

# Induction de lexiques bilingues

à partir de corpus comparables et parallèles.

---

Laurent Jakubina

14 Décembre 2017

Université de Montréal - Département d'informatique et de recherche opérationnelle - RALI

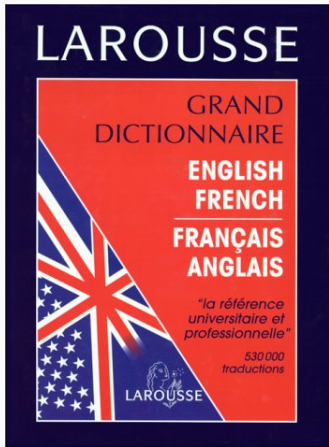
1. Introduction - Induction de lexiques Bilingues...
2. ...à partir de Corpus Parallèles  
[Article n°3 \(Jakubina and Langlais, 2016a\)](#)
3. ...à partir de Corpus Comparables  
via alignement de mots en corpus comparables  
[Article n°1 \(Jakubina and Langlais, 2015\)](#)  
comme tâche d'évaluation des embeddings interlingues  
[Article n°2 \(Jakubina and Langlais, 2016b\)](#)  
via reclassement supervisé de candidats à la traduction  
[Article n°4 \(Jakubina and Langlais, 2017\)](#)
4. Travaux futurs & Conclusion

## Introduction - Induction de lexiques Bilingues...

---

Lexique Bilingue (LB) ? (syn: dictionnaire de traduction)

Lexique Bilingue (LB) ? (syn: dictionnaire de traduction)



## Lexique Bilingue (LB) ? (syn: dictionnaire de traduction)

### 2 • *abstergent*

---

**abstergent** adj. n. (chim.), détergent m.

**abstract** n., résumé m.

**abundance** n., abondance f., teneur (phys. nucl.), - **ratio** (phys. nucl.), rapport de teneurs

**abyss** n. (géol., océano.), abîme m., gouffre m.

**abyssal** adj. (océano.), abyssal

**ac (alternating current)** (électr.), courant alternatif

**ac vector representation** n. (électr.), représentation f. de Fresnel (*du vecteur tournant*)

**Acarina** n. pl. (arthr.), acariens m.

**acaudate(d)** adj. (anat., zoo.), acaudé

**acaulescent** adj. (bot.), acaule

**acaulous - acaulescent**

**accelerating potential** n. (phys.), potentiel m. accélérateur

**accrete (to)** v. t. (géol.), s'accroître par accrétion

**accretion** n. (géol.), accrétion f.

**accumulation** n., amas m.

**accumulator** n. (électr., instr.), accumulateur m. (*électrique*)

**accuracy** adj., n., exactitude n. f., précision n. f., - **class** (instr.), classe de précision (*d'un instrument*)

**accurate** adj., exact, précis

**accurately** adv., avec exactitude

**acentric(al)** adj. (math.), acentré

**acephalous** adj. (path.), acéphale

**acetabulum** n. (anat., zoo.), acétabule f., acétabulum m.

**acetal** n. (chim.), acétal m. (*diéther*)

**acetaldehyde** n. (chim.), acétaldéhyde m.

**acetaldol** n. (chim.), acétaldol m.

**acetamide** n. (chim.), acétamide m.

Induction de Lexiques Bilingues (**ILB**) ?

---

## Induction de Lexiques Bilingues (**ILB**) ?

---

Ressources (textuelles) Bilingues/Plurilingues



---

...		...	
patchouli	patchouli		
stickler	pointilleuse	rigoriste	
nightly	nocturne	nuitamment	
gel	gel	gélifier	prendre
astounding	abasourdissant	effarant	renversant
...		...	

---



## Induction de Lexiques Bilingues (ILB) ?

---

...	...
patchouli	?
stickler	?
nightly	?
gel	?
astounding	?
...	...

---

+

Bilingues/Plurilingues  
Ressources (textuelles)



---

...
patchouli
pointilleuse
nuitamment
prendre
effarant
...

---

**Corpus** (: ensemble de documents, regroupés selon une perspective précise) **disponible en deux langues** (source et cible)

3 Types, selon le niveau de comparabilité source-cible  
(Sharoff et al., 2013)



**Corpus** (: ensemble de documents, regroupés selon une perspective précise) **disponible en deux langues** (source et cible)

3 Types, selon le niveau de comparabilité source-cible  
(Sharoff et al., 2013)

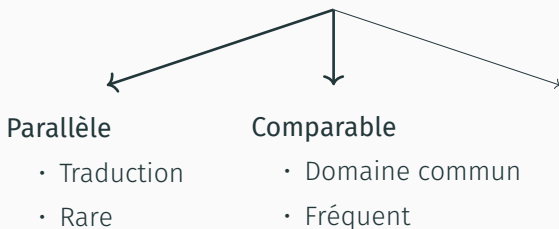


## Parallèle

- Traduction
- Rare

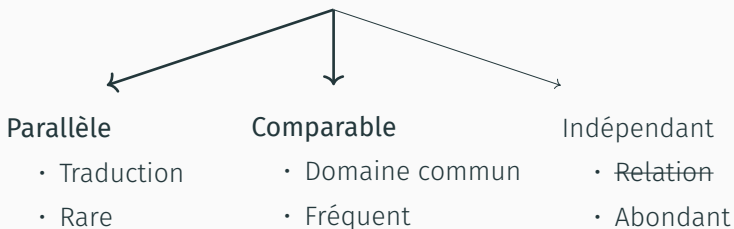
**Corpus** (: ensemble de documents, regroupés selon une perspective précise) **disponible en deux langues** (source et cible)

3 Types, selon le niveau de comparabilité source-cible  
(Sharoff et al., 2013)



**Corpus** (: ensemble de documents, regroupés selon une perspective précise) **disponible en deux langues** (source et cible)

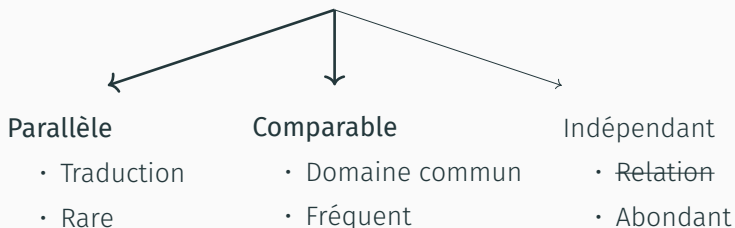
3 Types, selon le niveau de comparabilité source-cible  
(Sharoff et al., 2013)



# Ressources Bilingues (/ Plurilingues)

**Corpus** (: ensemble de documents, regroupés selon une perspective précise) **disponible en deux langues** (source et cible)

3 Types, selon le niveau de comparabilité source-cible  
(Sharoff et al., 2013)



⇒ selon le type, la phase d'induction (⇒) de l'**ILB** diffère...

# Corpus Parallèles : Exemple

Total household expenditures for this category reached P215.6 billion (B) in 2000 compared to only P125.6B in 1994. Among the major contributors are flour and noodles. NSO-FIES Total household expenditure on flour reached P1.04B in 2000 per NSO Family Income and Expenditures Survey (FIES). Bread consumption increased by 28% from P18B in 1997 to P23B in 2000. Spending on biscuits totaled P4.7B in 2000 up by P400 million (M) from 1997. Spending on noodles (categorized into bihon, miki, miswa, sotanghon, macaroni, etc.) reached P3.7B in 2000. Instant noodles, on the other hand, are valued at P6B in 2000, with growth estimated at 16% p.a. Household spending for flour is very small.

Les dépenses des ménages pour cette catégorie ont atteint 215,6 milliards de pesos en 2000, comparés à seulement 125,6 milliards de pesos en 1994. Parmi les principaux contributeurs, l'on retrouve la farine et les nouilles. NSO-FIES Selon l'enquête du Bureau national des statistiques portant sur le revenu familial et les dépenses (FIES), les dépenses totales des ménages consacrées à la farine ont atteint 1,04 milliard de pesos en 2000. La consommation de pain a augmenté de 28%, de 18 milliards de pesos en 1997 à 23 milliards en 2000. Les dépenses sur les biscottes ont totalisé 4,7 milliards de pesos en 2000, ayant connu une augmentation de 400 millions de pesos depuis 1997. Les dépenses au chapitre des nouilles (classées sous vermicelle chinois, nouilles de riz miki, de blé, asiatiques, macaroni, etc.) ont atteint 3,7 milliards de pesos en 2000. D'autre part, les dépenses sur les nouilles instantanées sont évaluées à 6 milliards de pesos en 2000, avec une croissance estimée de 16% p.a. Les dépenses des ménages concernant la farine sont vraiment minimes.

# Corpus Comparables : Exemple

Not logged in | Talk | Contributions | Create account | Log in

Article | Talk | Read | Edit | View history | Search

## Cafeteria

From Wikipedia, the free encyclopedia

*This article is about the food service location. For other uses, see Cafeteria (disambiguation).*

A **cafeteria** is a type of food service location in which there is little or no waiting staff table service, whether a restaurant or within an institution such as a large office building or school; a school dining location is also referred to as a **dining hall** or **canteen** (in British English).<sup>[1]</sup> Cafeterias are different from coffeehouses, despite being the Spanish translation of the English term.



A corporate cafeteria in Bangalore, India, December 2003.



Non connecté | Discussion | Contributions | Créer un compte

Se connecter

Article | Discussion | Lire | Modifier | Plus | Rechercher

## Caf  teria

*Cet article concerne lieu de restauration. Pour le genre d'aigle, voir Cafeteria.*

**Cet article est une   bauche concernant la restauration.**

Vous pouvez partager vos connaissances en l'am  liorant (**comment ?**) selon les recommandations des projets correspondants.


Une **caf  teria** est un lieu de restauration o   il n'y a pas (ou tr  s peu) de service    table. Le consommateur se sert g  n  ralement comme dans un libre-service,    l'aide de plateaux individuels.



Caf  teria dans une   cole de Montr  al.



# Corpus Comparables : Exemple



WIKIPÉDIA  
The Free Encyclopedia

Not logged in | Talk | Contributions | Create account | Log in


Article | Talk | Read | Edit | View history | Search

## Cafeteria


From Wikipedia, the free encyclopedia


*This article is about the food service location. For other uses, see Cafeteria (disambiguation).*

A **cafeteria** is a type of food service location in which there is little or no waiting staff table service whether a restaurant or within an institution such as a large office building or school a school dining location is also referred to as a **dining hall** or **canteen** (in British English).<sup>[1]</sup> Cafeterias are different from coffeehouses, despite being the Spanish translation of the English term.



A corporate cafeteria in Bangalore, India, December 2008.





WIKIPÉDIA  
L'encyclopédie libre

Non connecté | Discussion | Contributions | Créer un compte | Se connecter

Article | Discussion | Lire | Modifier | Plus | Rechercher


## Cafétéria

*Cet article concerne lieu de restauration. Pour le genre d'aigue, voir Cafeteria.*

**Cet article est une ébauche concernant la restauration.**

Vous pouvez partager vos connaissances en l’améliorant (comment ?) selon les recommandations des projets correspondants.

Une **cafétéria** est un lieu de restauration où il n'y a pas (ou très peu) de service à table. Le consommateur se sert généralement comme dans un libre-service, à l'aide de **plateaux** individuels.



Cafétéria dans une école de Montréal.

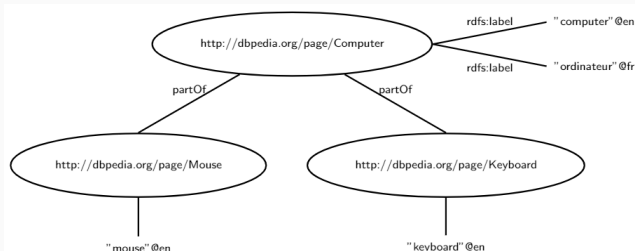
Accueil | Portails thématiques | Article au hasard | Contact

Contribuer | Débuter sur Wikipédia | Aide | Communauté | Modifications récentes | Faire un don

Imprimer / exporter | Créer un livre | Télécharger comme PDF | Version imprimable

Outils | Pages liées | Suivi des pages liées | Importer un fichier

Objectif: utilisation de l'**ILB** pour le Web Sémantique et **DBPEDIA**...



**rdfs:label** : Chaîne de caractère "human-readable"

⇒ majoritairement en anglais (Gómez-Pérez et al., 2013)

⇒ **En Français ?** Seulement 1 concept sur 5 (20%)<sup>1</sup>

<sup>1</sup>[wiki.dbpedia.org/Datasets/DatasetStatistics](http://wiki.dbpedia.org/Datasets/DatasetStatistics)

# Contexte de recherche prédoc

**Solution (possible)** : "Induire" un label français pour les concepts de DBpedia, par identification d'une traduction au label anglais

- Concepts de **DBPEDIA** = Titres d'articles de **WIKIPÉDIA** anglais
- Pas de label français = pas de lien inter-langues dans **WIKIPÉDIA**

<i>n</i> -gram	Distribution [%]
1-gram	158 901 [30.44]
2-gram	185 975 [35.63]
3 – 4-gram	141 872 [27.16]
6 – 9-gram	26 855 [5.14]
10 – 28-gram	8 292 [1.58]

Distribution du nombre de mots qui composent les titres de Wikipedia

Corpus Parallèles / Corpus Comparables

...à partir de Corpus Parallèles

---

**ILB** via les modèles IBM d'alignement de mots (Traduction Automatique Statistique) (Brown et al., 1990)

Exemple d'application en production: TRANSEARCH<sup>2</sup>

---

<sup>2</sup>Développé au RALI.

# TransSearch: Identificateur de traductions (en CP)

The screenshot shows the TransSearch interface with the following elements:

- Logo: TRANSSEARCH H3 BETA
- Logo: TERMINOTIX
- Logo: rali
- Navigation: REQUÊTES, MON COMPTE, PRÉFÉRENCES, AIDE, QUITTER
- User: UTILISATEUR : lapalme
- Search bar: Signet / Favori personnalisé: TransSearch (qu'est-ce que c'est ?)
- Collection: Collection de documents : Les Hansards canadiens
- Expression: Expression : take+ .. ride
- Buttons: Chercher, Requête bilingue
- Results header: 92 traductions de take+ .. ride dans 106 occurrences
- Table of results with columns for source text, frequency, and target text.

Source Text	Frequency	Target Text	Frequency
dindons de la farce	4	dindons de la farce	4
monté un bateau	3		
faire avoir	3	Emissions continue to rise and taxpayers are being <b>taken along for the ride</b> .	
se fasse rouler	2		
fait berné	2	They are left with nothing. Now they are here illegally with no documentation. Canadians are being <b>taken for a ride</b> .	
se fait jouer	2		
moqués de	2		
fait	2		
les a	2	This would affect close to 400,000 Canadians, 80,000 of them Quebecers, who have been the ones <b>taken for a ride</b> .	
se sont fait avoir	2		
le public pour attirer la	1		
a fait une balade	1		
nous rouler dans ce projet	1	I think that this is a prime example of a tainted system in which people who cannot afford to invest in sectors eligible for tax credits are urged to do so through all kinds of scams and end up being <b>taken for a ride</b> .	
nous tous	1		
en train de monter un bateau à la population canadienne	1		
tête des contribuables que se paie le	1		
passer une petite vite	1		
bourrer de l'autre côté de la chambre en	1		
ont pris la voiture que pour faire une balade	1		

Target text details from the screenshot:

- Les émissions continuent d'augmenter et c'est le contribuable qui est **le dindon de la farce**.
- Ces personnes se trouvent ici illégalement, elles n'ont aucun document et nous, les Canadiens, sommes les **dindons de la farce**.
- Il s'agit d'une mesure qui toucherait près de 400 000 Canadiens, dont 80 000 Québécois, qui ont été les **dindons de la farce**.
- Je pense que c'est un exemple patent d'un système vicié, où des gens qui n'ont pas les moyens d'investir dans des domaines où on peut obtenir des crédits d'impôt se voient, par toutes sortes de subterfuges, invités à le faire et, en bout de ligne, ils se trouvent à être **les dindons de la farce**.

(Bourdaillet et al., 2010) : Précision de 78.3% (Transpot)

⇒ Évaluer la couverture du vocabulaire de nos ressources parallèles

# Couverture de nos ressources parallèles - Référence

- = Liste de paires de titres (anglais; français) de **WIKIPÉDIA**
  - En relation de traduction = relié par lien inter-langue
  - Taille après filtres Entités nommées (et ponctuations)  $\approx$  200.000

## Exemples

Greeting	Salutation
Activities of daily living	Acte de la vie quotidienne
Boston Port Act	Acte du port de Boston
Layout Versus Schematic	Logiciel de vérification de schéma
Output gap	Écart de production
South Sulawesi languages	Langues sulawesi du Sud

## Couverture de nos ressources parallèles - Données

Corpus	#Paires <sup>Phrases</sup>	#Mots <sup>en</sup>	#Mots <sup>fr</sup>
	(en milliers)		
HEALTH	809	133	142
EUROPARL	2 007	132	142
HANSARDS	8 802	242	253
GIGAWORD	22 520	3 002	2 775



# Couverture de nos ressources parallèles - Approche

GIGAWORD<sup>en</sup>

GIGAWORD<sup>fr</sup>

21677278: *French and English translation cards [...]*

*greeting*

21677278: *Des fiches de salutation en anglais et en français [...]*

*salutation*

Exemple de phrases parallèles dans le corpus GIGAWORD respectant la **contrainte d'alignement** pour le terme *greeting* et sa traduction *salutation*.

## Couverture de nos ressources parallèles - Résultats

	HANSARDS	EUROPARL	HEALTH	GIGAWORD
<i>n</i>	%	%	%	%
1-4	1.1	0.9	1.2	4.0
5-10	0.4	0.3	0.4	1.5
11-100	0.7	0.5	0.6	2.9
>100	0.6	0.3	0.2	2.0
0	97.4	98.0	97.6	89.5

# Couverture de nos ressources parallèles - Résultats

	HANSARDS	EUROPARL	HEALTH	GIGAWORD
<i>n</i>	%	%	%	%
1-4	1.1	0.9	1.2	4.0
5-10	0.4	0.3	0.4	1.5
11-100	0.7	0.5	0.6	2.9
>100	0.6	0.3	0.2	2.0
0	97.4	98.0	97.6	89.5

## Couverture de nos ressources parallèles - Résultats

	HANSARDS	EUROPARL	HEALTH	GIGAWORD
<i>n</i>	%	%	%	%
1-4	1.1	0.9	1.2	4.0
5-10	0.4	0.3	0.4	1.5
11-100	0.7	0.5	0.6	2.9
>100	0.6	0.3	0.2	2.0
0	97.4	98.0	97.6	89.5

⇒ Couverture insuffisante de nos corpus parallèles pour l'**ILB**

pistes de solution

- **ILB** sur Corpus Comparables : suite de la présentation
- Plus de données parallèles :

## Article n°3

Jakubina, L. and Langlais, P. (2016a). Bad luc@wmt 2016: a bilingual document alignment platform based on lucene.

*In Proceedings of the First Conference on Machine Translation*, pages 703–709, Berlin, Germany. Association for Computational Linguistics

...à partir de Corpus Comparables

---

via l'alignement de mots en Corpus Comparables

ACL 1995 : (Rapp, 1995) et (Fung, 1995)

## Hypothèse émise par (Rapp, 1995)

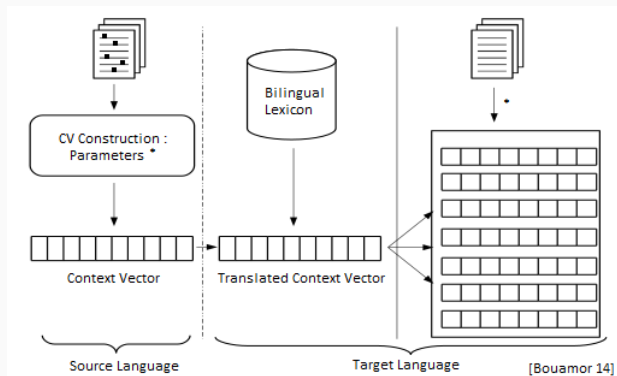
*Si deux mots (A et B) cooccurrent plus souvent que par le hasard dans une langue source, alors leur traduction respective (A' et B') doivent cooccurrer plus souvent que par le hasard dans la langue cible.*

⇒ Les mots en relations de traductions partagent des motifs de cooccurrence similaires

- Réalise une expérience de simulation pour le prouver
- Approche impraticable point de vue computationnel



# (Rapp, 1999) : Approche Standard / RAPP



## Étapes de l'Approche Standard

A attiré l'attention d'un grand nombre de chercheurs et de ce fait, a bénéficié de beaucoup d'améliorations (Sharoff et al., 2013)

Échantillonnage: espace de construction/recherche/voc réduit

## Exemples

- (Vulić et al., 2011): "...only lemmatized noun forms ; voc  $\approx$  8.500"
- (Morin and Prochasson, 2011): "...words occurring less than twice discarded..."
- (Gaussier et al., 2004): "...frequency less than 5 discarded ; voc  $\approx$  40.000"

⇒ Quid d'une "mise en production", à l'ère du **Big Data** ?

### Article n°1

Jakubina, L. and Langlais, P. (2015). Projective methods for mining missing translations in dbpedia.

*In Proceedings of the Eighth Workshop on Building and Using Comparable Corpora*, pages 23–31, Beijing, China. Association for Computational Linguistics

## Approches

- Reproduction de **RAPP**
  - Implémentée pour le **non**-échantillonnage
  - Espace de recherche = 2.942.067 mots cibles (français)
- Reproduction de (Bouamor et al., 2013)
  - Analyse Sémantique Explicite (Gabrilovich and Markovitch, 2009)
  - Supérieure à **RAPP**

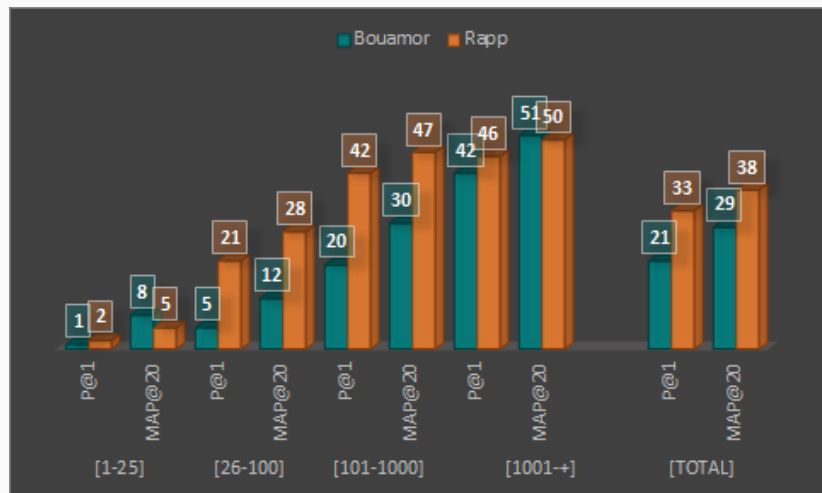
## Données

- **WIKIPÉDIA**<sup>3</sup> Anglais et Français 2013 non-échantillonné
- Liste de référence basée sur les titres 1-gramme de **WIKIPÉDIA**
  - Gammes de fréquences (très rare, rare, fréquent, très fréquent)

---

<sup>3</sup>DBpedia en est le silo RDF.

# (Jakubina and Langlais, 2015) - Résultats



- Biais en faveur des mots fréquents
- Performances **RAPP**  $\geq$  (Bouamor et al., 2013)
  - Protocoles différents: corpus spécialisés, pré-traitements, liste de référence ( $\approx 100$ ), ...
- Implémentation multi-threadée (8-16) avec **LUCENE** et similarité vectorielle par morceaux
  - mise-à-mal par les mots hyper-fréquents (ex: France, freq=1.000.000)
  - $\Rightarrow$  échantillonnage toujours nécessaire
- Types d'erreurs (effet thésaurus)

# (Jakubina and Langlais, 2015) - Exemples de sortie

## RAPP

myringotomy [1-25]	permette	devra	scopie	nécessitait	...
syllabication [26-100]	modifier	suffit	vouloir	<u>syllabique</u>	...
numerology [101-1000]	numérologie	<u>occultisme</u>	ésotérisme	divinatoire	...
entertainment [1001+]	beatmakers	manglobe	spycraft	déduplication	...

## (Bouamor et al., 2013)

myringotomy	<u>laryngologie</u>	oto	rhino	traitement	otite
syllabication	<u>langues</u>	consonne	langue	lettre	phonétique
numerology	oeuvre	<u>gematria</u>	angels	<u>nombres</u>	chire
entertainment	entertainment	divertissement	vidéo	sony	jeu

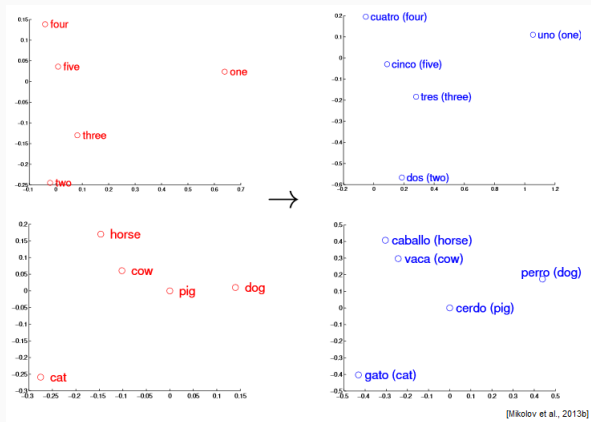
Idée: Combiner les listes de candidats<sup>4</sup> ?

<sup>4</sup>Potentielle traduction pour 431 termes sur 528 (81.6%) [100+...]

Comme tâche d'évaluation des embeddings interlingues

# MIKO (Mikolov et al., 2013b) : ILB comme tâche d'évaluation

*Nous observons des similarités dans la projection géométrique des embeddings de mots provenant de différentes langues.*





# MIKO (Mikolov et al., 2013b) : Observation(s)

- Pas de positionnement vis-à-vis de **RAPP**
- Quid des **Mots Rares** ?
- Échantillonnage
  - "...vocabularies consist of the words that occurred at least five times in the corpus." + "...short phrases..."
  - Espace de recherche de **MIKO**  $\approx$  110.000 mots

## Article n°2

Jakubina, L. and Langlais, P. (2016b). A comparison of methods for identifying the translation of words in a comparable corpus: Recipes and limits.

*CICLing: International Conference on Computational Linguistics and Intelligent Text Processing*

## Approches

- **RAPP** (50 variantes<sup>5</sup>)
- (Prochasson and Fung, 2011) (2 variantes)
  - Proposé pour gérer les **Mots Rares**
- **MIKO** (Mikolov et al., 2013b) (130 variantes)
  - Construction avec **WORD2VEC** (Mikolov et al., 2013a)
  - Projection (Dinu and Baroni, 2014)

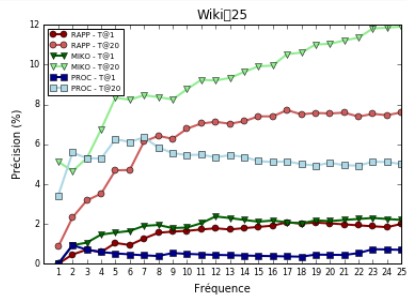
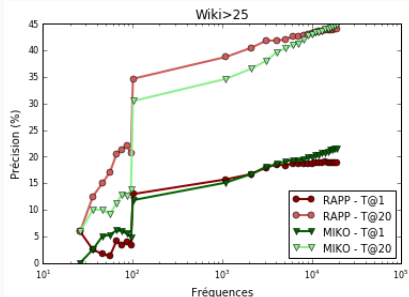
## Données

- **WIKIPÉDIA** anglais et français 2013 non-échantillonné
- 2 Listes de référence:  $Wiki_{>25}$  et  $Wiki_{\leq 25}$  [titres  $\Rightarrow$  mots]
  - 6.8 millions des mots (92%) de **WIKIPÉDIA** anglais apparaissent moins de 26 fois

---

<sup>5</sup>selon les valeurs des méta-paramètres ; Syn: signaux.

# (Jakubina and Langlais, 2016b) - Résultats



## RAPP

- Taille Fenêtre: 3
- Normalisation: PMI

## MIKO

- Modèle CBOW
- Taille Fenêtre: 5

## RAPP

- Taille Fenêtre: 15
- Normalisation: ORD

## MIKO

- Modèle Skip-Gram
- Taille Fenêtre: 10

## (Jakubina and Langlais, 2016b) - Impacts

- Performances **RAPP**  $\approx$  Performances **MIKO** (biais de fréquence)
  - (Levy et al., 2015) prouve que les modèles sont proches
- "Recettes" : différents jeux de méta-paramètres selon la fréquence du terme à traiter
  - Meilleures performances sur les **Mots Rares**
- (Prochasson and Fung, 2011) : Résultats décevants
  - Problème des **Mots Rares** reste ouvert
- Complémentarité des approches validée (combinaison-oracle)
  - Performances doublées sur les **Mots Rares**

## Candidats corrects mais pas dans notre liste de référence

### RAPP

donut (beigne)	<u>aromatisé</u>	<u>donut</u>	<u>beignet</u>
brilliantly (brillamment)	<u>imaginatif</u>	<u>captivant</u>	<u>rusé</u>
gentle (doux)	<u>enjoué</u>	<u>serviable</u>	affable
pathologically (pathologiquement)	<u>cordonale</u>	pathologique	<u>diagnostiqué</u>

### MIKO

donut (beigne)	liper	babalous	savonette
brilliantly (brillamment)	<u>éclatant</u>	<u>pathétique</u>	émouvant
gentle (doux)	<u>colérique</u>	<u>enjoué</u>	espiègle
pathologically (pathologiquement)	<u>psychosexuel</u>	<u>psychoactif</u>	piloerection

**Idée:** Approche basée sur la combinaison (pour les **Mots Rares**) ?

via reclassement supervisé de candidats à la traduction

# RERANKER : Reclassement supervisé après Combinaison

## Weighting [EN] : pondération [FR]

<u>Rapp</u>		<u>Miko</u>		<u>Faru</u>	
maximise	0.069	<b>pondération</b>	<b>0.753</b>	<b>pondération</b>	<b>0.428</b>
quantifié	0.066	<i>pondéré</i>	0.688	calculé	0.389
adaptatif	0.065	dépendra	0.648	quantification	0.387
<b>pondération</b>	<b>0.066</b>	calculé	0.622	luminance	0.336
<i>pondéré</i>	<i>0.063</i>	maximise	0.62	<i>pondéré</i>	<i>0.319</i>
calculé	0.054	quantifié	0.58	quantifié	0.307
dépendra	0.054	luminance	0.567	maximise	0.268
liquidité	0.052	dénominateur	0.55	inférieure	0.267

Reclassement (avec features) après Combinaison

Résultats : **pondération (0.020)** ; maximise (0,0027) ;  
*pondéré (0,0021)* ; quantifié (0,0019) ; calculé (0,0010) ; ...

- Axe de recherche prometteur mais peu exploré (Sharoff et al., 2013; Irvine and Callison-Burch, 2013)
- **Objectif** : Améliorer l'**ILB** sur les **Mots Rares**

## Article n°4

Jakubina, L. and Langlais, P. (2017). Reranking translation candidates produced by several bilingual word similarity sources.

*In Proceedings of the 15th Conference of the European Chapter of the Association for Computational Linguistics: Volume 2, Short Papers*, pages 605–611, Valencia, Spain. Association for Computational Linguistics



## Approche

- **RERANKER** avec la librairie LTR **RANKLIB**<sup>6</sup>
  - 3 ensembles de features: rang (2), string (5), fréquence (4)
  - 700 entrées pour s'entraîner ; 300 pour tester

## Données

- 3 Listes d'**ILB**: **RAPP**, **MIKO** et **FARU** (Faruqui and Dyer, 2014)
  - Listes adaptés selon la fréquence des mots sources<sup>7</sup>
  - 100 candidats par liste
- 3 Listes de référence:  $Wiki_{>25}$ ,  $Wiki_{\leq 25}$  et  $Euro_{5-6k}$ <sup>8</sup>

<sup>6</sup><https://sourceforge.net/p/lemur/wiki/RankLib/>

<sup>7</sup>"Recettes" de l'article n°2.

<sup>8</sup>EUROPARL : Débats parlementaires européens.

# (Jakubina and Langlais, 2017) - Résultats

	INDIVIDUEL			RECLASSÉ				COMBINÉ & RECLASSÉ		
	@1	@5	@20	@1	@5	@20		@1	@5	@20
<i>Wiki</i> <sub>&gt;25</sub>										oracle: 69.3
<b>RAPP</b>	20.0	33.0	43.0	36.3 <sup>2.5</sup>	48.8 <sup>1.9</sup>	53.8 <sup>1.9</sup>				
<b>MIKO</b>	17.0	32.6	41.6	38.1 <sup>1.9</sup>	49.0 <sup>1.5</sup>	54.3 <sup>1.3</sup>	R+M+F	<b>45.6<sup>2.2</sup></b>	<b>59.6<sup>1.1</sup></b>	<b>64.0<sup>1.8</sup></b>
<b>FARU</b>	13.3	26.0	33.3	34.3 <sup>1.5</sup>	44.0 <sup>2.6</sup>	47.9 <sup>2.1</sup>				
<i>Wiki</i> <sub>≤25</sub>										oracle: 28.6
<b>RAPP</b>	2.6	4.3	7.3	8.6 <sup>1.2</sup>	9.4 <sup>0.8</sup>	10.2 <sup>1.0</sup>				
<b>MIKO</b>	1.6	4.6	10.6	16.6 <sup>2.2</sup>	19.0 <sup>1.5</sup>	20.1 <sup>1.4</sup>	R+M+F	<b>21.3<sup>1.9</sup></b>	<b>24.4<sup>1.7</sup></b>	<b>25.7<sup>1.9</sup></b>
<b>FARU</b>	1.6	2.6	5.0	7.9 <sup>2.2</sup>	8.7 <sup>2.5</sup>	8.9 <sup>2.7</sup>				
<i>Euro</i> <sub>5-6k</sub>										oracle: 84.4
<b>RAPP</b>	16.6	31.8	41.2	34.6 <sup>5.7</sup>	48.6 <sup>1.2</sup>	51.9 <sup>1.2</sup>				
<b>MIKO</b>	42.0	59.0	67.8	47.0 <sup>2.3</sup>	68.1 <sup>2.7</sup>	73.0 <sup>1.7</sup>	R+M	<b>49.5<sup>3.7</sup></b>	<b>68.7<sup>1.5</sup></b>	76.1 <sup>1.0</sup>
<b>FARU</b>	30.6	47.7	59.8	41.2 <sup>3.9</sup>	58.0 <sup>3.5</sup>	66.0 <sup>3.5</sup>				

# (Jakubina and Langlais, 2017) - Résultats

	INDIVIDUEL			RECLASSÉ				COMBINÉ & RECLASSÉ		
	@1	@5	@20	@1	@5	@20		@1	@5	@20
<i>Wiki</i> <sub>&gt;25</sub>										oracle: 69.3
<b>RAPP</b>	20.0	33.0	43.0	36.3 <sup>2.5</sup>	48.8 <sup>1.9</sup>	53.8 <sup>1.9</sup>	R+M+F	<b>45.6</b> <sup>2.2</sup>	<b>59.6</b> <sup>1.1</sup>	<b>64.0</b> <sup>1.8</sup>
<b>MIKO</b>	17.0	32.6	41.6	38.1 <sup>1.9</sup>	49.0 <sup>1.5</sup>	54.3 <sup>1.3</sup>				
<b>FARU</b>	13.3	26.0	33.3	34.3 <sup>1.5</sup>	44.0 <sup>2.6</sup>	47.9 <sup>2.1</sup>				
<i>Wiki</i> <sub>≤25</sub>										oracle: 28.6
<b>RAPP</b>	2.6	4.3	7.3	8.6 <sup>1.2</sup>	9.4 <sup>0.8</sup>	10.2 <sup>1.0</sup>	R+M+F	<b>21.3</b> <sup>1.9</sup>	<b>24.4</b> <sup>1.7</sup>	<b>25.7</b> <sup>1.9</sup>
<b>MIKO</b>	1.6	4.6	10.6	16.6 <sup>2.2</sup>	19.0 <sup>1.5</sup>	20.1 <sup>1.4</sup>				
<b>FARU</b>	1.6	2.6	5.0	7.9 <sup>2.2</sup>	8.7 <sup>2.5</sup>	8.9 <sup>2.7</sup>				
<i>Euro</i> <sub>5-6k</sub>										oracle: 84.4
<b>RAPP</b>	16.6	31.8	41.2	34.6 <sup>5.7</sup>	48.6 <sup>1.2</sup>	51.9 <sup>1.2</sup>	R+M	<b>49.5</b> <sup>3.7</sup>	<b>68.7</b> <sup>1.5</sup>	76.1 <sup>1.0</sup>
<b>MIKO</b>	42.0	59.0	67.8	47.0 <sup>2.3</sup>	68.1 <sup>2.7</sup>	73.0 <sup>1.7</sup>				
<b>FARU</b>	30.6	47.7	59.8	41.2 <sup>3.9</sup>	58.0 <sup>3.5</sup>	66.0 <sup>3.5</sup>				

# (Jakubina and Langlais, 2017) - Résultats

	INDIVIDUEL			RECLASSÉ				COMBINÉ & RECLASSÉ		
	@1	@5	@20	@1	@5	@20		@1	@5	@20
<i>Wiki</i> <sub>&gt;25</sub>										oracle: 69.3
<b>RAPP</b>	20.0	33.0	43.0	36.3 <sup>2.5</sup>	48.8 <sup>1.9</sup>	53.8 <sup>1.9</sup>				
<b>MIKO</b>	17.0	32.6	41.6	38.1 <sup>1.9</sup>	49.0 <sup>1.5</sup>	54.3 <sup>1.3</sup>	R+M+F	<b>45.6<sup>2.2</sup></b>	<b>59.6<sup>1.1</sup></b>	<b>64.0<sup>1.8</sup></b>
<b>FARU</b>	13.3	26.0	33.3	34.3 <sup>1.5</sup>	44.0 <sup>2.6</sup>	47.9 <sup>2.1</sup>				
<i>Wiki</i> <sub>≤25</sub>										oracle: 28.6
<b>RAPP</b>	2.6	4.3	7.3	8.6 <sup>1.2</sup>	9.4 <sup>0.8</sup>	10.2 <sup>1.0</sup>				
<b>MIKO</b>	1.6	4.6	10.6	16.6 <sup>2.2</sup>	19.0 <sup>1.5</sup>	20.1 <sup>1.4</sup>	R+M+F	<b>21.3<sup>1.9</sup></b>	<b>24.4<sup>1.7</sup></b>	<b>25.7<sup>1.9</sup></b>
<b>FARU</b>	1.6	2.6	5.0	7.9 <sup>2.2</sup>	8.7 <sup>2.5</sup>	8.9 <sup>2.7</sup>				
<i>Euro</i> <sub>5-6k</sub>										oracle: 84.4
<b>RAPP</b>	16.6	31.8	41.2	34.6 <sup>5.7</sup>	48.6 <sup>1.2</sup>	51.9 <sup>1.2</sup>				
<b>MIKO</b>	42.0	59.0	67.8	47.0 <sup>2.3</sup>	68.1 <sup>2.7</sup>	73.0 <sup>1.7</sup>	R+M	<b>49.5<sup>3.7</sup></b>	<b>68.7<sup>1.5</sup></b>	76.1 <sup>1.0</sup>
<b>FARU</b>	30.6	47.7	59.8	41.2 <sup>3.9</sup>	58.0 <sup>3.5</sup>	66.0 <sup>3.5</sup>				

# (Jakubina and Langlais, 2017) - Résultats

	INDIVIDUEL			RECLASSÉ				COMBINÉ & RECLASSÉ		
	@1	@5	@20	@1	@5	@20		@1	@5	@20
<i>Wiki</i> <sub>&gt;25</sub>										oracle: 69.3
<b>RAPP</b>	20.0	33.0	43.0	36.3 <sup>2.5</sup>	48.8 <sup>1.9</sup>	53.8 <sup>1.9</sup>				
<b>MIKO</b>	17.0	32.6	41.6	38.1 <sup>1.9</sup>	49.0 <sup>1.5</sup>	54.3 <sup>1.3</sup>	R+M+F	<b>45.6<sup>2.2</sup></b>	<b>59.6<sup>1.1</sup></b>	<b>64.0<sup>1.8</sup></b>
<b>FARU</b>	13.3	26.0	33.3	34.3 <sup>1.5</sup>	44.0 <sup>2.6</sup>	47.9 <sup>2.1</sup>				
<i>Wiki</i> <sub>≤25</sub>										oracle: 28.6
<b>RAPP</b>	2.6	4.3	7.3	8.6 <sup>1.2</sup>	9.4 <sup>0.8</sup>	10.2 <sup>1.0</sup>				
<b>MIKO</b>	1.6	4.6	10.6	16.6 <sup>2.2</sup>	19.0 <sup>1.5</sup>	20.1 <sup>1.4</sup>	R+M+F	<b>21.3<sup>1.9</sup></b>	<b>24.4<sup>1.7</sup></b>	<b>25.7<sup>1.9</sup></b>
<b>FARU</b>	1.6	2.6	5.0	7.9 <sup>2.2</sup>	8.7 <sup>2.5</sup>	8.9 <sup>2.7</sup>				
<i>Euro</i> <sub>5-6k</sub>										oracle: 84.4
<b>RAPP</b>	16.6	31.8	41.2	34.6 <sup>5.7</sup>	48.6 <sup>1.2</sup>	51.9 <sup>1.2</sup>				
<b>MIKO</b>	42.0	59.0	67.8	47.0 <sup>2.3</sup>	68.1 <sup>2.7</sup>	73.0 <sup>1.7</sup>	R+M	<b>49.5<sup>3.7</sup></b>	<b>68.7<sup>1.5</sup></b>	76.1 <sup>1.0</sup>
<b>FARU</b>	30.6	47.7	59.8	41.2 <sup>3.9</sup>	58.0 <sup>3.5</sup>	66.0 <sup>3.5</sup>				

# (Jakubina and Langlais, 2017) - Résultats

	INDIVIDUEL			RECLASSÉ				COMBINÉ & RECLASSÉ		
	@1	@5	@20	@1	@5	@20		@1	@5	@20
<i>Wiki</i> <sub>&gt;25</sub>										oracle: 69.3
<b>RAPP</b>	20.0	33.0	43.0	36.3 <sup>2.5</sup>	48.8 <sup>1.9</sup>	53.8 <sup>1.9</sup>	R+M+F	<b>45.6</b> <sup>2.2</sup>	<b>59.6</b> <sup>1.1</sup>	<b>64.0</b> <sup>1.8</sup>
<b>MIKO</b>	17.0	32.6	41.6	38.1 <sup>1.9</sup>	49.0 <sup>1.5</sup>	54.3 <sup>1.3</sup>				
<b>FARU</b>	13.3	26.0	33.3	34.3 <sup>1.5</sup>	44.0 <sup>2.6</sup>	47.9 <sup>2.1</sup>				
<i>Wiki</i> <sub>≤25</sub>										oracle: 28.6
<b>RAPP</b>	2.6	4.3	7.3	8.6 <sup>1.2</sup>	9.4 <sup>0.8</sup>	10.2 <sup>1.0</sup>	R+M+F	<b>21.3</b> <sup>1.9</sup>	<b>24.4</b> <sup>1.7</sup>	<b>25.7</b> <sup>1.9</sup>
<b>MIKO</b>	1.6	4.6	10.6	16.6 <sup>2.2</sup>	19.0 <sup>1.5</sup>	20.1 <sup>1.4</sup>				
<b>FARU</b>	1.6	2.6	5.0	7.9 <sup>2.2</sup>	8.7 <sup>2.5</sup>	8.9 <sup>2.7</sup>				
<i>Euro</i> <sub>5-6k</sub>										oracle: 84.4
<b>RAPP</b>	16.6	31.8	41.2	34.6 <sup>5.7</sup>	48.6 <sup>1.2</sup>	51.9 <sup>1.2</sup>	R+M	<b>49.5</b> <sup>3.7</sup>	<b>68.7</b> <sup>1.5</sup>	76.1 <sup>1.0</sup>
<b>MIKO</b>	42.0	59.0	67.8	47.0 <sup>2.3</sup>	68.1 <sup>2.7</sup>	73.0 <sup>1.7</sup>				
<b>FARU</b>	30.6	47.7	59.8	41.2 <sup>3.9</sup>	58.0 <sup>3.5</sup>	66.0 <sup>3.5</sup>				

# (Jakubina and Langlais, 2017) - Résultats

INDIVIDUEL				RECLASSÉ			COMBINÉ & RECLASSÉ			
	@1	@5	@20	@1	@5	@20		@1	@5	@20
<i>Wiki</i> <sub>&gt;25</sub>								oracle: 69.3		
RAPP	20.0	33.0	43.0	36.3 <sup>2.5</sup>	48.8 <sup>1.9</sup>	53.8 <sup>1.9</sup>	R+M+F	45.6 <sup>2.2</sup>	59.6 <sup>1.1</sup>	64.0 <sup>1.8</sup>
MIKO	17.0	32.6	41.6	38.1 <sup>1.9</sup>	49.0 <sup>1.5</sup>	54.3 <sup>1.3</sup>		21.3 <sup>1.9</sup>	24.4 <sup>1.7</sup>	25.7 <sup>1.9</sup>
FARU	13.3	26.0	33.3	34.3 <sup>1.5</sup>	44.0 <sup>2.6</sup>	47.9 <sup>2.1</sup>				
<i>Wiki</i> <sub>≤25</sub>								oracle: 28.6		
RAPP	2.6	4.3	7.3	8.6 <sup>1.2</sup>	9.4 <sup>0.8</sup>	10.2 <sup>1.0</sup>	R+M+F	21.3 <sup>1.9</sup>	24.4 <sup>1.7</sup>	25.7 <sup>1.9</sup>
MIKO	1.6	4.6	10.6	16.6 <sup>2.2</sup>	19.0 <sup>1.5</sup>	20.1 <sup>1.4</sup>				
FARU	1.6	2.6	5.0	7.9 <sup>2.2</sup>	8.7 <sup>2.5</sup>	8.9 <sup>2.7</sup>				
<i>Euro</i> <sub>5-6k</sub>								oracle: 84.4		
RAPP	16.6	31.8	41.2	34.6 <sup>5.7</sup>	48.6 <sup>1.2</sup>	51.9 <sup>1.2</sup>	R+M	49.5 <sup>3.7</sup>	68.7 <sup>1.5</sup>	76.1 <sup>1.0</sup>
MIKO	42.0	59.0	67.8	47.0 <sup>2.3</sup>	68.1 <sup>2.7</sup>	73.0 <sup>1.7</sup>				
FARU	30.6	47.7	59.8	41.2 <sup>3.9</sup>	58.0 <sup>3.5</sup>	66.0 <sup>3.5</sup>				

# (Jakubina and Langlais, 2017) - Résultats

	INDIVIDUEL			RECLASSÉ			COMBINÉ & RECLASSÉ			
	@1	@5	@20	@1	@5	@20	@1	@5	@20	
<i>Wiki</i> <sub>&gt;25</sub>							oracle: 69.3			
RAPP	20.0	33.0	43.0	36.3 <sup>2.5</sup>	48.8 <sup>1.9</sup>	53.8 <sup>1.9</sup>	R+M+F	45.6 <sup>2.2</sup>	59.6 <sup>1.1</sup>	64.0 <sup>1.8</sup>
MIKO	17.0	32.6	41.6	38.1 <sup>1.9</sup>	49.0 <sup>1.5</sup>	54.3 <sup>1.3</sup>				
FARU	13.3	26.0	33.3	34.3 <sup>1.5</sup>	44.0 <sup>2.6</sup>	47.9 <sup>2.1</sup>				
<i>Wiki</i> <sub>≤25</sub>							oracle: 28.6			
RAPP	2.6	4.3	7.3	8.6 <sup>1.2</sup>	9.4 <sup>0.8</sup>	10.2 <sup>1.0</sup>	R+M+F	21.3 <sup>1.9</sup>	24.4 <sup>1.7</sup>	25.7 <sup>1.9</sup>
MIKO	1.6	4.6	10.6	16.6 <sup>2.2</sup>	19.0 <sup>1.5</sup>	20.1 <sup>1.4</sup>				
FARU	1.6	2.6	5.0	7.9 <sup>2.2</sup>	8.7 <sup>2.5</sup>	8.9 <sup>2.7</sup>				
<i>Euro</i> <sub>5-6k</sub>							oracle: 84.4			
RAPP	16.6	31.8	41.2	34.6 <sup>5.7</sup>	48.6 <sup>1.2</sup>	51.9 <sup>1.2</sup>	R+M	49.5 <sup>3.7</sup>	68.7 <sup>1.5</sup>	76.1 <sup>1.0</sup>
MIKO	42.0	59.0	67.8	47.0 <sup>2.3</sup>	68.1 <sup>2.7</sup>	73.0 <sup>1.7</sup>				
FARU	30.6	47.7	59.8	41.2 <sup>3.9</sup>	58.0 <sup>3.5</sup>	66.0 <sup>3.5</sup>				



- **FARU** : performances similaires à **RAPP** et **MIKO** (biais de fréquence)
- Reclassement : ✓
- Combinaison (avec reclassement) : ✓
  - Améliorations de 15% en moyenne
  - Notamment, sur les **Mots Rares**
  - Avec peu de features
- Analyses complémentaires
  - Taille données d'entraînements
  - Impacts du reclassement et des features
  - Analyse d'erreurs plus détaillée<sup>9</sup>

---

<sup>9</sup>Morpho ; Relié (syno, anto, hypo, ...) ; etc.

# (Jakubina and Langlais, 2017) - Exemples de sorties

Wiki<sub>>25</sub> - RERANKER

donut (beigne)	donuts	donut	grillé	baozi
brilliantly (brillamment)	brillant	éclatant	éblouissant	imaginatif
gentle (doux)	colérique	enjoué	insouciant	affable
pathologically (pathologiquement)	pathologique	sigmoïdite	fibromatose	nosophobie

Wiki<sub>≤25</sub>

raillery | raillerie

<b>RERANKER</b>	<b>raillerie</b>	railler	raille	railleries	<b>1</b>
<b>RAPP</b>	pourra	désire	...	<b>raillerie</b>	5
<b>MIKO</b>	reprocher	naïveté	...	<b>raillerie</b>	70
<b>FARU</b>	intercommunal	plaidoiries	fourrière	plaida	∅

schizoidism | schizoïdie

<b>RERANKER</b>	<b>schizoïdie</b>	musicothérapeute	confusionnels	mentalisation	<b>1</b>
<b>RAPP</b>	refactorisat.	disciples	knowlse	boutetage	∅
<b>MIKO</b>	atrophier	hypomania.	psychoneurol.	neurochimiques	∅
<b>FARU</b>	hypomania.	psychoneurol.	...	<b>schizoïdie</b>	65

# (Jakubina and Langlais, 2017) - Exemples de sorties

Wiki<sub>≤25</sub>

botanize		herboriser				
<b>RERANKER</b>		herborise	herboriser	prodigua	attardant	2
<b>RAPP</b>		aoutourou	délai	...	herboriser	6
<b>MIKO</b>		attardant	herboriser	prodigua	hastrel	2
<b>FARU</b>		compromissions	lorois	caboteur	pontée	∅

uncrushable		infoissable				
<b>RERANKER</b>		infoissable	incrustait	raquettistes	paludicroque	1
<b>RAPP</b>		senente	imbroyables	pulvérix	attriteur	∅
<b>MIKO</b>		nattées	pauimage	infoissable	mouillettes	3
<b>FARU</b>		roninson	ospovat	talánov	mouraviova	∅

brotherliness		confraternité				
<b>RERANKER</b>		volonté	précepte	fraternel	désintéressé	3
<b>RAPP</b>		qoudous	tâche	...	fraternel	44
<b>MIKO</b>		joie	volonté	...	fraternel	68
<b>FARU</b>		observant	cueillies	concordent	moisson	∅

## Travaux futurs & Conclusion

---

# Axes de travaux futurs

- Combinaison - Reclassement
  - Reproduction sur d'autres paires de langues
  - Autres features / approches (Irvine and Callison-Burch, 2013)<sup>10</sup>
- Améliorations modèles / approches
  - Instanciation du modèle approprié (selon la fréquence)
  - Quid d'un modèle incluant les features de la combinaison ?
- Mots rares
  - Modèle(s) d'embeddings adapté(s)
- Réflexion méta sur l'**ILB**
  - Métriques d'évaluation
  - Échantillonnage(s) approprié(s) / intelligent(s)
- Retour sur... **DBPEDIA** ?

---

<sup>10</sup>contexte, temporel, orthographe, thématique et fréquence.

Déjà réalisé<sup>11</sup> : Allemand, Espagnol et Roumain (vs Anglais)

⇒ Approche stable sur d'autres paires de langues

---

<sup>11</sup>Publication en cours de rédaction.

# Mots rares - Modèle(s) d'embeddings adapté(s)

Améliorer les modèles des mots rares en combinant l'information de contexte avec d'autres informations:

- Rothe, S. and Schütze, H. (2015). Autoextend: Extending word embeddings to embeddings for synsets and lexemes
- Avraham, O. and Goldberg, Y. (2017). The interplay of semantics and morphology in word embeddings
- Pilehvar, M. T. and Collier, N. (2017). Inducing embeddings for rare and unseen words by leveraging lexical resources
- Bamler, R. and Mandt, S. (2017). Dynamic word embeddings
  - Embeddings intègrent le(s) sens du mot "dans le temps"
- Calixto, I., Liu, Q., and Campbell, N. (2017). Multilingual multi-modal embeddings for natural language processing
  - Embeddings améliorés en apprenant avec du texte ET des images

Possible aussi pour les mots fréquents

# Conclusion

- Tâche d'Induction de Lexiques Bilingues sur deux types de corpus bilingues
- Corpus parallèles (rare): approche efficace mais couverture insuffisante
- Corpus comparables (fréquent): modèles basé sur le contexte
  - Performances convaincantes sur les mots fréquents
  - Trouvent leurs limites avec les mots rares
- Améliorations satisfaisantes de l'**ILB** grâce à la combinaison et au reclassement supervisé
- Nombreux axes de travaux futurs ; Intérêts portés vers le reclassement et les mots rares



Merci !

Questions ?

- Adaptation de l'Analyse Semantique Explicite (ESA): utilise Wikipedia comme base de connaissance (Gabrilovich and Markovitch, 2007)
- Espace des mots  $\rightarrow$  Espace des titres (vecteur)
- Dictionnaire  $\rightarrow$  Lien interlangue Wikipedia

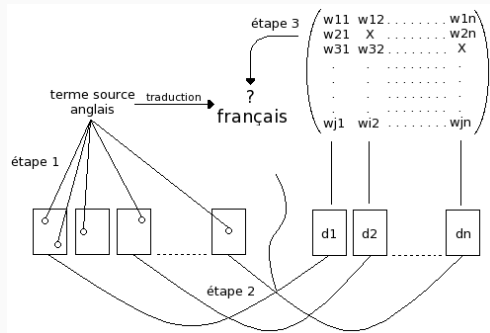


Figure 1: Approche ESA Bouamor

# ILB comme tâche d'évaluation

Apprendre une projection des embeddings anglais vers les embeddings français

Embeddings Anglais					Project <sup>o</sup> (Dinu and Baroni, 2014)				Embeddings Français			
...	...	...	...		...	...	...		...	...	...	...
<i>text</i>	0.22	0.09	...		1.73	4.83	...		<i>texte</i>	0.70	1.56	...
<i>cafeteria</i>	-0.32	-0.28	...	↔	-7.21	-3.93	...	↔	<i>un</i>	0.06	-1.44	...
<i>stereo</i>	1.97	-0.30	...		5.07	-3.41	...		<i>cafeteria</i>	0.00	0.36	...
<i>dummy</i>	0.28	0.24	...		5.42	8.51	...		<i>ville</i>	1.75	-1.43	...
<i>one</i>	-0.36	0.23	...		1.57	-1.25	...		<i>saint</i>	1.31	-0.03	...
...	...	...	...		...	...	...		...	...	...	...

(Apprendre **directement** des embeddings dits interlingues / bilingues)

Embeddings Bilingues					
...	...	...	...	...	...
<i>text/texte</i>	-0.33	0.17	-0.43	0.92	...
<i>cafeteria/cafeteria</i>	-0.63	0.45	-0.09	-0.28	...
<i>stereo/stereo</i>	-0.32	0.04	-0.65	0.08	...
<i>dummy/factice</i>	-0.12	-0.27	0.03	0.26	...
<i>one/un</i>	-0.04	-0.12	0.05	-0.47	...
...	...	...	...	...	...

## Resultats Article n°2

- Best variant for each approach according to TOP@1.
- An *Oracle* shows that approaches are complementary.
- Disappointment for the poor performance of the document approach which was specifically designed to handle rare words.

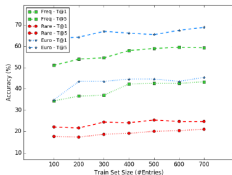
	<i>Wiki</i> <sub>≤25</sub>			<i>Wiki</i> <sub>&gt;25</sub>		
	@1	@5	@20	@1	@5	@20
<b>MIKO</b>	2,2	6,1	11,9	21,7	34,2	44,9
<b>RAPP</b>	2,0	4,3	7,6	19,0	32,7	44,3
Doc	0,7	2,3	5,0	10,0	19,0	24,0
<i>oracle</i>	4,6	10,5	19,0	31,8	46,8	57,6

# Analyse Reranker

## Analysis

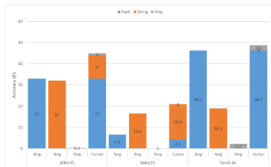
### Training Size

- Test = 300 test words (fixed).
- More training examples  
→ better results.
- But 200 training examples seem good enough.



### Feature Selection

- Overall: rank >> string >> frequency.
- Using all features is preferable.



### Reranker Analysis

- Always keeps solutions proposed at @1 by (at least) two approaches.
- Rare that the reranker picks the reference solution at @1 if no approach did it.
- Still the average rank of the reference solution decreases.

	Wiki <sub>&gt;25</sub>	Wiki <sub>≤25</sub>	Euro <sub>5-6k</sub>
Rapp	12.7	19.6	16.2
Miko	16.3	30.0	7.5
Faru	20.4	35.5	11.3
<i>list-oracle</i>	12.3	9.1	7.1
reranker	5.6	4.0	4.9

Average rank of the reference translation

### Error Analysis

- Manual inspection of the first candidate by our best reranker for the first 100 test forms for which the candidate translation differs from the reference one.
- % of JUNK errors is much higher on Wiki<sub>>25</sub>  
→ another illustration of the bias in favor of frequent terms.

	Wiki <sub>&gt;25</sub>	Wiki <sub>≤25</sub>	Euro <sub>5-6k</sub>
MORPHO	18	3	26
RELATED	16	4	23
synonyms	15	1	19
antonyms	1	2	2
hyponym			1
cohyponym		1	1
POLYSEMY	4	0	5
LOOSLY	14	15	20
ENGLISH	21	6	7
JUNK	27	72	19

## References

---

Ammar, W., Mulcaire, G., Tsvetkov, Y., Lample, G., Dyer, C., and Smith, N. A. (2016).

**Massively multilingual word embeddings.**

*arXiv preprint arXiv:1602.01925.*

Avraham, O. and Goldberg, Y. (2017).

**The interplay of semantics and morphology in word embeddings.**

Bamler, R. and Mandt, S. (2017).

**Dynamic word embeddings.**

Bouamor, D., Popescu, A., Semmar, N., and Zweigenbaum, P. (2013).

**Building specialized bilingual lexicons using large scale background knowledge.**

In *EMNLP*, pages 479–489.

Bourdaillet, J., Huet, S., Langlais, P., and Lapalme, G. (2010).

**TransSearch: from a bilingual concordancer to a translation finder.**

*Machine Translation*, 24(3-4):241–271.

Brown, P. F., Cocke, J., Pietra, S. A. D., Pietra, V. J. D., Jelinek, F., Lafferty, J. D., Mercer, R. L., and Roossin, P. S. (1990).

**A statistical approach to machine translation.**

*Computational linguistics*, 16(2):79–85.

Calixto, I., Liu, Q., and Campbell, N. (2017).

**Multilingual multi-modal embeddings for natural language processing.**

Coulmance, J., Marty, J.-M., Wenzek, G., and Benhalloum, A. (2015).

**Trans-gram, Fast Cross-lingual Word-embeddings.**

*In Proceedings of EMNLP*, pages 1109–1113.

Dinu, G. and Baroni, M. (2014).

**Improving zero-shot learning by mitigating the hubness problem.**

*CoRR*.



Faruqui, M. and Dyer, C. (2014).

**Improving Vector Space Word Representations Using Multilingual Correlation.**

*In Proceedings of EACL.*

Fung, P. (1995).

**Compiling bilingual lexicon entries from a non-parallel english-chinese corpus.**

*In Third Workshop on Very Large Corpora, pages 173–183.*

Gabrilovich, E. and Markovitch, S. (2007).

**Computing semantic relatedness using wikipedia-based explicit semantic analysis.**

*IJCAI'07, pages 1606–1611.*

Gabrilovich, E. and Markovitch, S. (2009).

**Wikipedia-based semantic interpretation for natural language processing.**

*Journal of Artificial Intelligence Research*, 34(2):443.

Gaussier, E., Renders, J.-M., Matveeva, I., Goutte, C., and Déjean, H. (2004).

**A geometric view on bilingual lexicon extraction from comparable corpora.**

In *Proceedings of the 42nd Annual Meeting on Association for Computational Linguistics*, page 526. Association for Computational Linguistics.

Gouws, S., Bengio, Y., and Corrado, G. (2015).

**BilBOWA: Fast Bilingual Distributed Representations without Word Alignments.**

*In Proceedings of the 32nd ICML*, pages 748–756.

Gómez-Pérez, A., Vila-Suero, D., Montiel-Ponsoda, E., Gracia, J., and Aguado-de Cea, G. (2013).

**Guidelines for multilingual linked data.**

WIMS '13, pages 3:1–3:12.

Irvine, A. and Callison-Burch, C. (2013).

**Supervised Bilingual Lexicon Induction with Multiple Monolingual Signals.**

*In Proceedings of NAACL-HLT*, pages 518–523.

Jakubina, L. and Langlais, P. (2015).

**Projective methods for mining missing translations in dbpedia.**

*In Proceedings of the Eighth Workshop on Building and Using Comparable Corpora*, pages 23–31, Beijing, China. Association for Computational Linguistics.

Jakubina, L. and Langlais, P. (2016a).

**Bad luc@wmt 2016: a bilingual document alignment platform based on lucene.**

*In Proceedings of the First Conference on Machine Translation*, pages 703–709, Berlin, Germany. Association for Computational Linguistics.

Jakubina, L. and Langlais, P. (2016b).

**A comparison of methods for identifying the translation of words in a comparable corpus: Recipes and limits.**

*CICLing: International Conference on Computational Linguistics and Intelligent Text Processing.*

Jakubina, L. and Langlais, P. (2017).

**Reranking translation candidates produced by several bilingual word similarity sources.**

*In Proceedings of the 15th Conference of the European Chapter of the Association for Computational Linguistics: Volume 2, Short Papers*, pages 605–611, Valencia, Spain. Association for Computational Linguistics.

Levy, O., Goldberg, Y., and Dagan, I. (2015).

**Improving distributional similarity with lessons learned from word embeddings.**

*Transactions of ACL*, pages 211–225.

Mikolov, T., Chen, K., Corrado, G., and Dean, J. (2013a).

**Efficient estimation of word representations in vector space.**

Mikolov, T., Le, Q. V., and Sutskever, I. (2013b).

**Exploiting similarities among languages for machine translation.**

*CoRR*.

Morin, E. and Prochasson, E. (2011).

**Bilingual lexicon extraction from comparable corpora enhanced with parallel corpora.**

*In Proceedings of the 4th workshop on building and using comparable corpora: comparable corpora and the web*, pages 27–34.

Pilehvar, M. T. and Collier, N. (2017).

**Inducing embeddings for rare and unseen words by leveraging lexical resources.**

Prochasson, E. and Fung, P. (2011).

**Rare word translation extraction from aligned comparable documents.**

HLT '11, pages 1327–1335.

Rapp, R. (1995).

**Identifying Word Translations in Non-parallel Texts.**

*In Proceedings of the 33rd ACL*, pages 320–322.

Rapp, R. (1999).

**Automatic identification of word translations from unrelated English and German corpora.**

*In Proceedings of the 37th annual meeting of the Association for Computational Linguistics on Computational Linguistics*, pages 519–526. Association for Computational Linguistics.

Rothe, S. and Schütze, H. (2015).

**Autoextend: Extending word embeddings to embeddings for synsets and lexemes.**



Sharoff, S., Rapp, R., and Zweigenbaum, P. (2013).

**Overviewing Important Aspects of the Last Twenty Years of Research in Comparable Corpora.**

In *Building and Using Comparable Corpora*, pages 1–17. Springer Berlin Heidelberg.

Vulić, I., De Smet, W., and Moens, M.-F. (2011).

**Identifying word translations from comparable corpora using latent topic models.**

In *Proceedings of the 49th Annual Meeting of the Association for Computational Linguistics: Human Language Technologies: short papers-Volume 2*, pages 479–484. Association for Computational Linguistics.

Vulic, I., Kiela, D., Clark, S., and Moens, M.-F. (2016).

**Multi-modal representations for improved bilingual lexicon learning.**

*In The 54th Annual Meeting of the Association for Computational Linguistics, page 188.*