

Zimina Maria

*Université de la Sorbonne Nouvelle - Paris 3*

*Equipe LEXICO (SYLED)*

---

# **Alignement de textes bilingues par classification ascendante hiérarchique**

*Aligner* deux textes consiste à mettre en relation des unités textuelles qui se correspondent :

1. The Declaration was adopted at once.

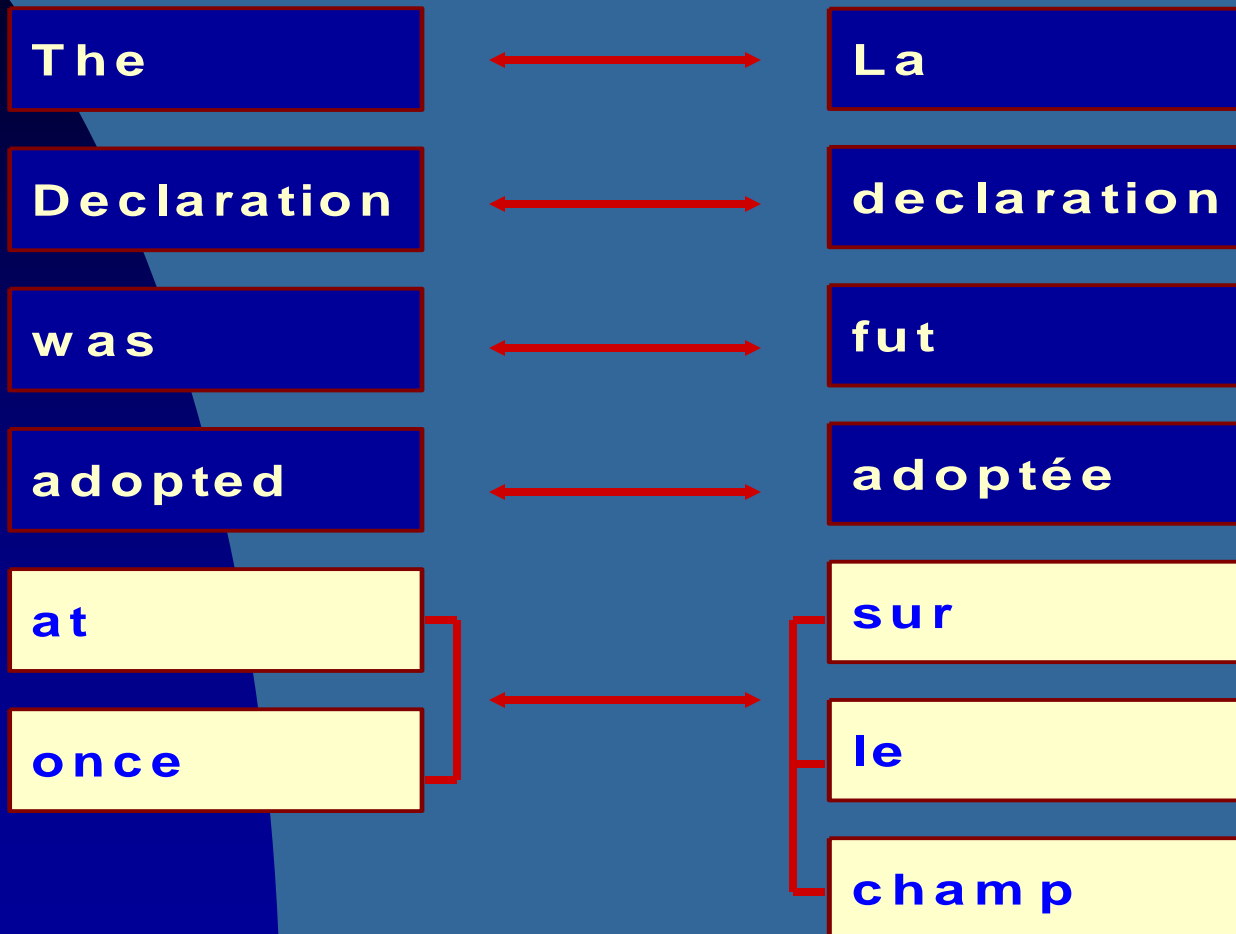
1. La déclaration fut adoptée sur le champ.

2. All the states pronounced unanimately for the abolition of the former treaty as it was judged out of date.

2. Tous les Etats se prononcèrent de manière unanime pour le retrait du Traité.

3. Il fut jugé qu'il ne répondait plus aux conjonctures actuelles.

Les unités en correspondance se positionnent à plusieurs niveaux : *paragraphes, phrases, syntagmes, mots.*



*Domaines d'application* : traduction, terminologie, lexicographie, industries de la langue, publication etc...

- L'alignement permet de réutiliser les ressources de traductions existantes.
- Il facilite la création des lexiques multilingues.
- Les techniques d'alignement aident à la maintenance de documentation technique dans des langues différentes.

Des méthodes analogues sont également utilisées dans l'appariement d'un texte et de sa traduction phonétique ou d'un texte et de son enregistrement.

L 'alignement permet d 'élaborer les *concordances multilingues*.

RECHERCHER ⇒

countervail

One of the strongest weapons the U.S. has been using to restrict entry of foreign goods onto its market has been **countervail duties**

Les **droits compensateurs** constituent l'une des principales armes que les États-Unis ont utilisée pour restreindre les importations de produits étrangers sur leur marché.

Un *bi-texte* est un ensemble de deux textes et des liens qui relient les unités en correspondance.

- Pour générer automatiquement un *bi-texte* on a recours aux *algorithmes d'alignement*.
- La plupart des algorithmes font appel aux méthodes statistiques.
- Au stade actuel, ces algorithmes ne sont pas encore capables de rendre explicites toutes les correspondances de traduction.

Méthodes d'alignement par *classification ascendante hiérarchique* de formes graphiques et segments répétés.



**Corpus :**

Convention de sauvegarde des Droits  
de l'homme et des libertés  
fondamentales

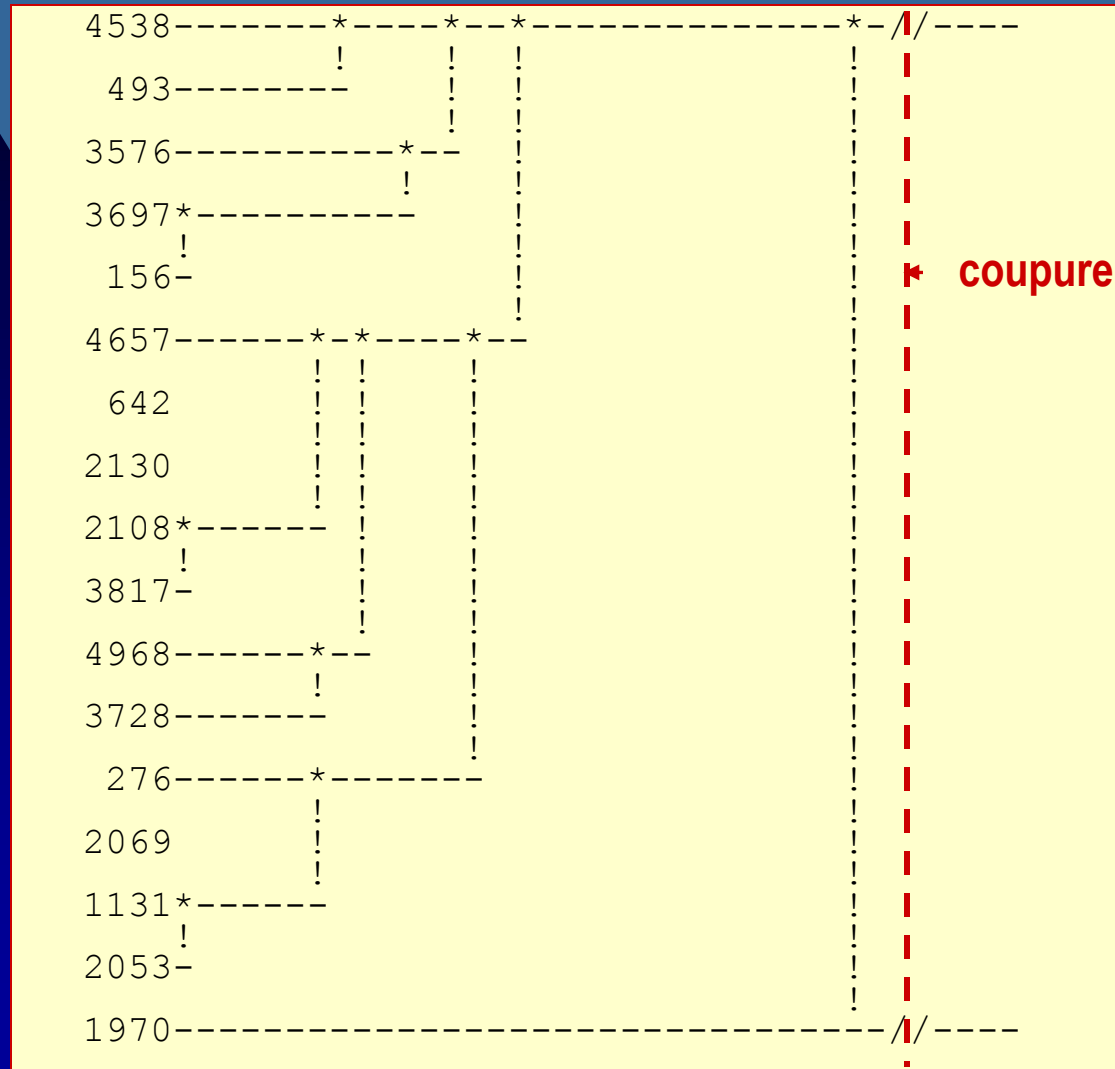
■ Corpus anglais  
273 685 occurrences

■ Corpus français  
285 961 occurrences

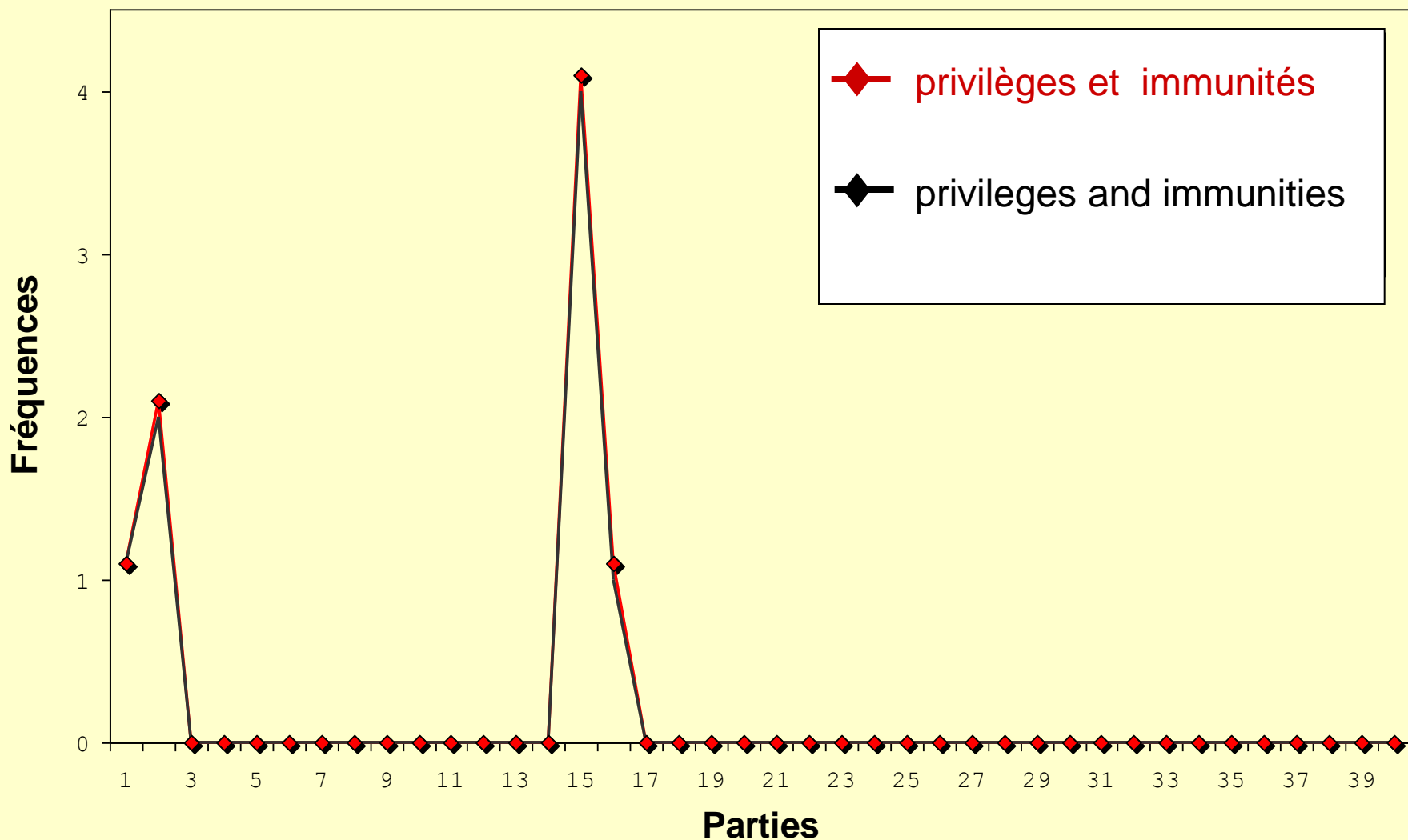




La *coupure* du dendrogramme permet de considérer une classification résumée aux classes inférieures de la hiérarchie qui regroupent des éléments extrêmement proches.



*Profils graphiques* des segments répétés agrégés dans les mêmes classes.



La classification a dégagé un nombre important de classes qui regroupent des *individus - correspondances de traduction*.

CLASSE 2 6 9 2

*égalité des* ■ = français  
*equality of* ■ = anglais

CLASSE 2 5 2 6

*replies to*  
*réponses à*

CLASSE 2 2 6 7

*speculate as to*  
*spéculer sur*

CLASSE 2 4 2 2

*37 above*  
*37 ci*

CLASSE 2 5 9 4

*differences in*  
*différences entre*

CLASSE 2 1 5 8

*lecture des*  
*reading out*

CLASSE 2 7 0 8

*accepts that*  
*admet que*

CLASSE 2 6 5 6

*né en*  
*born in*

CLASSE 2 5 2 3

*côté de la*  
*côté de*  
*side of*  
*side of the*

CLASSE 2 3 7 1

*verdict d*  
*verdict of*

CLASSE 2 6 6 9

*contributed to*  
*contribué à*

CLASSE 2 6 5 7

*considerations of*  
*considérations d*

# Filtrage des résultats

- apparier les formes et segments avec des *fréquences générales homogènes* :

■ = français  
■ = anglais

CLASSE 2622

248	<i>puni de</i>	F = 11
1414	<i>punishable by</i>	F = 12
1402	<i>increased by</i>	F = 5



CLASSE 2571

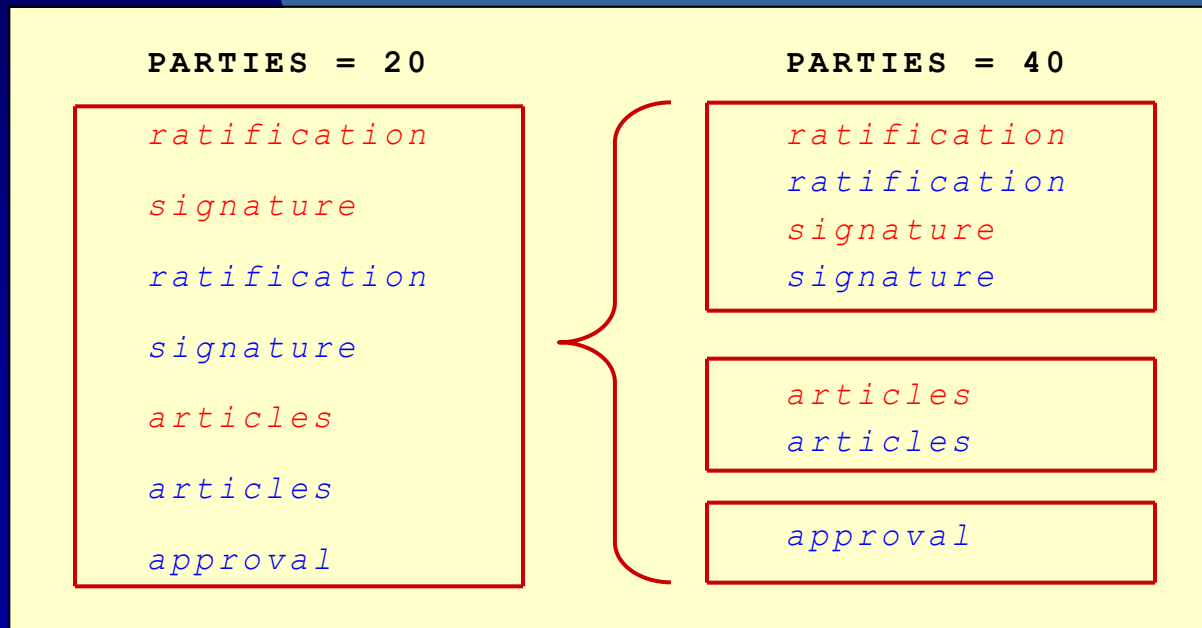
772	<i>résine de cannabis</i>	F = 12
1618	<i>cannabis resin</i>	F = 12
1599	<i>schedule 3</i>	F = 15
1895	<i>ought to</i>	F = 5



# Filtrage des résultats

<i>signature</i>	14	16	5	0	1	0	0	0	2	0	1	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	1	0	0	0	0	1	0	0	1	0	0	1	0	0							
<i>signature</i>	30	5	1	0	2	1	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	1	0	0	1	1	1	0	0	0	0	0	0	0	0	0						
<i>signature</i>	13	16	5	0	0	0	0	0	1	0	1	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0			
<i>signature</i>	29	5	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0		
<i>articles</i>	15	21	1	0	0	1	0	0	0	1	1	0	0	1	0	1	1	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	
<i>articles</i>	36	1	1	0	1	1	1	1	1	1	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	1	0	1	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	
<i>articles</i>	13	21	1	0	0	1	0	0	0	0	1	0	0	1	0	1	2	0	1	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0
<i>articles</i>	34	1	1	0	1	1	2	1	0	0	1	0	0	1	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0

- tenir compte de l'influence du *découpage* du corpus :



■ = *français*

■ = *anglais*

## Pistes de développement

- utiliser les classes d'individus agrégés pour aligner les phrases correspondantes
- intégrer la classification aux systèmes d'alignement des phrases pour *augmenter la résolution* jusqu'au niveau des mots ou des syntagmes

L'intérêt de cette méthode réside dans son degré de *flexibilité* : les correspondances de traduction sont recherchées librement dans l'ensemble de deux textes.

Cette recherche peut être automatisée.  
Elle est peu coûteuse en calcul.