



**HAL**  
open science

# Appariement de phrases courtes pour la traduction automatique par l'exemple

Julien Gosme

► **To cite this version:**

Julien Gosme. Appariement de phrases courtes pour la traduction automatique par l'exemple. MAnifestation des JEunes Chercheurs en Sciences et Technologies de l'Information et de la Communication, Nov 2009, Avignon, France. pp.972. hal-00439892

**HAL Id: hal-00439892**

**<https://hal.science/hal-00439892>**

Submitted on 8 Dec 2009

**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.

# Appariement de phrases courtes pour la traduction automatique par l'exemple

Julien Gosme

Laboratoire GREYC, Université de Caen Basse-Normandie  
Contact: [julien.gosme@info.unicaen.fr](mailto:julien.gosme@info.unicaen.fr)

---

## Résumé

La constitution de ressources linguistiques pour les systèmes de traduction automatique fondée sur les données est une tâche critique. Ces systèmes de traduction ont besoin de corpus de phrases alignées pour chaque couple de langues. La constitution de telles ressources est généralement effectuée à la main par des traducteurs. Nous proposons une méthode automatisant la constitution de corpus bilingues de phrases courtes en employant une représentation vectorielle bien connue en recherche d'information. Un dictionnaire bilingue est nécessaire par couple de langues considéré. Nous utilisons le Web afin de constituer des corpus de documents monolingues sur lesquels appliquer l'appariement de phrases courtes. Le coût humain total de la constitution d'un corpus bilingue de phrase est extrêmement réduit: seule une validation manuelle est nécessaire après appariement. Une expérience en français-anglais permet d'estimer la précision de la méthode d'appariement. 800 phrases traductions ont été collectées avec une précision supérieure ou égale à 0,8 à partir de 100 000 phrases collectées dans chaque langue.

## Abstract

Constitution of linguistic resources for data-driven translation system is often a critical task. Those translation systems need corpus of aligned sentences for each language pairs. The constitution of those linguistic resources are often done by human translators. We propose a method designed to automate the constitution of bilingual corpus of short sentences. This tool is based on a vector representation, commonly used in the information retrieval field. A bilingual dictionary is required for each language pair considered. The Web provides monolingual corpus of sentences. The pairing method is then applied to those sentences. The human cost of the overall method of constitution of bilingual corpus of sentence is extremely reduced: once the pairing applied, only a human validation is required on the bicorpus. An experiment on French-English data shows the precision of the pairing method. 800 sentences have been collected with a precision of 0.80 among 100,000 collected sentences for each language.

**Mots-clés :** appariement de phrases, alignement de phrases, constitution automatique de ressources linguistiques, traduction automatique fondée sur les données

**Keywords:** sentence pairing, sentence alignment, automated constitution of linguistic resources, data-driven machine translation

---

## 1. Introduction

La traduction automatique fondée sur les données comme la traduction automatique statistique et la traduction automatique par l'exemple est un paradigme qui nécessite des quantités toujours plus importantes de données pour améliorer ses performances.

Les meilleures performances sont obtenues avec les couples les plus documentés tel que l'anglais-français puisqu'il existe d'importants corpus bilingues comme Hansard ou Europarl.

Un couple de langue est qualifié de peu documenté lorsqu'il n'existe aucun corpus de documents alignés disponible pour ces langues.

Cependant des langues séparément bien documentées comme le chinois, le français ou le japonais peuvent former des couples peu documentés. C'est le cas du couple chinois-français. Bien

évidemment un couple de langues constitué d'au moins une langue peu documentée est peu documenté. Par exemple, le couple japonais-vietnamien est peu documenté car le vietnamien est peu documenté : le nombre de sites Web bilingues est très petit.

La méthode habituelle pour obtenir des phrases alignées, celle de Gale et Church [3], utilise des textes parallèles, c'est-à-dire exactement traductions l'un de l'autre. Les textes sont segmentés afin d'obtenir des passages alignés. Généralement, les passages alignés sont constitués d'une phrase, parfois plus. Notons que cette segmentation ne met aucun passage de côté : l'intégralité du texte est conservé.

Étant donné le manque de textes parallèles pour les couples de langues peu documentés, l'alignement de phrases de Gale et Church [3] ne peut pas s'y appliquer. Si les langues ont une bonne visibilité sur le Web, alors il est désirable de définir une méthode pour la constitution de corpus de phrases alignées.

Notre idée est d'apparier des phrases aspirées depuis des sites Web monolingues, sans imposer de conditions strictes sur l'alignement des textes. Si un dictionnaire bilingue existe pour le couple de langues considéré, alors nous pouvons l'utiliser pour apparier des phrases de ces deux langues.

En pratique, nous collecterons donc des phrases depuis le Web puis nous les apparierons à l'aide d'une méthode vectorielle. Dans cet article, nous nous intéressons au couple de langues français-anglais afin de valider la méthode d'appariement. Nous étendrons notre étude à d'autres couples de langues par la suite.

## 2. Description du problème

Le Web est déjà une ressource exploitée pour la constitution de corpus monolingues et bilingues. Par exemple, Nie [6] a constitué un corpus de document anglais-français pour la recherche d'information multilingue. La sélection des documents parallèles candidats est basée sur la similarité des URLs d'un même document en plusieurs langues.

Resnik [7] a développé Strands, un outils de recherche et de collecte de documents parallèles sur le Web. La structure et la longueur des documents sont utilisées pour identifier les documents parallèles.

Ma [5] a développé un outil dans le même but. La différence principale avec les deux outils précédents est la méthode utilisée pour identifier les documents parallèles : ils sont identifiés selon leur similarité lexicale, estimée à l'aide d'un lexique de traduction.

Notre méthode d'appariement se rapproche plus de l'outil de Ma que des deux autres.

Comme les outils précédents se focalisent sur l'extraction de documents depuis le Web, il est nécessaire d'appliquer un alignement de phrases pour utiliser ces données dans des systèmes de traduction automatique fondée sur les données.

Notre méthode ne rencontre pas ce problème. Au lieu de collecter des documents, d'en extraire les textes parallèles puis d'en extraire des phrases alignées, nous proposons de collecter des documents et d'en extraire directement des phrases appariées.

Ceci est un avantage car la plupart des sites Web multilingues ne sont pas composés de documents alignés mais plutôt de documents comparables. Deux documents comparables ne sont pas exactement traductions l'un de l'autre, certains passages d'un document peuvent n'avoir aucun équivalent dans un document comparable dans une autre langue. Comme exemple de corpus comparables, nous pouvons citer les articles de l'encyclopédie Wikipédia sur les mêmes sujets car les textes y sont souvent adaptés à la culture ou à l'histoire du public destinataire, ou les articles de journaux de langues différentes faisant référence aux mêmes événements aux mêmes dates.

Nous nous intéresserons donc aux textes comparables. Contrairement aux trois outils précédemment cités, nous nous concentrerons sur la précision des résultats et non sur le rappel. En effet, comme le Web est gigantesque et qu'il est en constant accroissement, il est préférable d'explorer une quantité plus importante de sites et collecter plus de phrases plutôt que de réduire la précision pour augmenter la rappel.

Nous choisissons de construire un outil multilingue requérant un minimum de ressources linguistiques pour les langues traitées. Notre outil requiert seulement un dictionnaire bilingue par couple de langues traité.

TAB. 1 – Correspondances entre la terminologie de l'appariement de phrases et la terminologie de la recherche d'information.

Appariement de phrases	Recherche d'information
mot	terme
phrase	document
matrice mot-phrase	matrice terme-document
phrase source	document requête
phrase cible	document similaire

### 3. Aspirateur de document Web monolingue

Nous avons écrit notre propre aspirateur Web afin d'en contrôler toutes les étapes. Seul le texte des pages collectées est utilisé, nous épurons les mises en forme des pages Web. Nous avons inclus une étape de segmentation en phrases.

Nous avons également inclus un module d'identification des langues au niveau des documents aspirés. La méthode employée est celle des profil de n-grammes de Cavnar et Trenkle [1]. Nous substituons la distance définie par Cavnar et Trenkle par un calcul d'angle entre profils de n-grammes. Nous avons déterminé empiriquement que des profils de trigrammes d'octets sur des documents codés en UTF-8 donnent des résultats satisfaisants. Comme la méthode de Cavnar et Trenkle a besoin d'exemples de textes pour chaque langue considérée, nous avons choisi d'utiliser la page intitulée « Wikipédia » dans l'encyclopédie éponyme en anglais et en français comme jeu d'entraînement pour nos profils de trigrammes.

### 4. Emprunt d'une technique de recherche d'information et adaptation au problème de l'appariement de phrases

L'idée centrale de notre méthode est d'adapter une technique bien connue de la recherche d'informations à nos besoins.

Dans cette technique, les relations entre termes et documents sont représentées sous forme vectorielle. En particulier, la représentation de Salton [8] permet de représenter aussi bien les termes que les documents dans un même espace vectoriel d'index.

Le parallèle que nous établissons entre les terminologies de l'appariement de phrases et de la recherche d'information est présenté dans le tableau 1.

Afin d'adapter cette méthode à l'appariement de phrases, nous substituons la liste de termes par un dictionnaire bilingue. Toute phrase (les documents dans la terminologie de la recherche d'information) peut alors être représentée par un vecteur dans l'espace vectoriel du dictionnaire. Grâce à cette méthode, nous représentons les phrases sources dans l'espace vectoriel de la langue source et les phrases cibles dans l'espace vectoriel de la langue cible. Pour comparer phrases sources et phrases cibles, il est nécessaire de projeter les phrases sources de l'espace vectoriel source sur l'espace vectoriel cible. À cette fin, nous utilisons de nouveau le dictionnaire bilingue. Il peut en effet être utilisé comme matrice de changement de base. La projection est alors simplement le produit d'un vecteur-phrase source par la matrice de changement de base.

Une méthode classique utilisée en recherche d'information pour estimer la similarité entre ensemble d'objets dont on connaît le nombre d'occurrences est celle du cosinus (voir Dumais [2] et Landauer [4]). Nous utilisons cette méthode comme estimation de la qualité des appariements produits.

### 5. Exposé didactique

L'exemple suivant montre chaque étape de l'appariement en détail. Le couple de langues utilisé dans cet exposé didactique est l'anglais-français.

Un petit corpus de phrases anglaises et françaises est donné dans le tableau 2.

Les comparaisons de chaînes de caractères sont effectuées sans prendre en compte la casse des caractères ni la ponctuation.

TAB. 2 – Petit exemple de corpus de phrases anglaises et françaises.

Phrase	Langue
Le premier janvier.	français
Je voudrais du pain.	français
Je voudrais une tasse de thé.	français
I would like some bread.	anglais
I would like a cup of tea.	anglais

### 5.1. Dictionnaire bilingue français-anglais et matrices de changement de bases

Le dictionnaire bilingue français-anglais utilisé dans cet exemple est donné dans le tableau 3.

Le dictionnaire est traité de la même manière que les phrases du corpus, c'est-à-dire sans prendre en compte la casse ni la ponctuation.

TAB. 3 – Exemple de dictionnaire bilingue français-anglais.

Français	Anglais
premier	first
janvier	january
je	I
vouloir	would
du	some
pain	bread
coupe	cup
tasse	cup
thé	tea
manger	eat
boire	drink

Le dictionnaire bilingue peut être représenté dans les espaces vectoriels suivants :

$$(v_{\text{premier}}, v_{\text{janvier}}, v_{\text{je}}, \dots), \text{ pour le français}$$

$$(v_{\text{first}}, v_{\text{january}}, v_{\text{I}}, \dots), \text{ pour l'anglais}$$

Seuls les mots communs au corpus et au dictionnaire sont utilisés pour construire la matrice de changement de bases. Dans cet exemple, les deux dernières entrées sont écartées. De cette manière le temps de calcul se trouve réduit sans aucun effet sur la précision. La représentation matricielle du dictionnaire bilingue est donc :



$$\begin{array}{cccccc}
 \text{first} & \text{january} & & & & \\
 & \text{I} & \text{would} & \text{some} & \text{bread} & \text{cup} \\
 \left( \begin{array}{cccccc}
 \cdot & \cdot & \frac{1}{5} & \frac{1}{5} & \frac{1}{5} & \frac{1}{5} \\
 \cdot & \cdot & \frac{1}{7} & \frac{1}{7} & \cdot & \frac{1}{7}
 \end{array} \right) & \begin{array}{l} \text{I would like some bread.} \\ \text{I would like a cup of tea.} \end{array}
 \end{array}$$

### 5.3. Changement de bases vectorielles

Afin de comparer les vecteurs-phrases exprimés dans des bases vectorielles différentes, les vecteurs-phrases français sont projetés dans l'espace vectoriel anglais.

Il suffit alors de multiplier la matrice des vecteurs-phrases français  $P^f$  par la matrice de changement de bases du français à l'anglais  $A^{f \rightarrow e}$ . Le résultat de cette projection est la matrice  $\tilde{P}^e = P^f \times A^{f \rightarrow e}$  :

$$\begin{array}{cccccc}
 \text{first} & \text{january} & & & & \\
 & \text{I} & \text{would} & \text{some} & \text{bread} & \text{cup} \\
 \left( \begin{array}{cccccc}
 \frac{1}{3} & \frac{1}{3} & \cdot & \cdot & \cdot & \cdot \\
 \cdot & \cdot & \frac{1}{4} & \cdot & \frac{1}{4} & \frac{1}{4} \\
 \cdot & \cdot & \frac{1}{6} & \cdot & \cdot & \frac{1}{6}
 \end{array} \right) & \begin{array}{l} \text{Le premier janvier.} \\ \text{Je voudrais du pain.} \\ \text{Je voudrais une tasse de thé.} \end{array}
 \end{array}$$

Cette projection est une traduction brute des phrases considérées comme sacs de mots, c'est-à-dire sans tenir compte des positions des mots. Les angles formés entre les vecteurs de la matrice  $P^e$  et les vecteurs de la matrice  $\tilde{P}^e$  donnent une indication sur la qualité de l'appariement des deux phrases correspondantes. Les angles sont des nombres compris entre 0 et  $\frac{\pi}{2}$ . La matrice suivante contient les angles formés entre tous les vecteurs-phrases :

$$\begin{array}{cc}
 \text{I would like some bread.} \\
 \text{I would like a cup of tea.} \\
 \left( \begin{array}{cc}
 \frac{\pi}{2} & \frac{\pi}{2} \\
 0,72 & \frac{\pi}{2} \\
 1,46 & 1,23
 \end{array} \right) & \begin{array}{l} \text{Le premier janvier.} \\ \text{Je voudrais du pain.} \\ \text{Je voudrais une tasse de thé.} \end{array}
 \end{array}$$

Les couples de phrases dont les vecteurs forment un angle de  $\frac{\pi}{2}$  n'ont rien en commun, ils peuvent donc être exclus des résultats. Les meilleurs appariements sont ceux dont les vecteurs forment un angle minimum. Le résultat est la liste des couples de phrases ayant le plus petit angle à la fois en colonne et en ligne. En d'autres termes, le résultat est l'intersection entre les meilleurs appariements calculés en colonne et les meilleurs appariements calculés en ligne.

Le résultat de cet exemple est le bicorpus de phrases du tableau 4.

TAB. 4 – Bicorpus de phrases résultat d'appariement.

Angle	Français	Anglais
0,72	Je voudrais du pain.	I would like some bread.
1,23	Je voudrais une tasse de thé.	I would like a cup of tea.

Afin d'améliorer encore la précision, la même méthode est appliquée dans l'autre sens. Lors de ce

TAB. 5 – Précision de l'appariement sur l'échantillon pour chaque de décile du score.

Intervalles de scores	Nombre de couples valides (i)	Nombre total de couples (ii)	Précision (i)/(ii)
]1,45; 1,57]	4	27	0,15
]1,30; 1,45]	4	29	0,14
]1,15; 1,30]	6	17	0,35
]1,00; 1,15]	7	12	0,58
]0,85; 1,00]	5	7	0,71
]0,70; 0,85]	2	3	0,67
]0,55; 0,70]	3	3	1,00
]0,40; 0,55]	1	1	1,00
]0,25; 0,40]	1	1	1,00
TOTAL	32	100	0,32

second calcul, les phrases anglaises sont projetées dans l'espace vectoriel du français. Le résultat final est l'intersection des résultats des deux calculs.

## 6. Résultats d'expériences

Nous avons conduit une expérience de collecte et d'appariement de phrases courtes. La collecte de phrases est effectuée à partir du site de l'annuaire « Yahoo!<sup>1</sup> ». Le codage de caractères des documents aspirés est déterminé afin de recoder les documents en UTF8. Les balises sont supprimées à l'exception des balises marquant la fin d'un paragraphe qui sont remplacées par un saut de ligne. À cette étape du traitement, chaque document est composé d'un ensemble de phrases, une par ligne.

L'identification de la langue est effectuée au grain document par la méthode des profils de n-grammes de Cavnar et Trenkle citée plus haut.

Seuls les documents en anglais ou en français sont conservés. Nous avons fixé le nombre de phrases à collecter à 100 000 pour chaque langue.

Le dictionnaire français-anglais utilisé pour l'appariement est l'union des dictionnaires « Freelang<sup>2</sup> », « Wiktionary », « Omegawiki » et « Wikipedia interlangue<sup>3</sup> », légèrement modifié par nos soins. Le dictionnaire bilingue totalise 105 832 traductions entre 76 323 mots ou expressions en français et 62 694 mots ou expressions en anglais.

L'appariement permet d'obtenir 6 673 couples de phrases. Les angles obtenus sont compris entre 0,31 et 1,57. Afin d'estimer la qualité des couples de phrases obtenus, nous extrayons un échantillon aléatoire de 100 couples de phrases parmi les 6 673 que nous validons à la main.

Le tableau 5 détaille la précision de l'appariement en fonction du score des couples de phrases de l'échantillon. Les précisions sont calculées pour chaque décile du score.

Les données du tableau 5 permettent de calculer la corrélation entre le score d'un couple de phrases et le fait que ce couple soit une traduction. La valeur absolue de la corrélation est de 0,96. La précision baisse donc fidèlement avec le score d'appariement.

Nous pouvons noter que la précision est moins d'un tiers et les couples de score inférieur à 1,00 représentent 15% de l'échantillon. En extrapolant ce résultat de l'échantillon à l'ensemble des couples de phrases appariées, plus de 800 couples de phrases sur les 1 000 meilleurs sont des traductions exactes. Le tableau 6 fait la synthèse de l'expérience, du nombre de documents Web aspirés au nombre de couples de phrases traductions (par extrapolation).

Il faut noter que les phrases aspirés ne subissent aucune sélection visant à accroître le nombre de traductions. Obtenir 800 couples de phrases traductions à partir de 100 000 phrases par langue

<sup>1</sup> L'URL utilisée est <http://fr.yahoo.com>.

<sup>2</sup> Dictionnaire disponible à l'adresse <http://www.freelang.com>.

<sup>3</sup> Dictionnaires disponibles à l'adresse <http://www.dicts.info/udd1.php>.

TAB. 6 – Synthèse de l'expérience d'aspiration-appariement-validation.

	Quantité	Temps machine
Phrases aspirés en français	100 000	3 jours
Phrases aspirés en anglais	100 000	
Couples de phrases produite par appariement	6 673	10 heures
Couples de phrases traduction (précision supérieure ou égale à 0,8)	800	—

équivalent à un rendement de 8 pour 1 000. Ce rendement peut sembler faible, mais il est à considérer que les phrases alignées ont été obtenues sans aucune intervention humaine en moins de 4 jours. On pourrait donc escompter récolter 73 000 phrases en un an par extrapolation.

## 7. Conclusion

Nous avons introduit une méthode permettant de constituer des bicorpus de phrases courtes à partir du Web. Cette méthode utilise une représentation vectorielle des phrases, propice à leurs comparaisons interlingues. Nous avons montré que l'angle entre les vecteurs-phrases est un bon indice de la qualité des appariements produits. Ce qui distingue cette méthode est qu'elle ne se base pas sur des documents préalablement alignés. Elle est donc utile pour constituer des corpus de phrases dans des couples de langues peu documentées. Une expérience sur le couple français-anglais montre que cette méthode permet d'obtenir 800 phrases traductions à partir de 100 000 phrases aspirées par langue, soit un rendement de 8 pour 1 000. Le temps nécessaire à la collecte et à l'appariement auras durée moins de 4 jours. Ces résultats permettent de valider la méthode d'appariement et nous encourageant à étendre notre étude à d'autres couples de langues. En particulier nous espérons obtenir des résultats similaires à l'aide de cette méthode sur des couples de langues peu documentés.

## Bibliographie

1. W.B. Cavnar et J.M. Trenkle. N-gram-based text categorization. *Ann Arbor MI*, 48113:4001, 1994.
2. S.T. Dumais, T.A. Letsche, M.L. Littman, et T.K. Landauer. Automatic cross-language retrieval using latent semantic indexing. In *AAAI Spring Symposium on Cross-Language Text and Speech Retrieval*, 1997.
3. W.A. Gale et K.W. Church. A program for aligning sentences in bilingual corpora. *Computational Linguistics*, 19(1):75–102, 1993.
4. T.K. Landauer, P.W. Foltz, et D. Laham. An Introduction to Latent Semantic Analysis. *DISCOURSE PROCESSES*, 25:259–284, 1998.
5. X. Ma et M. Liberman. Bits: A method for bilingual text search over the web. In *Machine Translation Summit VII*, 1999.
6. J.Y. Nie, M. Simard, P. Isabelle, et R. Durand. Cross-language information retrieval based on parallel texts and automatic mining of parallel texts from the web. In *Proceedings of the 22nd Annual International ACM SIGIR Conference on Research and Development in Information Retrieval*, pages 74–81, 1999.
7. P. Resnik, Laboratory for Language, et Media Processing. Parallel strands: A preliminary investigation into mining the web for bilingual text. *Lecture notes in computer science*, pages 72–82, 1998.
8. G. Salton, A. Wong, et CS Yang. A vector space model for automatic indexing. *Communications of the ACM*, 18(11):613–620, 1975.