test nous a également permis de mesurer le temps de réaction des participants. Cependant, notre méthode pour la mise en place du test de la décision lexicale était différente de la méthode habituelle en raison de la présence d'une échelle avec un curseur et un bouton pour confirmer la réponse. Les données obtenues sur le temps de réaction sont très différentes des données fournies par la base Open Lexicon<sup>7</sup> de la ressource Lexique. Par exemple, le temps de réaction pour le mot 'abri' est égal à 570.13 ms dans l'Open Lexicon, dans notre expérience on a obtenu la moyenne du temps de réaction pour le mot 'abri' égal à 3570.064 ms. Nous avons calculé les valeurs de la corrélation entre les variables : moyennes et écarts types et la variable de temp de réaction. La corrélation entre les moyennes et les temps de réaction est faible ( $r = -0.216$ ) ainsi que la corrélation entre les écarts types et les temps de réaction ( $r = 0.251$ ), ceci nous permet de dire que dans notre cas les participants n'ont pas répondu plus ou moins vite selon la catégorie abstraite ou concrète d'un mot ou la valeur des écarts types.

#### C. Analyse de l'impact de la langue maternelle et des troubles du langage

Dans le questionnaire que les participants ont rempli après l'expérience, ils ont indiqué leur langue maternelle et signalé des éventuels troubles du langage oral ou écrit, le cas échéant. Comme un nombre suffisamment important de participants, 141 personnes, ont déclaré qu'ils avaient des problèmes liés au langage oral ou écrit, nous avons décidé d'analyser ces données séparément pour comprendre si les résultats obtenus pour ces participants pouvaient affecter l'image globale. Après nos analyses ces tests n'ont pas montré des différences significatives qui pourraient nous faire

<sup>7</sup> Open Lexicon. Consulté le 30 mars 2020, à l'adresse <http://www.lexique.org/shiny/openlexicon/>.

comprendre que les données de personnes ayant des troubles du langage auraient une incidence sur la totalité des données (moyennes globales et moyennes troubles du langage :  $t = 0,10436$ ,  $p = 0,917$  ; écarts types globales et écarts types troubles du langage :  $t = -1,0765$ ,  $p = 0,2828$  ; temps de réaction globales et temps de réaction troubles du langage :  $t = 0,23679$ ,  $p = 0,8131$ ).

Le groupe de participantes non-francophones natifs était moins nombreux, 74 personnes. Les résultats ont été également comparés avec ceux des francophones natifs. Les différences significatives n'ont pas été observées entre les moyens ( $t = -0.2010$ ,  $p = 0.8409$ ), écarts types ( $t = 0.971$ ,  $p = 0.3321$ ) et temps de réaction ( $t = 1.6954$ ,  $p\text{-value} = 0.092$ ) de deux groupes.

#### D. Analyse d'effet de l'ambiguïté abstrait/concret

En général, l'ambiguïté abstrait/concret que nous avons pu observer à la suite de l'implémentation des méthodes distributionnelles est confirmée par les résultats des jugements humains. Les stimuli dans les catégories dont les personnes doutaient appartenaient souvent à deux groupes sémantiques, abstraits et concrets, selon les résultats de l'annotation automatique. Comme dans Kwong (2013) nos données confirment l'hypothèse qu'en cas de polysémie, une personne choisit une signification concrète, plutôt qu'abstraite. Par exemple, les mots 'cadre', 'échelle', 'cote', 'espèce', 'réserve', 'secours' ont été classés comme concrets.

#### E. Analyse de l'impact de la fréquence des stimuli

Puisque nos stimuli ont été pris au hasard à partir de données annotées automatiquement, leurs fréquences sont différentes. Parmi nos stimuli, il y en a qui sont très rares selon la ressource Lexique 3 et d'autres qui sont fréquents en français. Une analyse de covariance a été effectuée pour confirmer que les participants ont choisi la catégorie d'un mot avec ou sans beaucoup d'hésitation indépendamment de sa fréquence<sup>8</sup>. La fréquence d'un stimulus n'a pas eu l'impact sur les trois variables : moyens ( $r = -0.0039$ ), temps de réaction ( $r = -0.041$ ) et écarts types ( $r = -0.0901$ ).

### 5.4.4. Discussion

Au cours de cette expérience, nous avons effectué deux étapes. La première étape est l'expansion automatique d'une liste de 20 noms abstraits et 42 noms concrets (à partir d'une première vérification manuelle sur un nombre réduit de mots). La deuxième étape consiste à vérifier les résultats de l'annotation automatique à l'aide d'une expérience de comparaison avec des jugements humains.

---

<sup>8</sup> Lexique.org Consulté le 7 mars 2020, à l'adresse <http://www.lexique.org/>

Les résultats de la première étape nous ont montré une différence dans la structure des champs sémantiques de deux catégories de noms abstraits et concrets. Nous avons effectué une extension de notre liste initiale de 61 noms à l'aide de deux méthodes distributionnelles. Dès le début de l'annotation automatique, il est devenu clair que la méthode des voisins distributionnels fonctionne beaucoup mieux pour l'annotation de mots abstraits et que la méthode des cooccurrences syntaxiques avait une légère supériorité quantitative pour les noms concrets. La sortie globale de la méthode de voisins distributionnels est trois fois plus grande que celle de la méthode de cooccurrences syntaxiques. À la fin, après l'élimination de répétitions et des non-mots de notre liste, nous sommes arrivés à 7.898 noms annotés automatiquement. Lors de la deuxième étape de notre expérience, nous avons proposé une expérience par la foule en ligne à la suite de laquelle 1.083 annotateurs humains ont annoté 120 stimuli. Ces stimuli représentaient également 4 catégories : voisins distributionnels concrets, voisins distributionnels abstraits, cooccurrences syntaxiques concrets, cooccurrences syntaxiques abstraits.

Les résultats de l'expérience ont montré que la méthode qui utilise les voisins distributionnels fonctionne mieux pour les noms abstraits et concrets. Nous avons également observé quelques résultats secondaires. Selon les analyses statistiques, les participants doutaient moins lorsqu'ils annotaient des mots concrets. Si un mot avait plusieurs sens, dont l'un abstrait et l'autre concret, les participants le réfèrent à la catégorie des mots concrets.

Étant donné que la méthode des voisins distributionnels a donné des meilleurs résultats dans la tâche de la propagation automatique de la liste initiale de mots et a confirmé sa conformité au jugement humain à un niveau assez élevé (77% de taux de précision par rapport à 57% pour les cooccurrences syntaxiques). Nous prévoyons de continuer à utiliser cette méthode de distribution afin d'élargir la liste des 7.898 mots déjà obtenus et pouvoir intégrer les informations obtenues dans le système de la simplification automatique des textes. La base actuelle est disponible sur demande.

## 6. Conclusion

Notre travail est composé de trois études expérimentales principales. La première expérience consiste à analyser l'impact de la fréquence des stimuli polysémiques sur leur reconnaissance comme abstraits ou concrets chez des normo lecteurs. Dans le but de trouver les traits formels qui permettent d'identifier le nom comme abstrait ou concret nous avons mené notre

deuxième expérience. Elle visait à observer le lien possible entre la structure morphologique des noms et leur niveau d'abstractivité. La troisième expérience avait pour but de créer automatiquement une base des noms abstraits et concrets en propageant automatiquement l'annotation des mots d'une liste initiale avec deux méthodes distributionnelles. Nous avons obtenu plus de 7.000 mots annotés en termes d'iconicité. Une comparaison avec un millier d'annotations humains confirme la faisabilité de la tâche, spécialement avec l'application de la méthode des voisins distributionnels.

Classifier et définir les caractéristiques des mots concrets et abstraits est une tâche complexe. Le degré d'iconicité de chaque mot dépend d'une part de sa perception par une seule personne et d'autre part du contexte. L'objectif principal de notre première expérience : observer la perception des mots abstraits et concrets par les normo lecteurs et essayer de déterminer si un mot appartient à une catégorie particulière tenant en compte un facteur tel que la fréquence.

Au cours de l'expérience, nous sommes arrivés à la conclusion que les résultats de cette étude ne montrent pas la dépendance entre la fréquence d'un mot de son appartenance à la catégorie de l'abstrait ou du concret. Cependant, en analysant l'accord entre les annotateurs pour chaque stimulus, on a aperçu que les mots qui ont reçu plus de réponses Concret ont également suscité un large accord entre les participants. Ce qui nous permet de penser que les mots concrets, possédant une iconicité, sont plus faciles à imaginer que des mots plus abstraits.

L'annotation sémantique reste une tâche difficile à automatiser, premièrement, car il est parfois impossible d'obtenir des résultats précis, même avec une annotation manuelle. L'étude des modèles de perception du vocabulaire, qu'elle soit abstraite ou concrète, est l'une de ces tâches. Dans la deuxième expérience on a observé la dépendance entre la structure morphologique complexe et l'abstractivité d'un nom. Nous avons élargi l'échelle des réponses possibles à 4, donnant ainsi plus de choix aux personnes.

Les paramètres de comparaison et d'analyse ont été choisis en accord entre les annotateurs, obtenus à l'aide de la méthode de calcul Kappa de Fleiss. Les moyennes des réponses des participants sur échelle de 1 à 4 nous montrent le niveau d'abstractivité du mot. Les valeurs des écarts-types montrent l'accord entre les annotateurs pour chaque stimulus.

Comme il été attendu, les moyennes des mots choisis comme les stimuli concrets sont plus hautes que pour les stimuli choisis comme abstrait. Cependant les valeurs des écarts types sont relativement proches pour ces deux groupes des mots. Il n'y a pas de grande différence entre les valeurs des stimuli construits et stimuli simples. Les résultats ne montrent pas la dépendance

attendue de la structure morphologique et l'abstractivité. Les écarts types sont assez grands et les moyennes sont entre 2 et 3. Portant en ajoutant à nos stimuli les mots très concrets (ex. barbe, chocolat) on a prouvé l'existence des mots hautement imaginables (iconiques) dont la valeur 4 dans notre échelle (très concret) était évidente pour tous les participants.

Dans la troisième expérience on a établi un corpus basé sur le trait abstrait / concret avec 7.898 noms et 2.958 unités uniques à partir d'une liste initiale annotée manuellement de 369 noms. Notre hypothèse de la différence dans la représentation sémantique des mots abstraits a été confirmée à l'aide de l'application de deux méthodes distributionnelles : voisins distributionnels et cooccurrences syntaxiques. On a reçu des résultats quantitativement différents pour les noms abstraits et concrets après une extension de nos listes initiales. Les voisins distributionnels étaient plus nombreux pour les noms abstraits, tandis que pour les noms concrets, les voisins distributionnels et les cooccurrences syntaxiques montraient le même résultat (en nombre). On a obtenu la plus grande sortie en noms abstraits, par rapport aux noms concrets, bien que on soit parti de listes initiales de la même taille (cela ne signifie pas qu'il existe plus de mots abstraits dans la langue française). Après avoir supprimé les répétitions, nous avons obtenu des listes de taille presque égale, même avec la prédominance de mots concrets. Ceci nous permet de juger que les mots abstraits ont des réseaux sémantiques plus étendus, qui incluent un plus grand nombre d'autres mots abstraits, par rapport aux mots concrets, dont les réseaux sémantiques sont plus restreints. La différence entre les méthodes distributionnelles (voisins distributionnels et cooccurrences sémantiques) nous semble pertinente pour envisager l'implication d'une méthode ou de l'autre en fonction de la catégorie (abstraite ou concrète) pour l'augmentation future de la taille de notre liste initiale.

Nous envisageons de poursuivre nos recherches dans le but d'enrichir la typologie abstrait/concret avec un nouveau facteur : la valence émotionnelle. Nous prévoyons de vérifier si les valeurs émotionnelles peuvent avoir un impact sur le traitement cognitif des mots abstraits et concrets et faciliter ainsi leur reconnaissance, leur lecture et leur compréhension.

## Bibliographie

- Abnar, S., Ahmed, R., Mijnheer, M., & Zuidema, W. (2018). Experiential, Distributional and Dependency-based Word Embeddings have Complementary Roles in Decoding Brain Activity. *Proceedings of the 8th Workshop on Cognitive Modeling and Computational Linguistics (CMCL 2018)*, 57-66. <https://doi.org/10.18653/v1/W18-0107>
- Billami, M. B, François, T. & Gala, N. (2018). ReSyf: a French lexicon with ranked synonyms. Actes de la 27th International Conference on Computational Linguistics (COLING 2018), Santa Fe, NewMexico, United States, 2570-2581.
- Bonin, P., Méot, A., Aubert, L.-F., Malardier, N., Niedenthal, P., & Capelle-Toczek, M.-C. (2003). Normes de concrétude, de valeur d'imagerie, de fréquence subjective et de valence émotionnelle pour 866 mots. *L'Année psychologique*, 103(4), 655-694. <https://doi.org/10.3406/psy.2003.29658>
- Brysbaert, M., Warriner, A. B., & Kuperman, V. (2014). Concreteness ratings for 40 thousand generally known English word lemmas. *Behavior Research Methods*, 46(3), 904-911. <https://doi.org/10.3758/s13428-013-0403-5>
- Chen, Z., He, Z., Liu, X., & Bian, J. (2018). Evaluating semantic relations in neural word embeddings with biomedical and general domain knowledge bases. *BMC Medical Informatics and Decision Making*, 18(S2), 65. <https://doi.org/10.1186/s12911-018-0630-x>
- Coltheart, V., Laxon, V. J., & Keating, C. (1988). Effects of word imageability and age of acquisition on children's reading. *British Journal of Psychology*, 79(1), 1-12. <https://doi.org/10.1111/j.2044-8295.1988.tb02270.x>
- Crutch, S. J., & Warrington, E. K. (2005). Abstract and concrete concepts have structurally different representational frameworks. *Brain*, 128(3), 615-627. <https://doi.org/10.1093/brain/awh349>

- Danguécan, A. N., & Buchanan, L. (2016). Semantic Neighborhood Effects for Abstract versus Concrete Words. *Frontiers in Psychology*, 7. <https://doi.org/10.3389/fpsyg.2016.01034>
- Dove, G. (2016). Three symbol ungrounding problems : Abstract concepts and the future of embodied cognition. *Psychonomic Bulletin & Review*, 23(4), 1109-1121. <https://doi.org/10.3758/s13423-015-0825-4>
- Ferrand, L. (2001). Normes d'associations verbales pour 260 mots « abstraits ». *L'année Psychologique*, 101(4), 683-721. <https://doi.org/10.3406/psy.2001.29575>
- Ferrand, L., & Alario, F.-X. (1998). Normes d'associations verbales pour 366 noms d'objets concrets. *L'année psychologique*, 98(4), 659-709. <https://doi.org/10.3406/psy.1998.28564>
- Fleiss, J. L. (1981). *Statistical Methods for Rates and Proportions*. 2nd Edition, John Wiley, New York, 38-46.
- Gala, N. et Brun, C. (2012) Propagation de polarités dans des familles de mots : impact de la morphologie dans la construction d'un lexique pour l'analyse d'opinions. Actes de Traitement Automatique des Langues Naturelles (TALN 12). Grenoble, juin 2012.
- Goriachun, D. (2019a). *Impact du degré d'abstraction du lexique dans la lecture et compréhension des mots* (rapport de stage non publié). Université d'Aix-Marseille, France.
- Goriachun, D. (2019b). *L'impact de la morphologie sur l'abstractivité des mots* (rapport de stage non publié). Université d'Aix-Marseille, France
- Goriachun, D. & Gala, N. (2020) Identifying Abstract and Concrete Words in French to Better Address Reading Difficulties. Workshop Tools and Resources to Empower People with Reading Difficulties (READI) at International conference on Language Resources and Evaluation (LREC 2020), poster session, pp. 33-40. Marseille, France.
- Gorman, A. M. (1961). Recognition memory for nouns as a function of abstractness and frequency. *Journal of Experimental Psychology*, 61(1), 23-29. <https://doi.org/10.1037/h0040561>

- Gwet, K. L. (2008). Computing inter-rater reliability and its variance in the presence of high agreement. *British Journal of Mathematical and Statistical Psychology*, 61(1), 29-48.  
<https://doi.org/10.1348/000711006X126600>
- Hessel, J., Mimno, D., & Lee, L. (2018). Quantifying the Visual Concreteness of Words and Topics in Multimodal Datasets. *Proceedings of the 2018 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long Papers)*, 2194-2205. <https://doi.org/10.18653/v1/N18-1199>
- Hultén, A., Vliet, M. van, Lammi, L., Kivisaari, S., Lindh-Knuutila, T., Faisal, A., & Salmelin, R. (2018). Cracking the problem of neural representations of abstract words : Grounding word meanings in language itself. *BioRxiv*, 391052. <https://doi.org/10.1101/391052>
- Jager, B., & Cleland, A. A. (2016). Polysemy Advantage with Abstract But Not Concrete Words. *Journal of Psycholinguistic Research*, 45(1), 143-156. <https://doi.org/10.1007/s10936-014-9337-z>
- James, C. T. (1975). The role of semantic information in lexical decisions. *Journal of Experimental Psychology: Human Perception and Performance*, 1(2), 130-136. <https://doi.org/10.1037/0096-1523.1.2.130>
- Jessen, F., Heun, R., Erb, M., Granath, D.-O., Klose, U., Papassotiropoulos, A., & Grodd, W. (2000). The Concreteness Effect : Evidence for Dual Coding and Context Availability. *Brain and Language*, 74(1), 103-112. <https://doi.org/10.1006/brln.2000.2340>
- Just, M. A., Newman, S. D., Keller, T. A., McEleney, A., & Carpenter, P. A. (2004). Imagery in sentence comprehension : An fMRI study. *NeuroImage*, 21(1), 112-124.  
<https://doi.org/10.1016/j.neuroimage.2003.08.042>
- Kroll, J. F., & Merves, J. S. (1986). Lexical access for concrete and abstract words. *Journal of Experimental Psychology: Learning, Memory, and Cognition*, 12(1), 92-107. <https://doi.org/10.1037/0278-7393.12.1.92>

- Kwong, O. Y. (2013). *New Perspectives on Computational and Cognitive Strategies for Word Sense Disambiguation*. Springer New York. <https://doi.org/10.1007/978-1-4614-1320-2>
- Lafourcade, M. (2007). *Making people play for Lexical Acquisition with the JeuxDeMots prototype*. 8.
- Lété, B. (2004). Chapitre 17 : Manulex : une base de données du lexique écrit adressé aux élèves. Dans : Elizabeth Calaque éd., *Didactique du lexique: Contextes, démarches, supports* (pp. 241-257). Louvain-la-Neuve, Belgique: De Boeck Supérieur. doi:10.3917/dbu.didle.2004.01.0241.
- Ljubešić, N., Fišer, D., & Peti-Stantić, A. (2018). Predicting Concreteness and Imageability of Words Within and Across Languages via Word Embeddings. *Proceedings of The Third Workshop on Representation Learning for NLP*, 217-222. <https://doi.org/10.18653/v1/W18-3028>
- Marslen-Wilson, W. D., Tyler, L. K., Waksler, R., & Older, L. (2013). *Abstractness and transparency in the mental lexicon*.
- Mohammad, S. M. (2016). 9 - Sentiment Analysis : Detecting Valence, Emotions, and Other Affectual States from Text. In H. L. Meiselman (Éd.), *Emotion Measurement* (p. 201-237). Woodhead Publishing. <https://doi.org/10.1016/B978-0-08-100508-8.00009-6>
- Morton, J., & Patterson, K. (s. d.). *Interpretation, or, an attempt at a new interpretation*. 15.
- New, B., Pallier, C., Ferrand, L., & Matos, R. (2001). Une base de données lexicales du français contemporain sur internet : LEXIQUE™//A lexical database for contemporary french : LEXIQUE™. *L'année psychologique*, 101(3), 447-462. <https://doi.org/10.3406/psy.2001.1341>
- Paivio, A. (1965). Abstractness, imagery, and meaningfulness in paired-associate learning. *Journal of Verbal Learning and Verbal Behavior*, 4(1), 32-38. [https://doi.org/10.1016/S0022-5371\(65\)80064-0](https://doi.org/10.1016/S0022-5371(65)80064-0)
- Paivio, A. (1986). *Mental representations : A dual coding approach*. Oxford University Press ; Clarendon Press.

- Paivio, A. (1990). *Mental Representations : A dual coding approach*. Oxford University Press.  
<http://www.oxfordscholarship.com/view/10.1093/acprof:oso/9780195066661.001.0001/acprof-9780195066661>
- Paivio, A. (1991). Dual coding theory : Retrospect and current status. *Canadian Journal of Psychology/Revue Canadienne de Psychologie*, 45(3), 255-287.  
<https://doi.org/10.1037/h0084295>
- Pierrejean, B., & Tanguy, L. (2019). Investigating the Stability of Concrete Nouns in Word Embeddings. *Proceedings of the 13th International Conference on Computational Semantics - Short Papers*, 65-70. <https://doi.org/10.18653/v1/W19-0510>
- Plaut, D. C., & Shallice, T. (1993). Deep dyslexia : A case study of connectionist neuropsychology. *Cognitive Neuropsychology*, 10(5), 377-500. <https://doi.org/10.1080/02643299308253469>
- Rabinovich, E., Sznajder, B., Spector, A., Shnayderman, I., Aharonov, R., Konopnicki, D., & Slonim, N. (2018). Learning Concept Abstractness Using Weak Supervision. *Proceedings of the 2018 Conference on Empirical Methods in Natural Language Processing*, 4854-4859.  
<https://doi.org/10.18653/v1/D18-1522>
- Reilly, M., & Desai, R. H. (2017). Effects of semantic neighborhood density in abstract and concrete words. *Cognition*, 169, 46-53. <https://doi.org/10.1016/j.cognition.2017.08.004>
- Sandberg, C., & Kiran, S. (2014). Analysis of abstract and concrete word processing in persons with aphasia and age-matched neurologically healthy adults using fMRI. *Neurocase*, 20(4), 361-388.  
<https://doi.org/10.1080/13554794.2013.770881>
- Schwanenflugel, P. J. (1991). Why are abstract concepts hard to understand? In *The psychology of word meanings* (p. 223-250). Lawrence Erlbaum Associates, Inc.

- Schwanenflugel, P. J., Harnishfeger, K. K., & Stowe, R. W. (1988). Context availability and lexical decisions for abstract and concrete words. *Journal of Memory and Language*, 27(5), 499-520.  
[https://doi.org/10.1016/0749-596X\(88\)90022-8](https://doi.org/10.1016/0749-596X(88)90022-8)
- Schwanenflugel, P. J., & Shoben, E. J. (1983). *Differential Context Effects in the Comprehension of Abstract and Concrete Verbal Materials*. 21.
- Schwanenflugel, P. J., & Stowe, R. W. (1989). Context availability and the processing of abstract and concrete words in sentences. *Reading Research Quarterly*, 24(1), 114-126.  
<https://doi.org/10.2307/748013>
- Shallice, T. (1988). *From neuropsychology to mental structure*. Cambridge University Press.  
<https://doi.org/10.1017/CBO9780511526817>
- Steady, L. M., & Compton, D. L. (2019). Examining the role of imageability and regularity in word reading accuracy and learning efficiency among first and second graders at risk for reading disabilities. *Journal of Experimental Child Psychology*, 178, 226-250.  
<https://doi.org/10.1016/j.jecp.2018.09.007>
- Swaab, T. Y., Baynes, K., & Knight, R. T. (2002). Separable effects of priming and imageability on word processing : An ERP study. *Cognitive Brain Research*, 15(1), 99-103.  
[https://doi.org/10.1016/S0926-6410\(02\)00219-7](https://doi.org/10.1016/S0926-6410(02)00219-7)
- Tellier, M., Stam, G., & Ghio, A. (2018). « Tout ça c'est abstrait » : Comment le degré d'abstraction d'un mot expliqué affecte-t-il la parole multimodale ? *XXXIe Journées d'Études sur la Parole*, 329-337. <https://doi.org/10.21437/JEP.2018-38>
- van der Plas, L. (2009). Combining syntactic co-occurrences and nearest neighbours in distributional methods to remedy data sparseness. *Proceedings of the Workshop on Unsupervised and Minimally Supervised Learning of Lexical Semantics - UMSLLS '09*, 45-53.  
<https://doi.org/10.3115/1641968.1641974>

Wilson, T., Wiebe, J., & Hoffmann, P. (2009). Recognizing Contextual Polarity : An Exploration of Features for Phrase-Level Sentiment Analysis. *Computational Linguistics*, 35(3), 399-433.  
<https://doi.org/10.1162/coli.08-012-R1-06-90>

## Annexes

### Annexe 1. Questionnaire avec des stimuli. Expérience 1

phrase	Abstrait	Concret
Permettez-moi d'ajouter un troisième <b>facteur</b> qui alimente le conflit.		
Je vais m'occuper d'une nouvelle <b>branche</b> de son cabinet.		
Le même principe doit guider nos travaux en ce qui concerne les armes légères ou de petit <b>calibre</b> .		
La Fondation a fait paraître une dizaine d'ouvrages et publie un <b>rapport</b> d'activité annuel.		
Le monde doit passer d'une décennie d' <b>impasse</b> à une décennie de décision.		
Dans cette position, ils doivent être fixés contre tout déplacement vers l' <b>arrière</b> .		
En plus de ça, une <b>branche</b> d'arbre tombe sur la maison de ma mère.		
Je suis là parce que tu as bouleversé la <b>balance</b> de la nature.		
Nous devons apporter notre <b>soutien</b> à ces personnes et à leurs familles.		
Il ne semble pas que cette <b>famille</b> ait d'autres ancêtres que latins.		
Il s'agit d'un large <b>spectre</b> de médicaments qui sont utilisés pour traiter les troubles de même nature.		
J'étais en colère, mais j'ai fait preuve de <b>retenue</b> .		
Si la <b>fortune</b> vous sourit, pourquoi continuez-vous à demander aussi fort mon appui ?		
Des chercheurs de <b>calibre</b> international ont participé à la conférence.		
Nous avons décidé de suivre ce <b>conseil</b> et de ne pas rechercher de coauteurs.		
Le <b>soutien</b> scolaire doit y être renforcé, en relation étroite avec les familles.		
Sa <b>raideur</b> , son intransigeance, tout ce pour quoi certains l'admirent.		
Les armes du service de police sont conservées dans l' <b>arsenal</b> .		
Lors de la comparaison des prix, il est tenu compte de la date à laquelle le <b>produit</b> laitier a été acheté.		
Maintenez le <b>contact</b> aussi longtemps que vous pourrez.		
Pendant neuf mois, j'ai vu cette <b>attente</b> déchirer ma famille.		
Tous les jours le <b>facteur</b> vient mais je ne reçois pas leur lettre.		
Au plan de la politique étrangère, c'est une toute autre <b>affaire</b> .		
Un <b>créneau</b> est, au Moyen Âge, une ouverture pratiquée au sommet d'un rempart.		
Le bas du <b>montant</b> de la porte doit être coupé avec une scie fine à hauteur de l'épaisseur de l'élément		
Le <b>courant</b> des rapides était parfait pour le kayak.		
Leur <b>hauteur</b> maximale ne devrait pas dépasser 30 mètres.		
Cette année a été très chargée pour ma <b>famille</b> avec deux mariages -et un nouvel enfant à venir.		
On distinguera le <b>circuit</b> long, qui part du producteur pour arriver au consommateur en passant par deux intermédiaires.		
La dernière fois qu'il est allé en prison, il a explosé le <b>plafond</b> de ses cartes de crédit.		
Nous avons pris le <b>circuit</b> le plus court lors de notre voyage.		
De plus, le <b>passage</b> de l'ère industrielle à l'ère de l'information ne se fait pas sans créer de soucis financiers.		

La coopération internationale et l' <b>entente</b> doivent guider notre responsabilité collective à cet égard.		
Une <b>raideur</b> ou une faiblesse extrême d'un ou plusieurs membres peut être notée.		
D'autre part, la <b>défense</b> a présenté des témoins qui ont eux aussi juré avoir vu quelqu'un.		
C'est un <b>produit</b> révolutionnaire qui changera le monde.		
Or, l'Europe, <b>berceau</b> de tant d'innovations capitales, est aujourd'hui en proie au doute.		
Il a coupé le <b>courant</b> pour changer l'ampoule en toute sécurité.		
Dans l'affirmative, veuillez décrire les dispositions ayant un <b>rapport</b> avec la résolution.		
Dans une <b>affaire</b> criminelle, tout est potentiellement utile.		
En règle générale, les frais d'abonnement sont inclus dans le <b>montant</b> du loyer mensuel.		
Le <b>spectre</b> de la méfiance a été brandi entre amis.		
Ils montent à l' <b>arrière</b> du camion et partent avec la personne qui les recrute.		
Dans une entrevue qu'il a accordée aujourd'hui, il nous dit que la <b>mémoire</b> est notre avenir.		
Ne pas laisser entrer en <b>contact</b> avec les yeux et la peau.		
Il faut protéger le <b>berceau</b> de l'enfance contre les maladies et les besoins qui l'assiègent.		
S'il s'agit d'un chemin de terre, l'usager sortant de l' <b>impasse</b> devra laisser la priorité aux autres véhicules.		
Le disque dur est la <b>mémoire</b> à long terme de l'ordinateur.		
Une éducation de qualité doit tenir compte du passé, concerner le <b>présent</b> et être tournée vers l'avenir.		
Le <b>passage</b> piéton est une intersection entre les usagers de la route, au même titre que les intersections entre deux véhicules.		
Cette machine nous indique si la <b>balance</b> est équilibrée.		
Les manifestants poursuivent mardi matin leurs opérations de <b>blocage</b> sur les autoroutes et les dépôts pétroliers.		
Il a même affirmé à un certain moment qu'il avait une <b>entente</b> signée sur son bureau.		
Ton père veut qu'on répare la <b>clôture</b> du terrain de derrière.		
En fait, la <b>cuisine</b> et la salle à manger sont conçues afin de séparer la famille des domestiques.		
À cause du trou dans mon <b>plafond</b> , une chouette est entrée.		
Selon le <b>conseil</b> , ces arguments sont dénués de fondement.		
Puissions-nous mettre notre liberté au <b>service</b> des autres.		
Préparez votre proposition à temps pour la faire livrer au moins un jour avant la date de <b>clôture</b> .		
Ce type d'aide garantit aux consommateurs la spécificité d'une <b>cuisine</b> locale basée sur du poisson frais.		
J'ai soudainement eu un <b>créneau</b> sur mon emploi du temps.		
Le Tribunal prévoit relever le défi grâce à un <b>arsenal</b> de stratégies opérationnelles et globales.		

Les exigences en matière d' <b>habitat</b> peuvent varier à chaque étape du cycle biologique.		
Afin d'éviter les longues files d' <b>attente</b> , les représentants sont invités à remplir le bulletin d'inscription.		
Un autre élément majeur de la transition vers l'indépendance concerne la <b>défense</b> .		
Tu as apporté la <b>fortune</b> nécessaire pour acheter ton titre.		
Un autre aspect inquiétant est la poursuite du <b>blocage</b> au sujet du statut de la région.		
Il vous a choisie parce que vous n'êtes pas à la <b>hauteur</b> .		
Tu ne te souviens pas où tu as eu ce <b>service</b> à thé ?		
La <b>retenue</b> d'eau se trouve en amont du barrage.		
Peut-être qu'elle me pardonnerait mon impolitesse si je lui envoyais un <b>présent</b> .		
Le logement et l' <b>habitat</b> jouent un rôle primordial pour la croissance et le développement des enfants.		

## Annexe 2. Données des participants. Expérience 1.

No	Age	Sexe	Lieu de naissance	Niveau d'études	Langue(s) maternelle	Langue(s) d'usage	Dyslexie
1	20	F	Digne-les-Bains	Bac+2	français	Espagnol (C), Anglais (C), LSF (C)	Non
2	19	F	Marseille	Bac+1	français	LSF (B), Portugais (A), Espagnol (B), Anglais (B)	Non
3	34	H	Montelimar	BEP	français	Anglais (B)	Non
4	22	F	Dijon	Bac+5	français	Anglais (C), Suedois (B), Espagnol (B), Italien (A)	Oui
5	22	F	Avignon	Bac+5	français	Anglais (B), Espagnol (A)	Non
6	27	H	Paris	Bac+5	français	Italien (B), Anglais (B), Russe (B)	Non
7	24	F	Russie	Bac+5	russe	Français (C), Anglais (C)	Non
8	29	F	Paris	Bac+5	français	Anglais (B), Espagnol (B), Arabe (B), Allemand (B)	Non
9	25	F	Châtons-en-Champagne	Bac+3	français	Espagnol (C), Anglais (C), Allemand (C), Italien (A), Persan (A)	Non
10	31	F	Cannes	Bac+7	français	Anglais (B), Espagnol (B)	Non
11	24	F	Maroc	Bac+6	français	Anglais (A)	Non
12	23	F	Amiens	Bac+5	français	Anglais (B)	Non
13	24	F	Hyères	Bac+3	français	Anglais (B), Persan (A)	Non
14	30	H	Paris	Bac+8	français	Anglais (C), Allemand (A), Espagnol (A)	Non
15	24	F	Aix-en-Provence	Bac+6	français	Anglais (C), Espagnol (A)	Non
16	21	F	Avignon	Bac+3	français	Anglais (C), Italien (B), Portugais (A), Provençal (A)	Non
17	27	H	Niaguis (Sénégal)	Bac+3	flaoussa et wolof	Français (C), Créol Bissau (B), Manding (B), Russe, Mandarin (A)	Non
18	22	F	Hyères	Bac+3	français	Anglais (B), Espagnol (B)	Non
19	21	F	Avignon	Bac+2	français	Provençal (A), Anglais (B)	Non
20	20	F	Suresnes	Bac+3	français	Anglais (C), Espagnol (C)	Non
21	24	F	Toulon	Bac+3	français	Anglais (C), LSF (B), Portugais (A), Turc	Non
22	20	F		Bac+2	français	Anglais (C), Espagnol (A), Provençal (A)	Non
23	57	H	Neuville-Entier		français		Non
24	22	H	Strasbourg	Bac+4	français	Anglais (B), LSF (A)	
25	54	H	Avignon	Bac+1	français		Non
26	22	F	Marseille	Bac+2	français	Espagnol (B)	Non
27	20	F	Montpellier	Bac+2	français	Espagnol (B), Anglais (B)	Non

### Annexe 3. Consignes pour les participants. Expérience 1.

#### Données anonymisées

Homme

Femme

Lieu de naissance \_\_\_\_\_

Niveau d'études (entourez)

Bac+1

Bac+2

Bac+3

Bac+4

Bac+5

Bac+6

Bac+7

Bac+8

Langue(s) maternelle(s) \_\_\_\_\_

Autres langues et niveaux (A -débutant-, B -intermédiaire- C -maîtrise-)

Éventuellement, êtes-vous dyslexique ?

Oui

Non

Âge \_\_\_\_\_

#### Consignes pour les participants

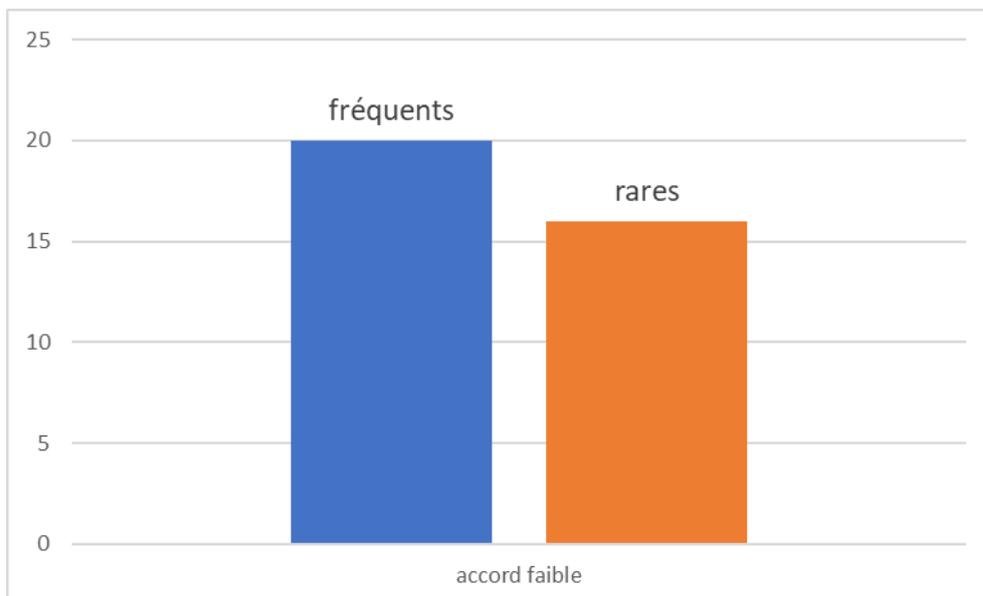
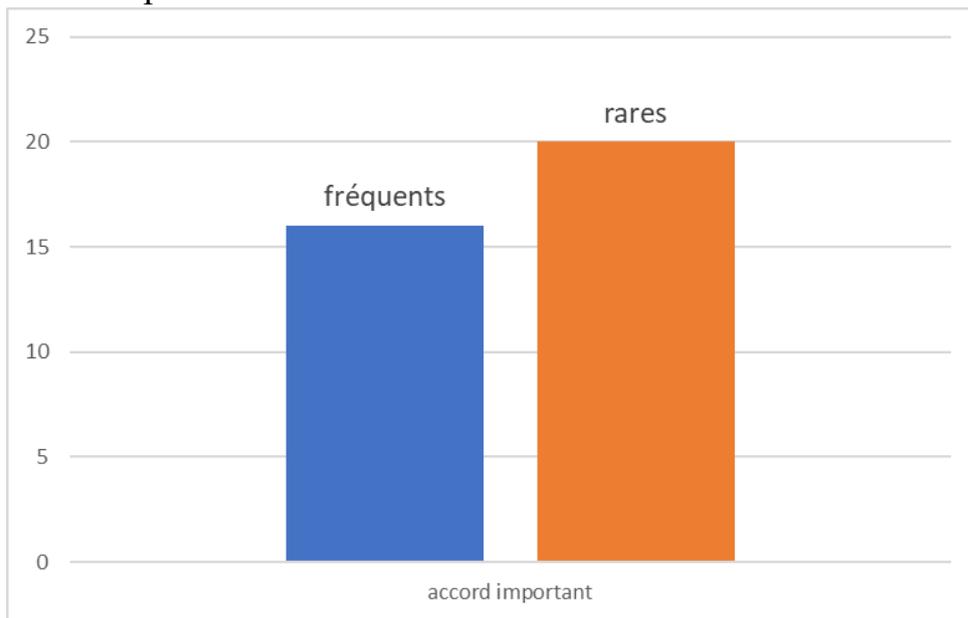
Dans les pages suivantes vous allez lire 72 phrases. Dans chacune il y a un mot cible qui est indiqué en caractères gras. Votre but est d'identifier si le mot en gras a un sens 'concret' ou bien 'abstrait'. Rentrez (A) ou (C) dans la colonne « Abstrait » ou « Concret » selon le cas.

Prenez le temps de bien lire et comprendre chaque phrase. Merci de donner la réponse consciencieusement, selon votre perception du mot dans son contexte (pas de dictionnaire ni trop de réflexion).

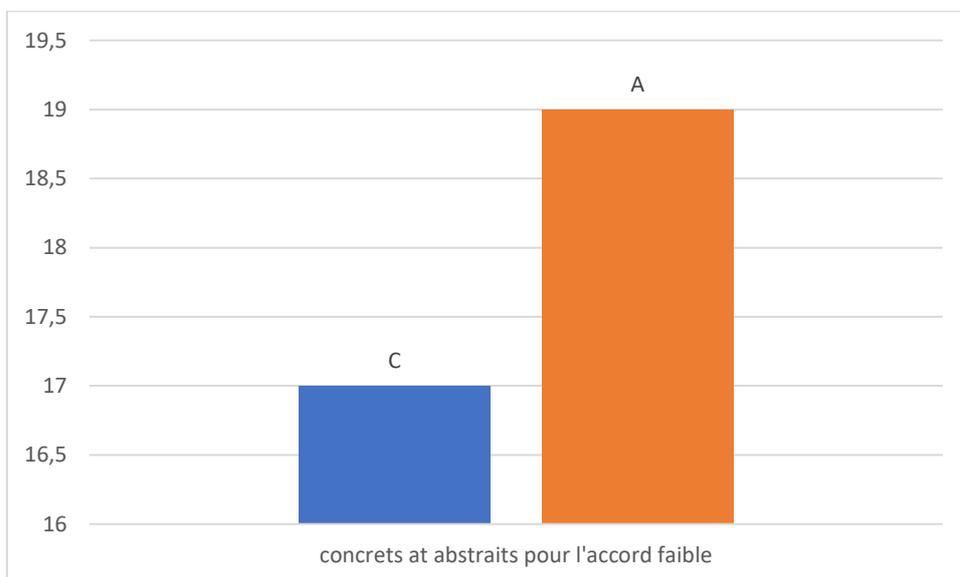
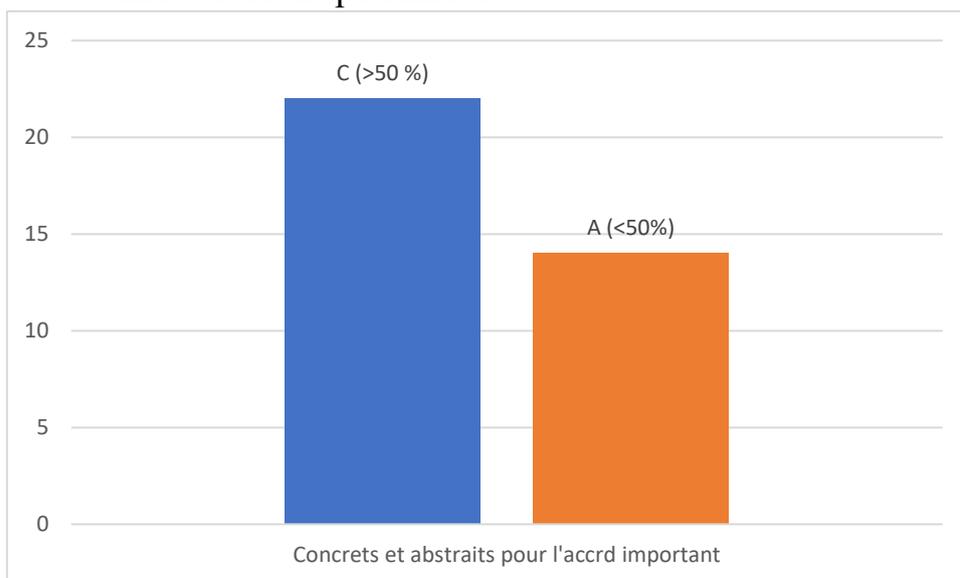
L'expérience ne devrait pas prendre plus d'une demi-heure.

Merci de votre participation !

Annexe 4. Le rapport des mots fréquents et rares à l'accord entre annotateurs.  
Expérience 1



## Annexe 5. Le rapport des mots concrets et abstraits à l'accord entre annotateurs. Expérience 1



## Annexe 6. Stimuli. Expérience 2

ABS	Freq	CONC	Freq	FILLERS	Freq
Mission	85	Station	30	Forêt	35
Rôle	70	Arrêt	51	Barbe	24
Fiction	6	Edition	10	Piscine	23
Conte	13	Revue	10	Coffre	39
Traitement	29	Logement	11	Sable	25
Soin	71	Studio	27	Jambe	113
Jugement	21	Mouvement	40	Chocolat	31
Critique	14	Geste	40	Guitare	13
Usage	14	Bagage	30	Tasse	21
Emploi	30	Valise	50	Singe	35
Avantage	25	Garage	23	Moyenne	35,9
Succès	40	Abri	25		
Patience	31	Licence	8		
Calme	60	Diplôme	17		
Puissance	34	Audience	14		
Pouvoir	117	Public	45		
Longueur	10	Secteur	27		
Distance	31	Zone	53		
Terreur	15	Erreur	124		
Crainte	14	Faute	169		
Moyenne	36,5	Moyenne	40,2		

# Annexe 7. Questionnaire. Expérience 2

## Données anonymisées

Merci de répondre aux questions suivantes

### 1. Vous êtes \*

Une seule réponse possible.

- Femme  
 Homme  
 Autre : \_\_\_\_\_

### 2. Age \*

\_\_\_\_\_

### 3. Niveau d'études

Une seule réponse possible.

- Bac+1  
 Bac+2  
 Bac+3  
 Bac+4  
 Bac+5  
 Bac+6  
 Bac+7  
 Bac+8  
 Autre : \_\_\_\_\_

### 4. Langue(s) maternelle(s) \*

\_\_\_\_\_

### 5. Autres langues et niveaux (A - débutant, B - intermédiaire, C - maîtrise)

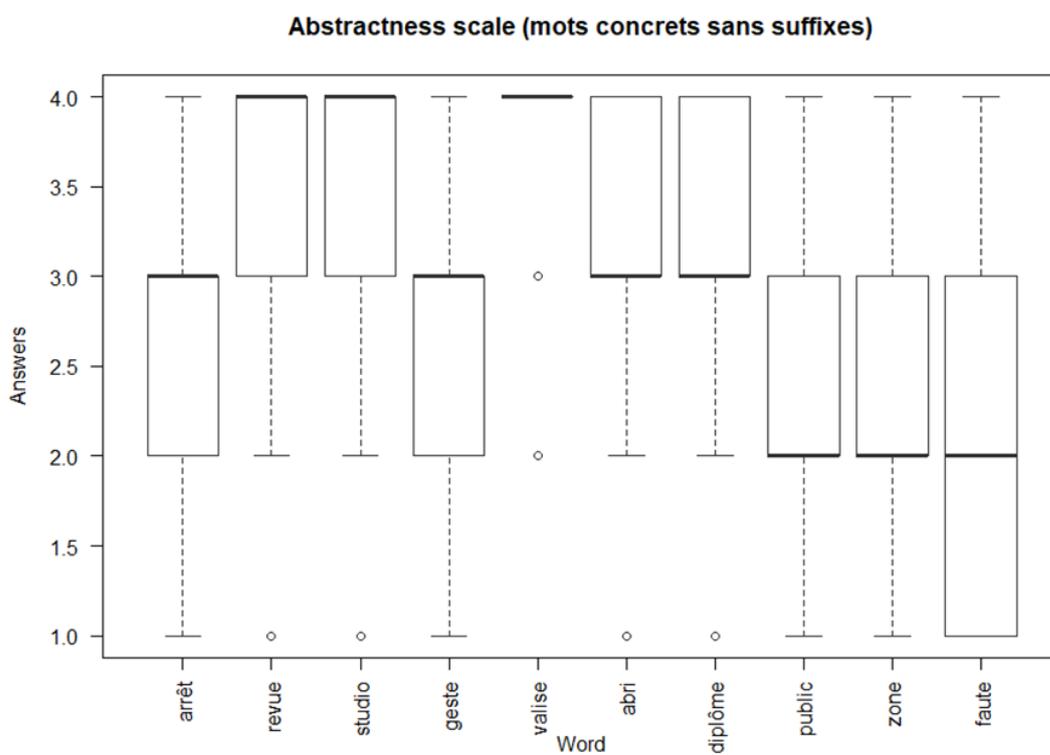
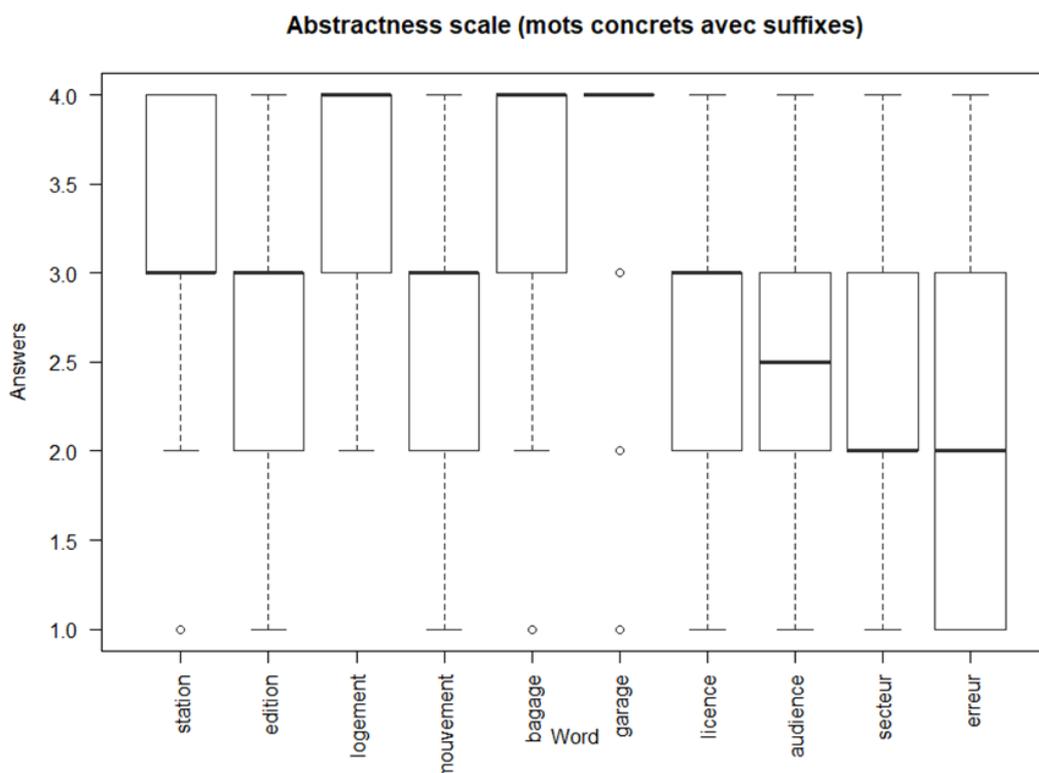
\_\_\_\_\_  
\_\_\_\_\_  
\_\_\_\_\_  
\_\_\_\_\_

### 6. Rentez votre choix \*

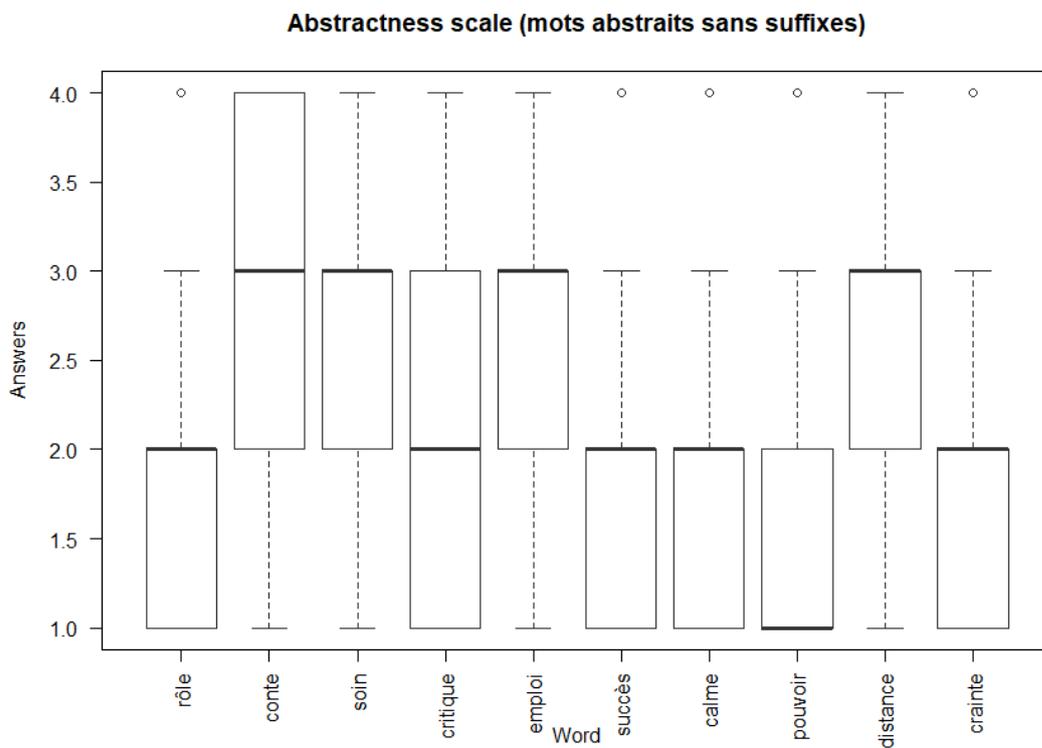
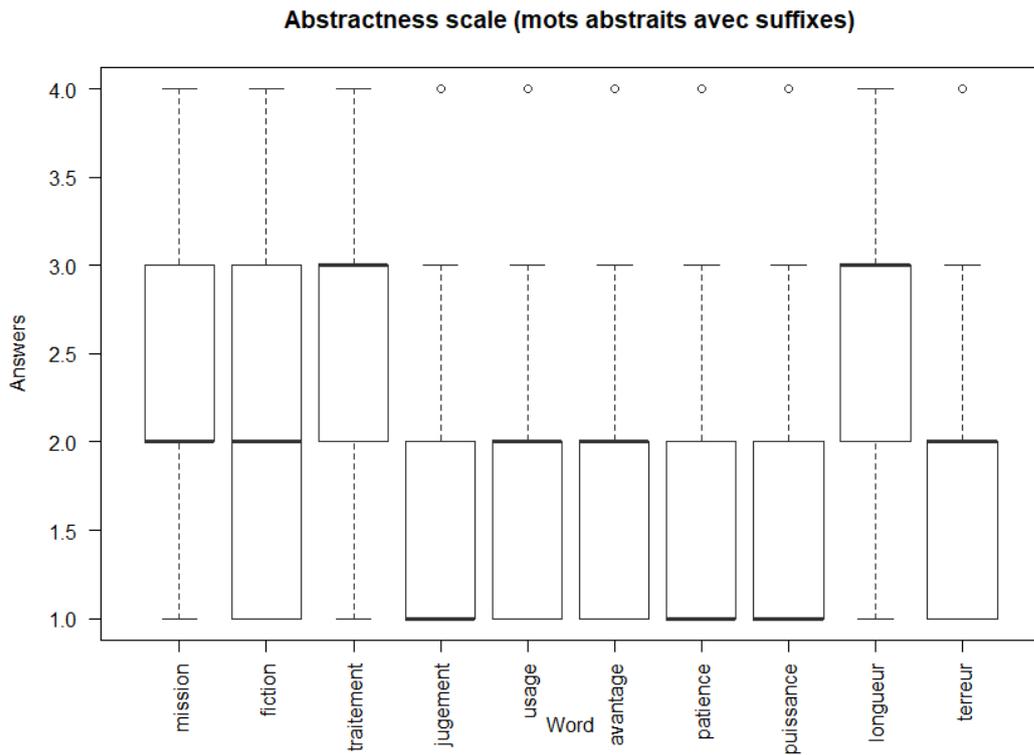
Une seule réponse possible par ligne.

	1 (abstrait)	2 (plutôt abstrait)	3 (plutôt concret)	4 (concret)
Chocolat	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>
Terreur	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>
Coffre	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>
Jugement	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>
Valise	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>
Calmé	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>
Emploi	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>
Abri	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>
Diplôme	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>
Usage	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>
Rôle	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>
Public	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>
Patience	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>
Longueur	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>
Faute	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>
Mouvement	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>
Audience	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>
Avantage	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>
Critique	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>
Bagage	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>
Licence	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>
Distance	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>
Mission	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>
Barbe	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>
Edition	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>
Soin	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>
Forêt	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>
Conte	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>
Guitare	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>
Traitement	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>
Succès	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>
Crainte	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>
Piscine	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>
Logement	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>
Secteur	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>
Singe	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>
Station	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>
Fiction	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>
Sable	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>
Tasse	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>
Erreur	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>

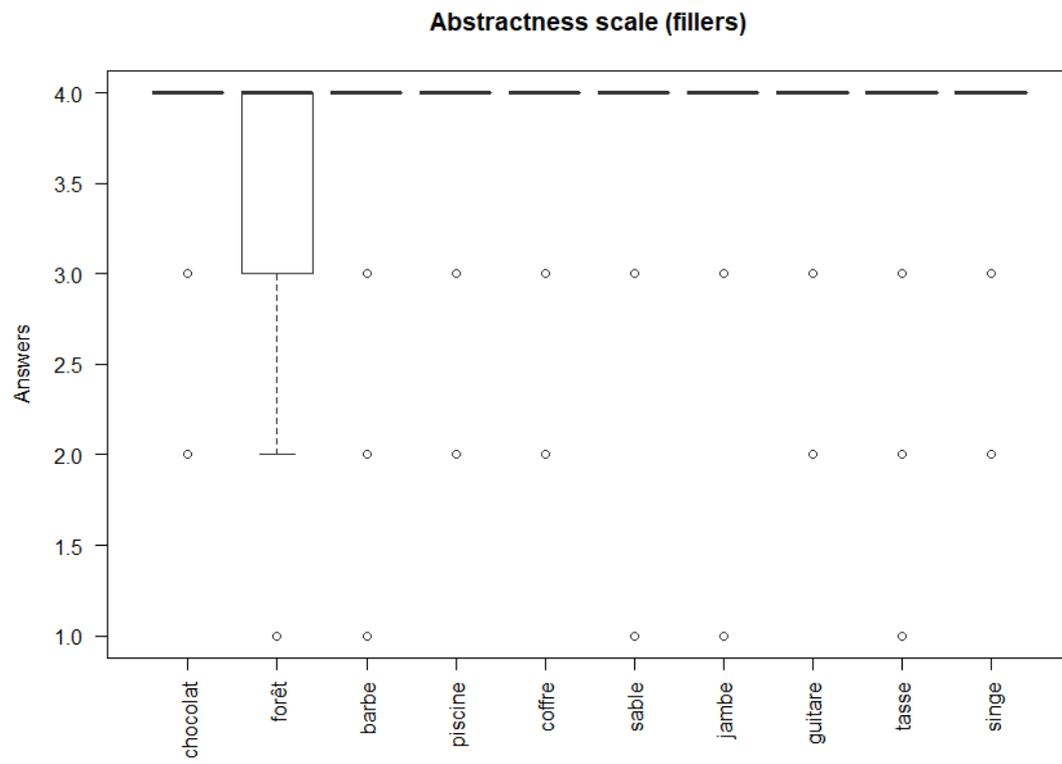
Annexe 8. Graphiques : Mots concrets avec des suffixes et mots concrets sans suffixes. Expérience 2



Annexe 9. Graphiques : Mots abstraits avec des suffixes et mots abstraits sans suffixes. Expérience 2



## Annexe 10. Graphique : Mots très concrets. Expérience 2



Annexe 11. Les moyennes et les écarts types. Expérience 2

MOT	Freq	ABS/CONC	MORPHO	MEAN	SD
Station	30	CONC	1	3,1	0,8
Edition	10	CONC	1	2,55	0,9
Logement	11	CONC	1	3,51	0,61
Mouvement	40	CONC	1	2,61	0,86
Bagage	30	CONC	1	3,44	0,78
Garage	23	CONC	1	3,76	0,61
Licence	8	CONC	1	2,65	1,03
Audience	14	CONC	1	2,61	0,95
Secteur	27	CONC	1	2,39	0,95
Erreur	124	CONC	1	2,04	0,91
Moyenne				2,866	0,84
MOT	Freq	ABS/CONC	MORPHO	MEAN	SD
Arrêt	51	CONC	0	2,71	0,87
Revue	10	CONC	0	3,37	0,84
Studio	27	CONC	0	3,51	0,74
Geste	40	CONC	0	2,86	0,84
Valise	50	CONC	0	3,88	0,4
Abri	25	CONC	0	3,2	0,76
Diplôme	17	CONC	0	3,06	0,81
Public	45	CONC	0	2,47	0,91
Zone	53	CONC	0	2,39	0,83
Faute	169	CONC	0	2,18	0,98
Moyenne				2,963	0,798
MOT	Freq	ABS/CONC	MORPHO	MEAN	SD
Mission	85	ABS	1	2,42	0,95
Fiction	6	ABS	1	2,06	0,93
Traitement	29	ABS	1	2,88	0,76
Jugement	21	ABS	1	1,59	0,78
Usage	14	ABS	1	1,8	0,87
Avantage	25	ABS	1	2,02	0,93
Patience	31	ABS	1	1,47	0,68
Puissance	34	ABS	1	1,67	0,87
Longueur	10	ABS	1	2,53	0,93
Terreur	15	ABS	1	1,65	0,77
Moyenne				2,009	0,847
MOT	Freq	ABS/CONC	MORPHO	MEAN	SD
Rôle	70	ABS	0	1,83	0,76
Conte	13	ABS	0	2,69	1,04
Soin	71	ABS	0	2,69	0,89

Critique	14	ABS	0	2,1	0,97
Emploi	30	ABS	0	2,73	0,93
Succès	40	ABS	0	1,82	0,89
Calme	60	ABS	0	1,69	0,78
Pouvoir	117	ABS	0	1,65	0,9
Distance	31	ABS	0	2,55	0,79
Crainte	14	ABS	0	1,65	0,77
Moyenne				2,14	0,872

## Annexe 12. Listes initiales

<b>19 noms abstraits de la liste initiale (Ferrand, 2001)</b>		<b>42 noms concrets de la liste initiale (Ferrand and Alario, 1998)</b>			
amitié	joie	arbre	chat	livre	poignée
colère	peur	avion	chemise	main	poisson
courage	santé	bateau	cheval	maison	porte
Crainte	sécurité	boîte	chien	manteau	pomme
effort	siècle	bouteille	cigarette	marteau	robe
espoir	succès	bras	église	montagne	sucre
gloire	tristesse	bureau	ferme	montre	table
haine	usage	café	feuille	mur	téléphone
idée	vérité	camion	fleur	oiseau	train
imagination		carte	journal	pain	voiture
		chaîne	lettre		

### Annexe 13. Exemples des données filtrées.

La 'relation' est la méthode par laquelle un mot a été obtenu : les voisins distributionnels (NN) ou la cooccurrences syntaxiques (SC). La catégorie correspond aux noms concrets (C) et abstraits (A), on note avec \* les erreurs de l'annotation automatique.

<b>Id Stimulus</b>	<b>Stimulus</b>	<b>Id Output</b>	<b>Output</b>	<b>Relation</b>	<b>Category</b>
1	aéroport	1	port	NN	C
1	aéroport	2	gare	NN	C
1	aéroport	3	parc	NN	C
1	aéroport	4	station	NN	C
1	aéroport	5	tarmac	SC	C
1	aéroport	6	atterrissage*	SC	C
1	aéroport	7	ravitaillement*	SC	C
1	aéroport	8	airbus	SC	C
2	ballon	9	balle	NN	C
2	ballon	10	objet	NN	C
2	ballon	11	vélo	NN	C
2	ballon	12	cassette	NN	C
2	ballon	13	nacelle	SC	C
2	ballon	14	manieur	SC	C
2	ballon	15	tour	SC	C
2	ballon	16	tentative*	SC	C
3	câble	17	téléphone	NN	C
3	câble	18	bouquet	NN	C
3	câble	19	télécommunication*	NN	C
3	câble	20	satellite	NN	C
3	câble	21	abonné*	SC	C
3	câble	22	gaine	SC	C
3	câble	23	abonnement	SC	C
3	câble	24	raccordement	SC	C
4	dessin	25	photo	NN	C
4	dessin	26	photographie	NN	C
4	dessin	27	peinture	NN	C
4	dessin	28	portrait	NN	C
4	dessin	29	ensemble	SC	C
4	dessin	30	dossier	SC	C
4	dessin	31	carton	SC	C
4	dessin	32	accompagné*	SC	C
5	enveloppe	33	dotation*	NN	C
5	enveloppe	34	subvention*	NN	C
5	enveloppe	35	prime*	NN	C

5	enveloppe	36	indemnité*	NN	C
5	enveloppe	37	protéine	SC	C
5	enveloppe	38	dos	SC	C
5	enveloppe	39	dépassement*	SC	C
5	enveloppe	40	déblocage*	SC	C
6	abus	41	recel	NN	A
6	abus	42	détournement	NN	A
6	abus	43	escroquerie	NN	A
6	abus	44	fraude	NN	A
6	abus	45	information	SC	A
6	abus	46	complicité	SC	A
6	abus	47	rencontre	SC	A
6	abus	48	juge*	SC	A
7	beauté	49	plaisir*	NN	A
7	beauté	50	richesse	NN	A
7	beauté	51	charme	NN	A
7	beauté	52	goût	NN	A
7	beauté	53	crème*	SC	A
7	beauté	54	évidence	SC	A
7	beauté	55	canon	SC	A
7	beauté	56	grain*	SC	A
8	chance	57	possibilité	NN	A
8	chance	58	capacité	NN	A
8	chance	59	avantage	NN	A
8	chance	60	potentiel	NN	A
8	chance	61	scepticisme	SC	A
8	chance	62	égalité	SC	A
8	chance	63	égalisation	SC	A
8	chance	64	illusion	SC	A
9	décision	65	choix	NN	A
9	décision	66	mesure	NN	A
9	décision	67	accord	NN	A
9	décision	68	déclaration	NN	A
9	décision	69	félicité	SC	A
9	décision	70	cassation	SC	A
9	décision	71	pourvoi	SC	A
9	décision	72	réaction	SC	A
9	émotion	73	inquiétude	NN	A
10	émotion	74	angoisse	NN	A
10	émotion	75	sentiment	NN	A
10	émotion	76	plaisir	NN	A
10	émotion	77	capteur*	SC	A

10	émotion	78	chantage	SC	A
10	émotion	79	larme*	SC	A
10	émotion	80	moment	SC	A

#### Annexe 14. Stimuli pour l'expérience en ligne.

rouble	affectation	civilisation
bourdonnement	mairie	cosmétique
ostracisme	généralisation	surprise
maman	crispation	échelle
service	milieu	itinéraire
information	soupir	possibilité
longueur	gardien	lancement
consommation	photographe	nappe
suspect	trésorerie	déjeuner
voie	catalogue	maison
pétale	pirate	athlète
stress	terminologie	orge
relais	trouble	amidon
rajeunissement	codirecteur	approfondissement
nation	modèle	fusion
parlement	socle	privatisation
victoire	réélection	purée
écho	fortune	perplexité
agriculture	ignorance	pincée
ingénierie	signe	manie
gare	progrès	serveur
compatibilité	hostilité	demandeur
sentiment	présentation	assurance
simplicité	réorientation	sûreté
échéance	zone	probabilité
secrétariat	diffusion	détresse
démence	voyage	cadre
jazz	exception	contrôleur
réserve	patronat	extrême
rassemblement	secours	estimation