

Journée d'études
« Outils de traitement de corpus textuels développés à Paris-Est »
18/06/2014 – Créteil

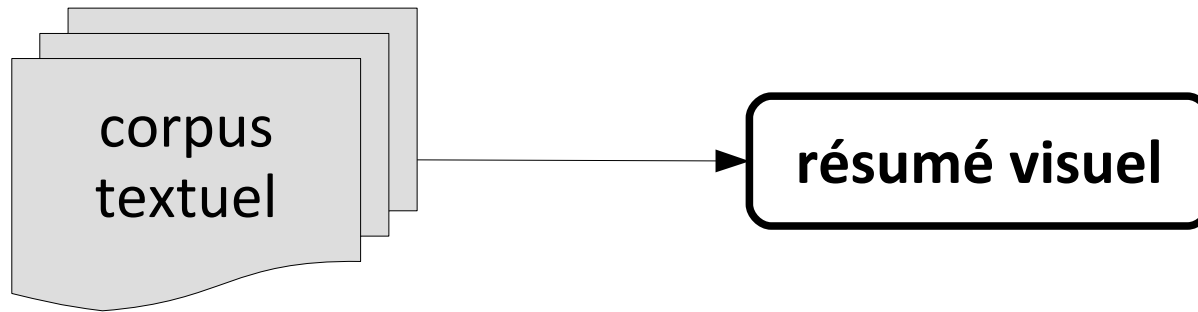
Visualisation et analyse de textes à partir d'arbres de mots avec TreeCloud

Philippe Gambette

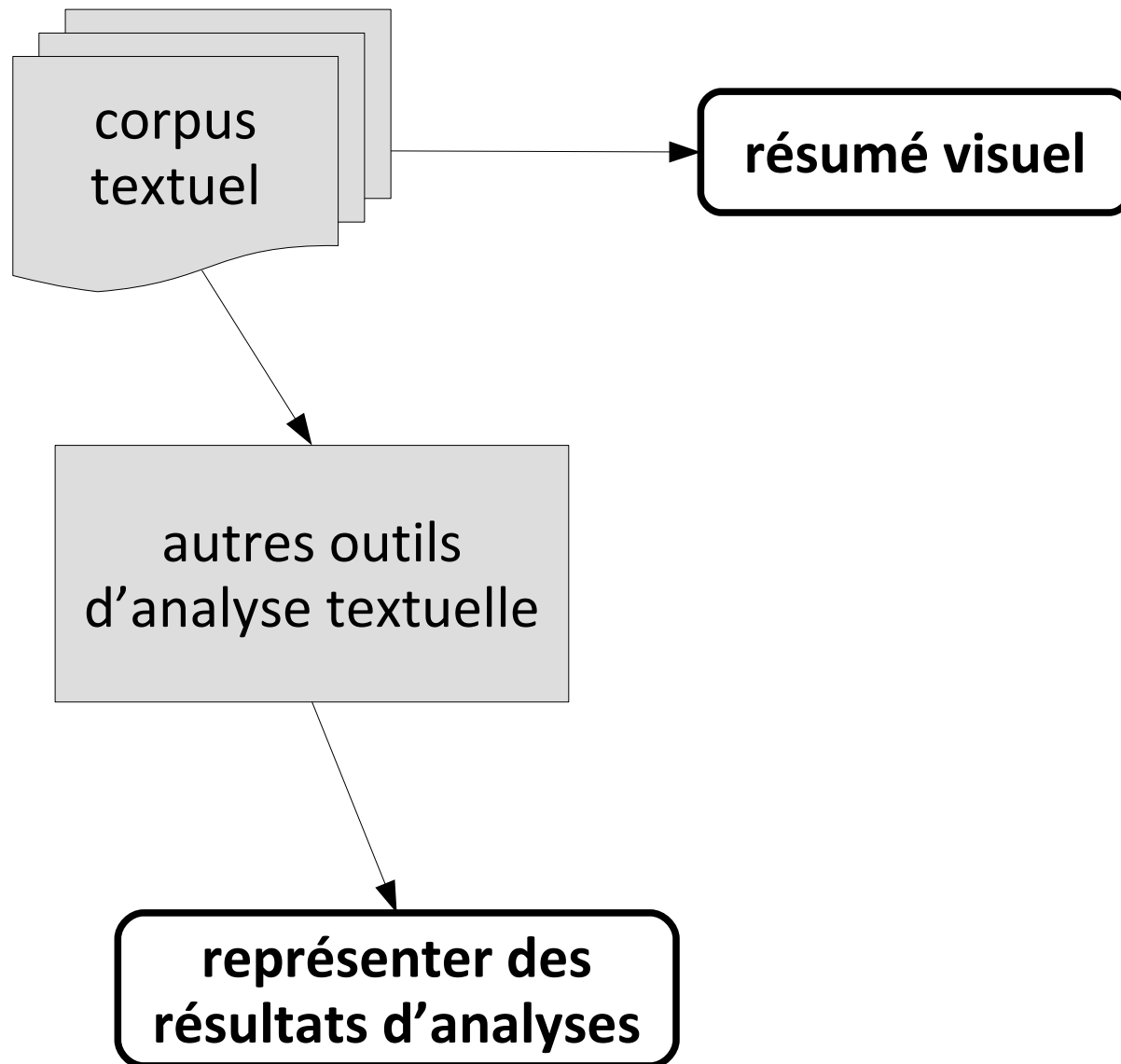
LIGM
Université Paris-Est
Marne-la-Vallée



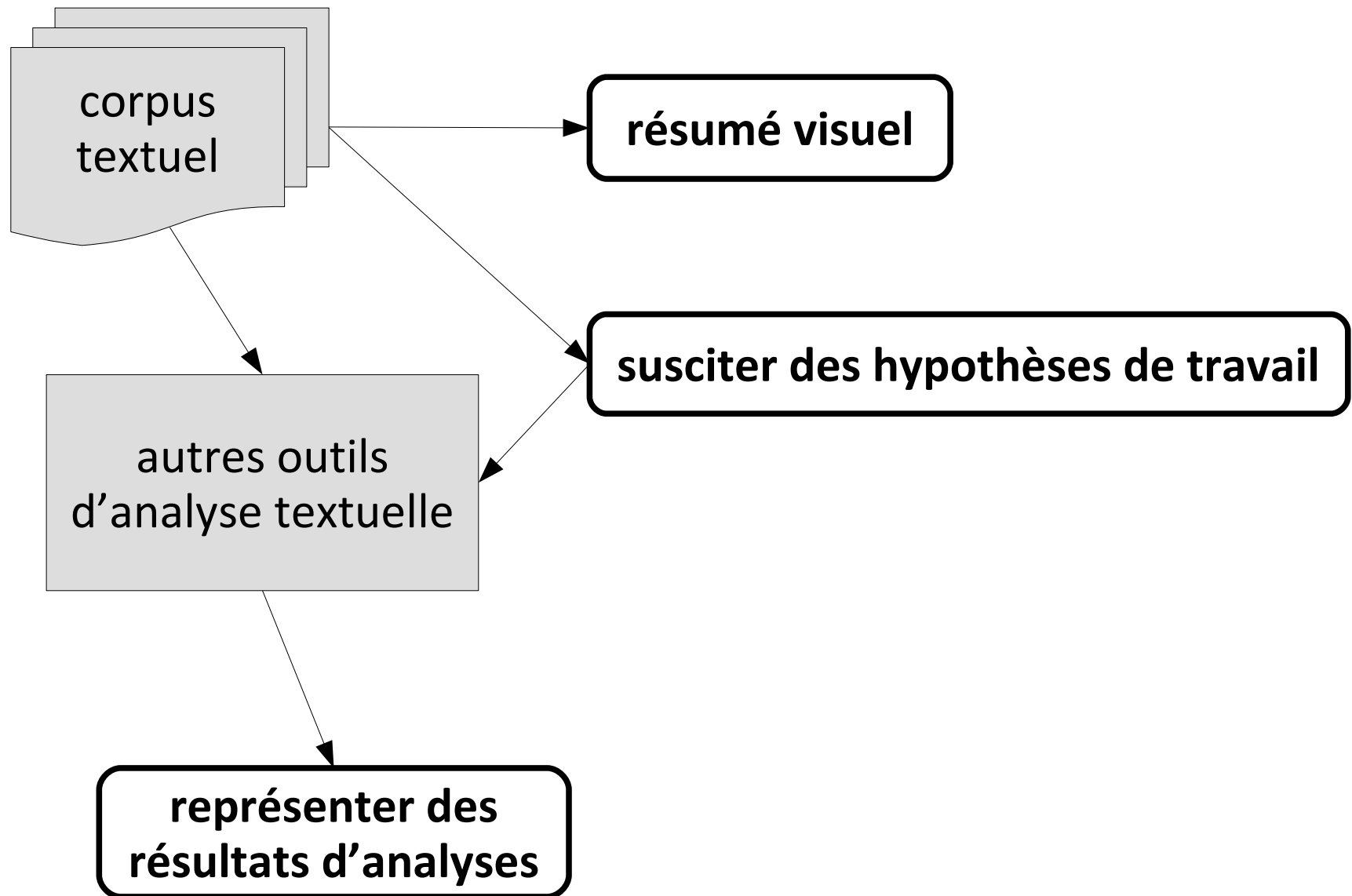
Le « nuage arboré », pour quoi faire ?



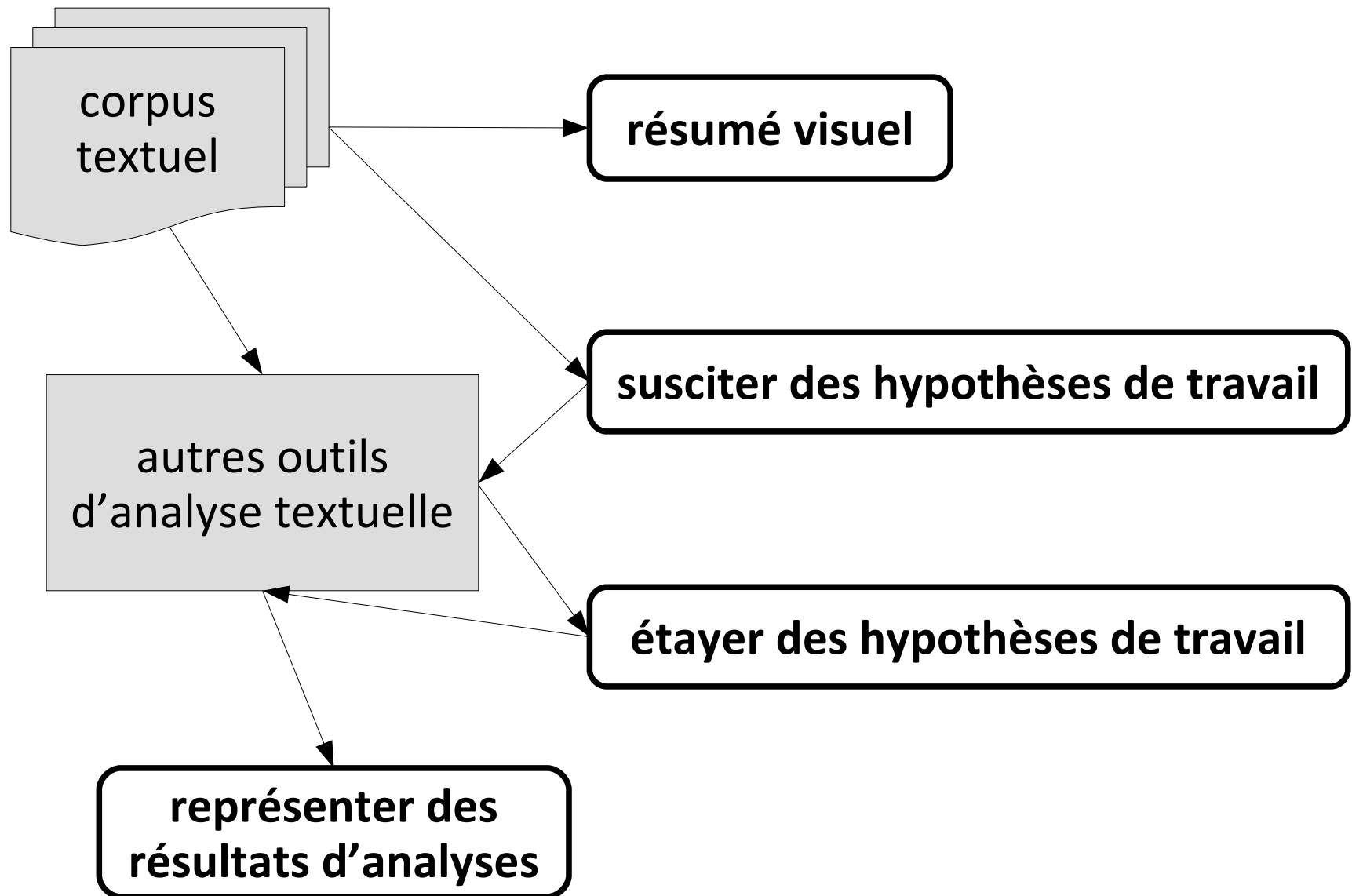
Le « nuage arboré », pour quoi faire ?



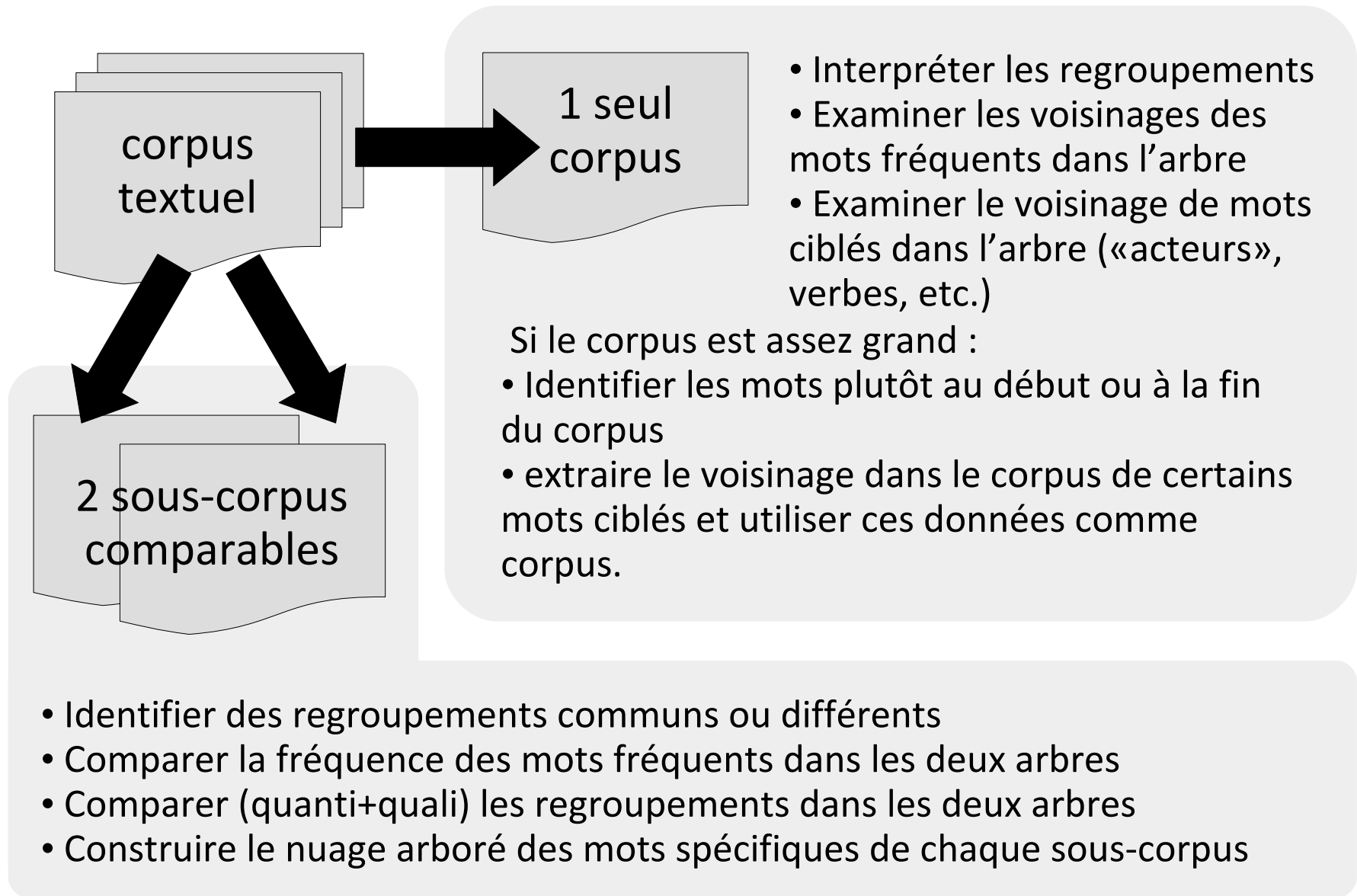
Le « nuage arboré », pour quoi faire ?



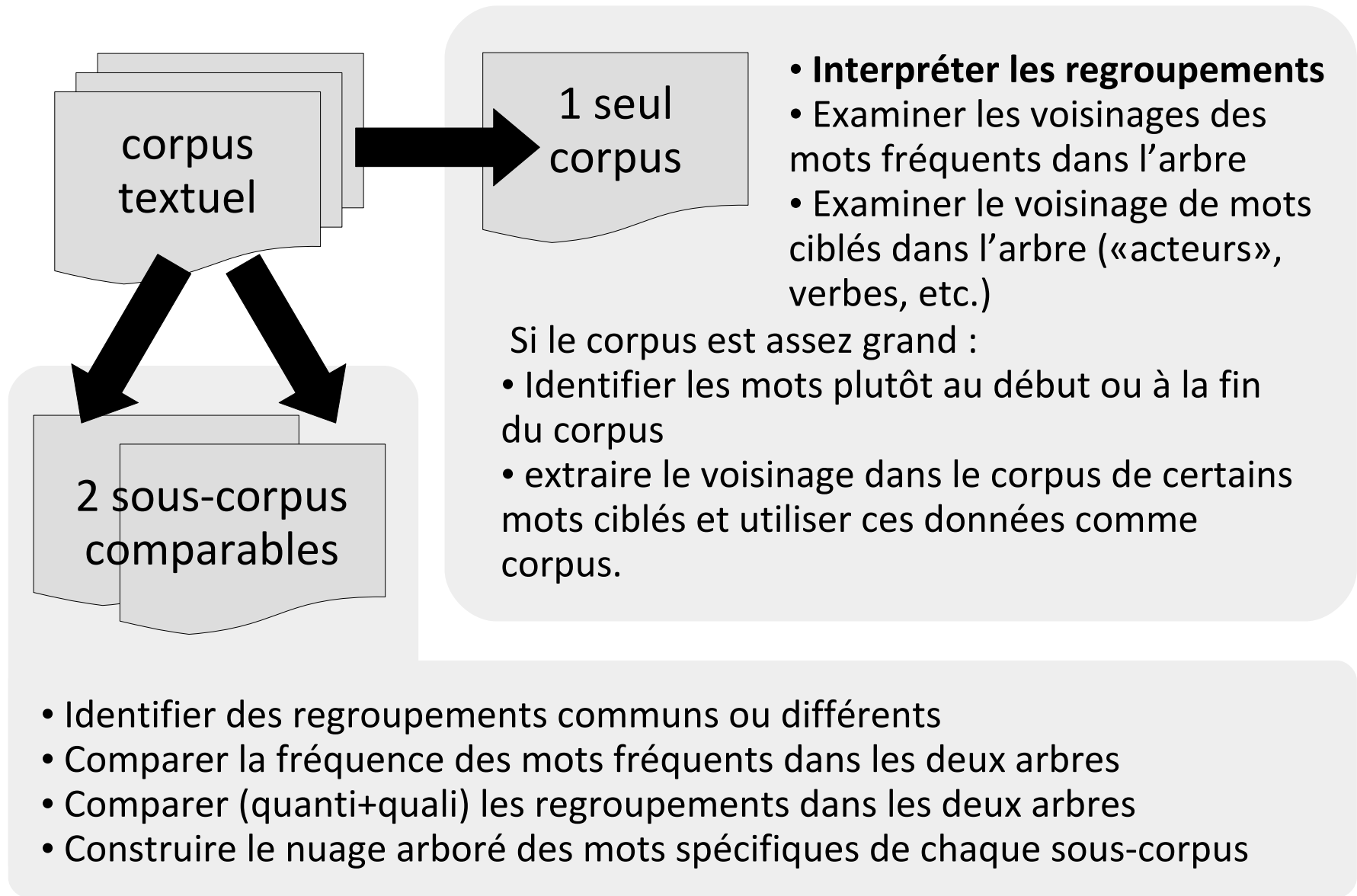
Le « nuage arboré », pour quoi faire ?



Exploration de corpus avec TreeCloud



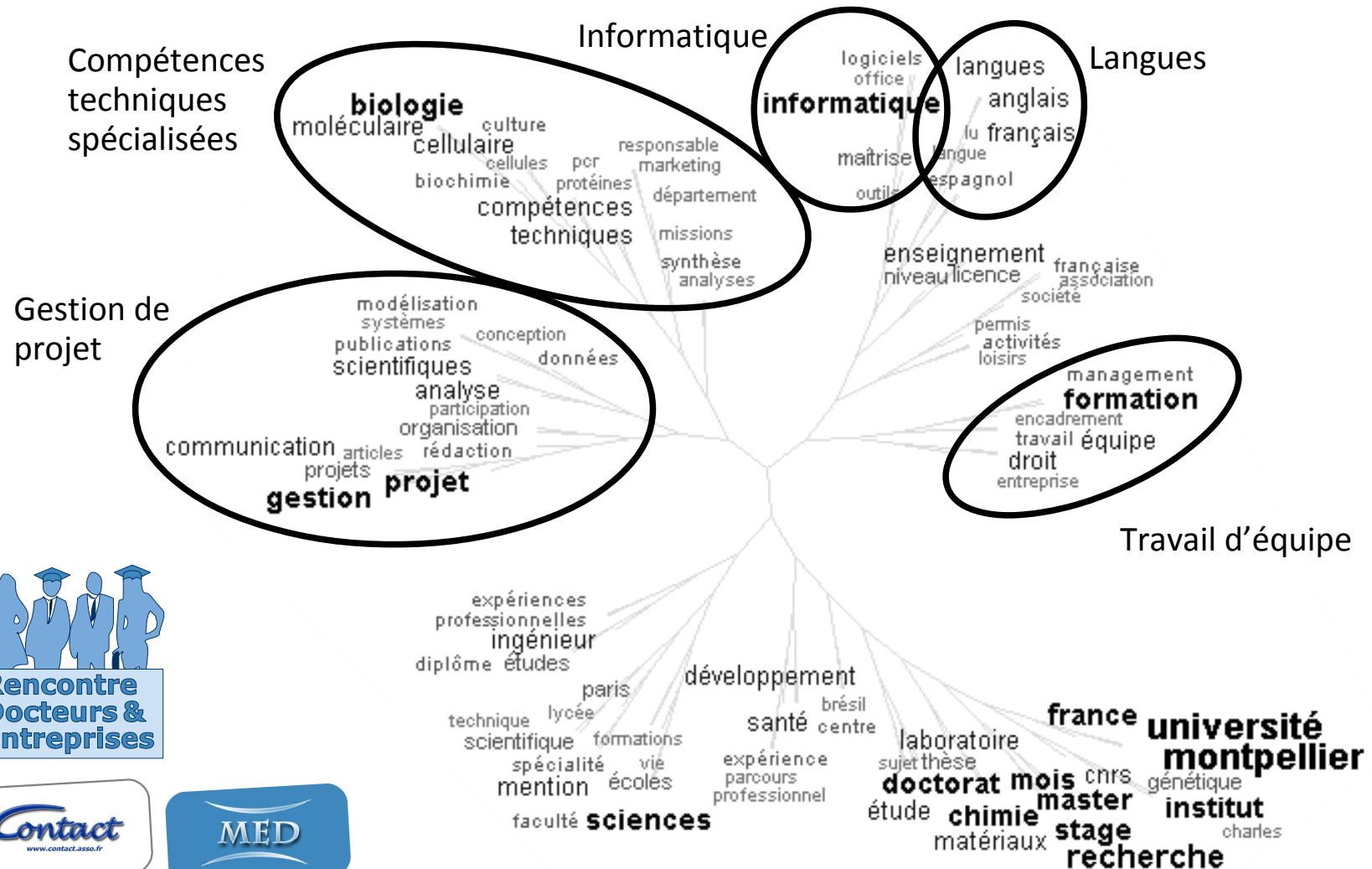
Exploration de corpus avec TreeCloud



Méthode : interpréter les regroupements

Dessiner des « patates »

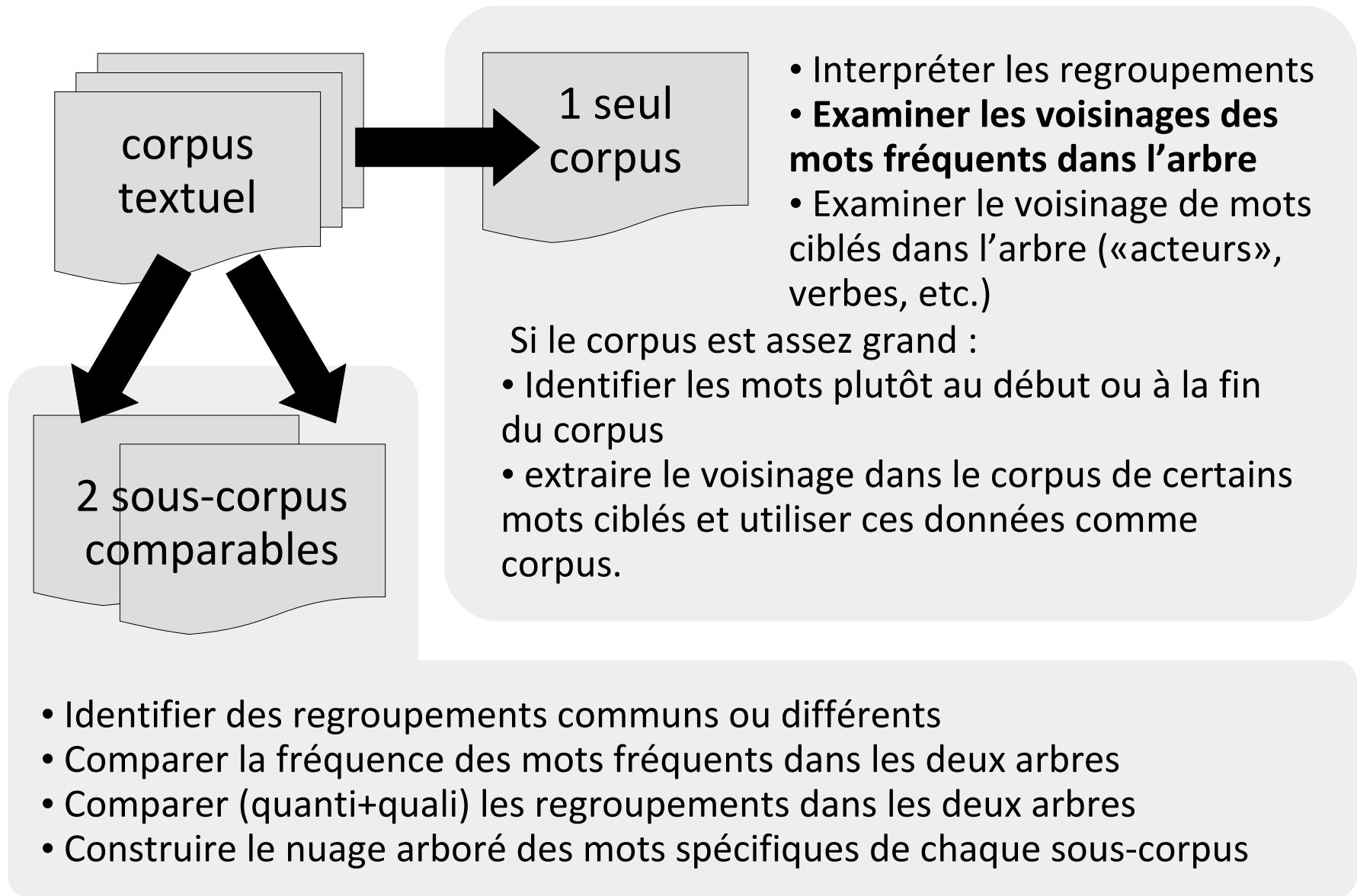
Corpus : une centaine de CV soumis à une rencontre docteurs-entreprises



Rencontre
Docteurs &
Entreprises

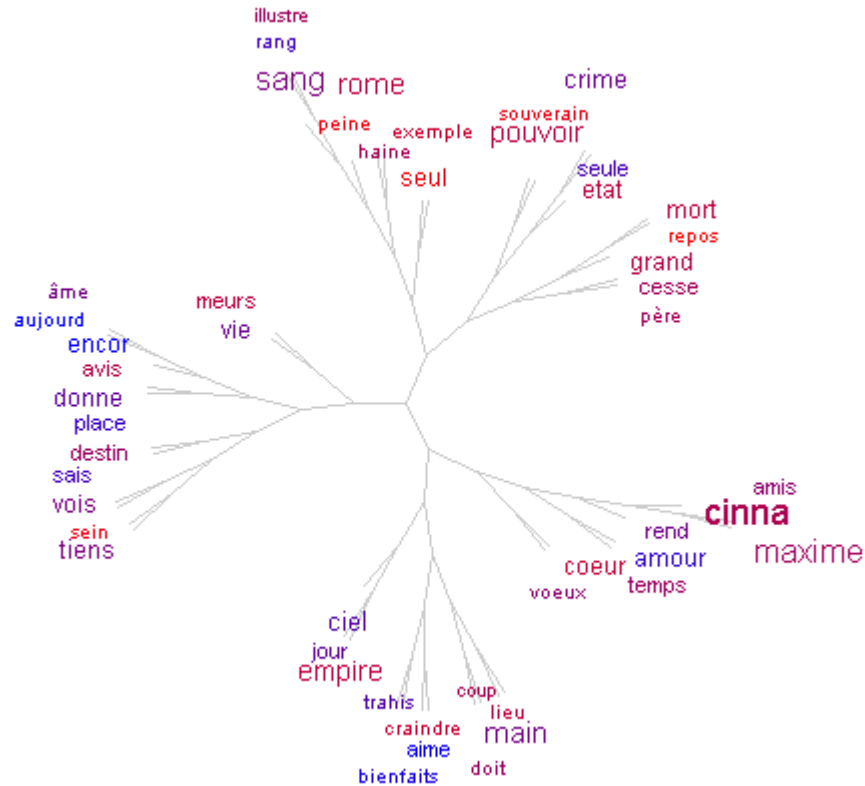


Exploration de corpus avec TreeCloud



Méthode : voisinage des mots fréquents

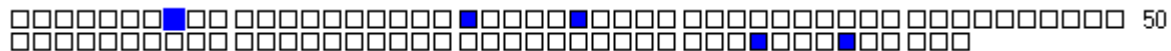
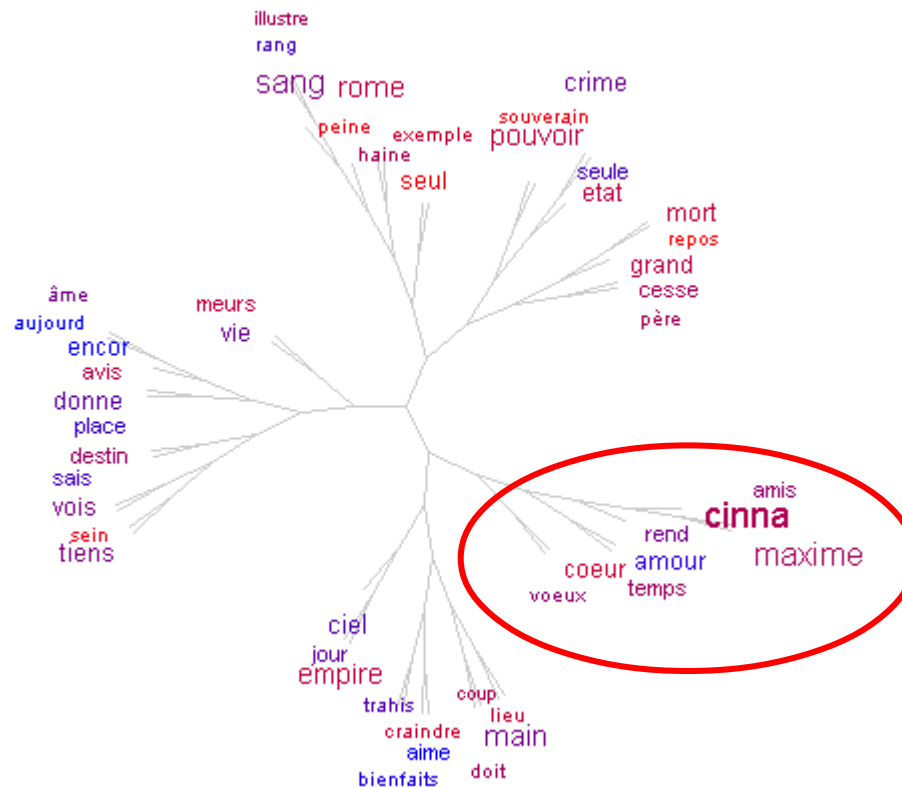
Amstutz & Gambette,
JADT 2010



Nuage arboré des 50 mots les plus fréquents des paroles d'Auguste dans Cinna

Méthode : voisinage des mots fréquents

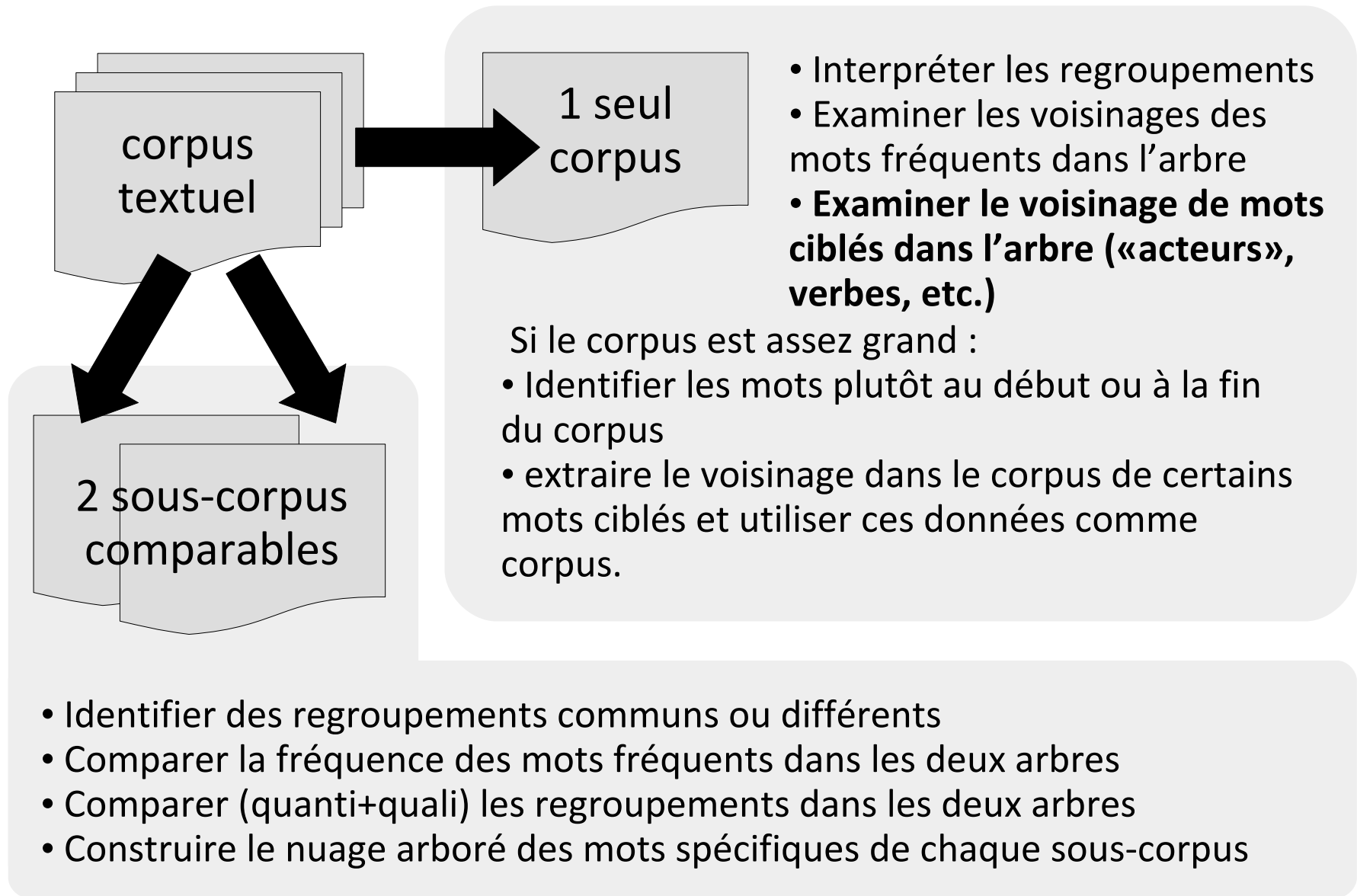
Amstutz & Gambette,
JADT 2010



Carte des sections Lexico3 et contextes de « amis » dans les paroles d'Auguste dans Cinna.

1. Voilà, mes chers **amis**, ce qui me met en peine.
2. Quoi ! mes plus chers **amis** ! quoi ! Cinna ! quoi ! Maxime !
3. Reprenez le pouvoir que vous m'avez commis, Si donnant des sujets il ôte les **amis**
4. Soyons **amis**, Cinna, c'est moi qui t'en convie
5. Il nous a trahis tous ; mais ce qu'il a commis Vous conserve innocents, et me rend mes **amis**.

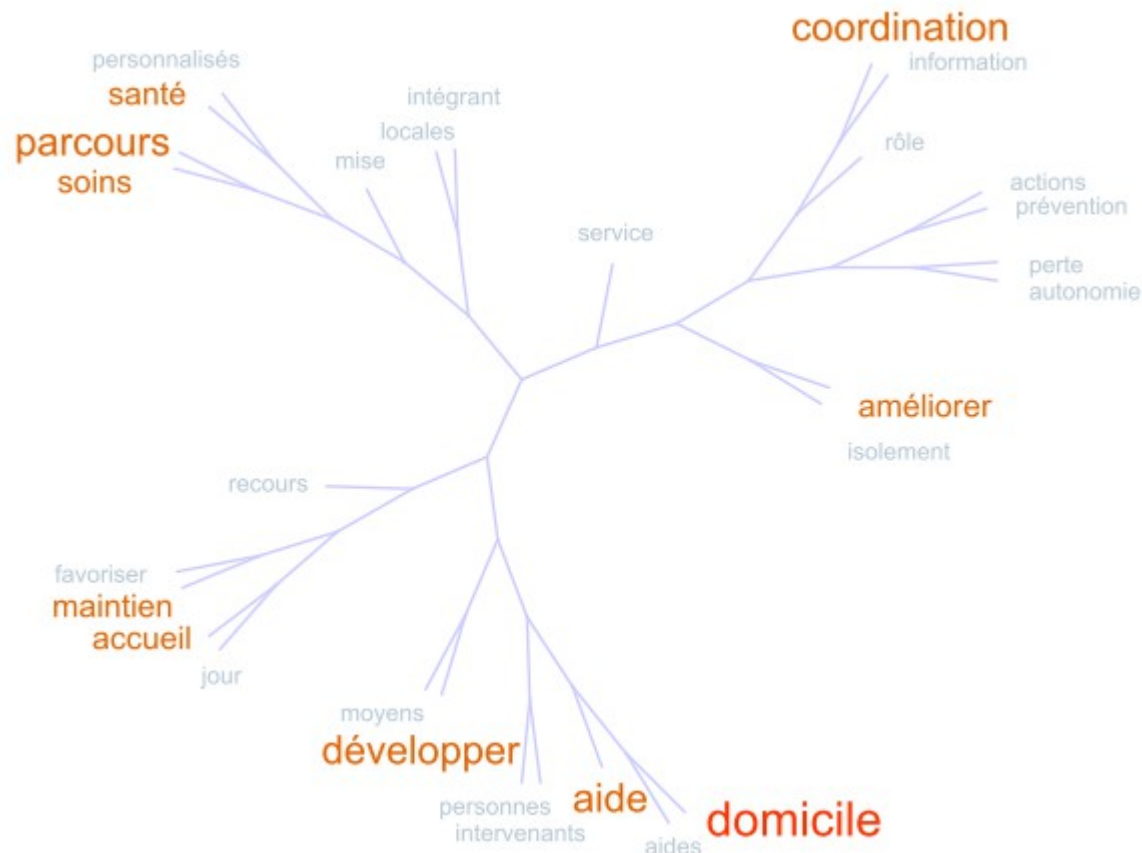
Exploration de corpus avec TreeCloud



Méthode : voisinage des verbes

Corpus : réponses à des questions ouvertes à des professionnels de la santé sur le parcours de santé des personnes âgées dans les Alpes de Haute-Provence

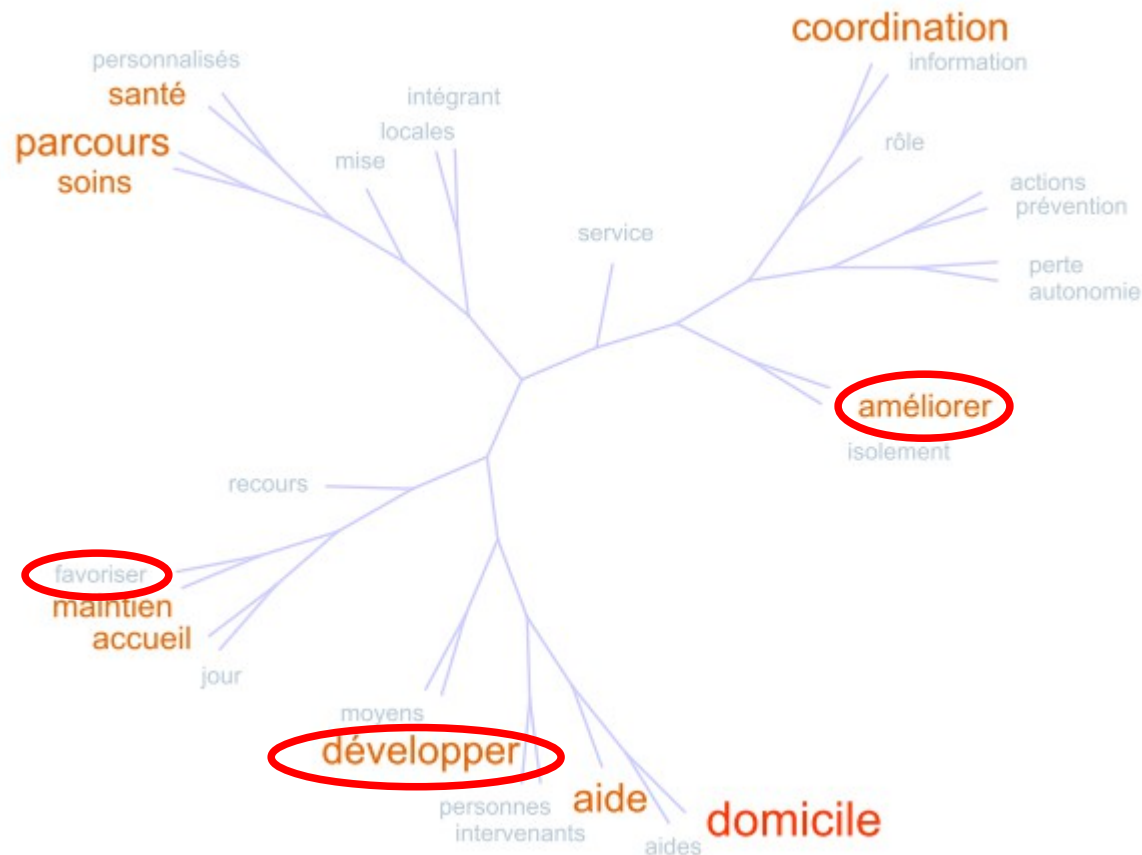
Suggestions d'améliorations



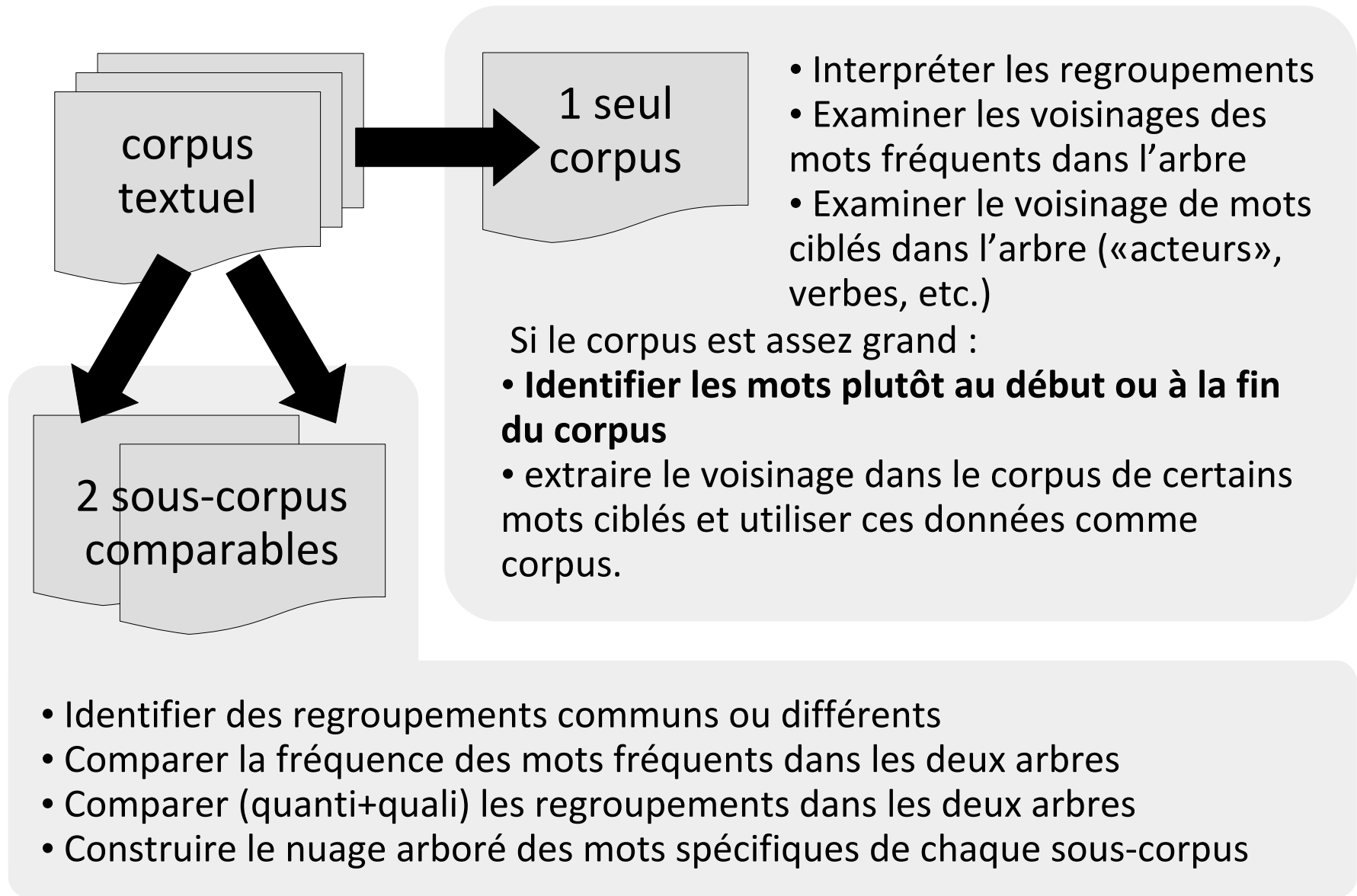
Méthode : voisinage des verbes

Corpus : réponses à des questions ouvertes à des professionnels de la santé sur le parcours de santé des personnes âgées dans les Alpes de Haute-Provence

Suggestions d'améliorations



Exploration de corpus avec TreeCloud



Exploration de corpus avec TreeCloud

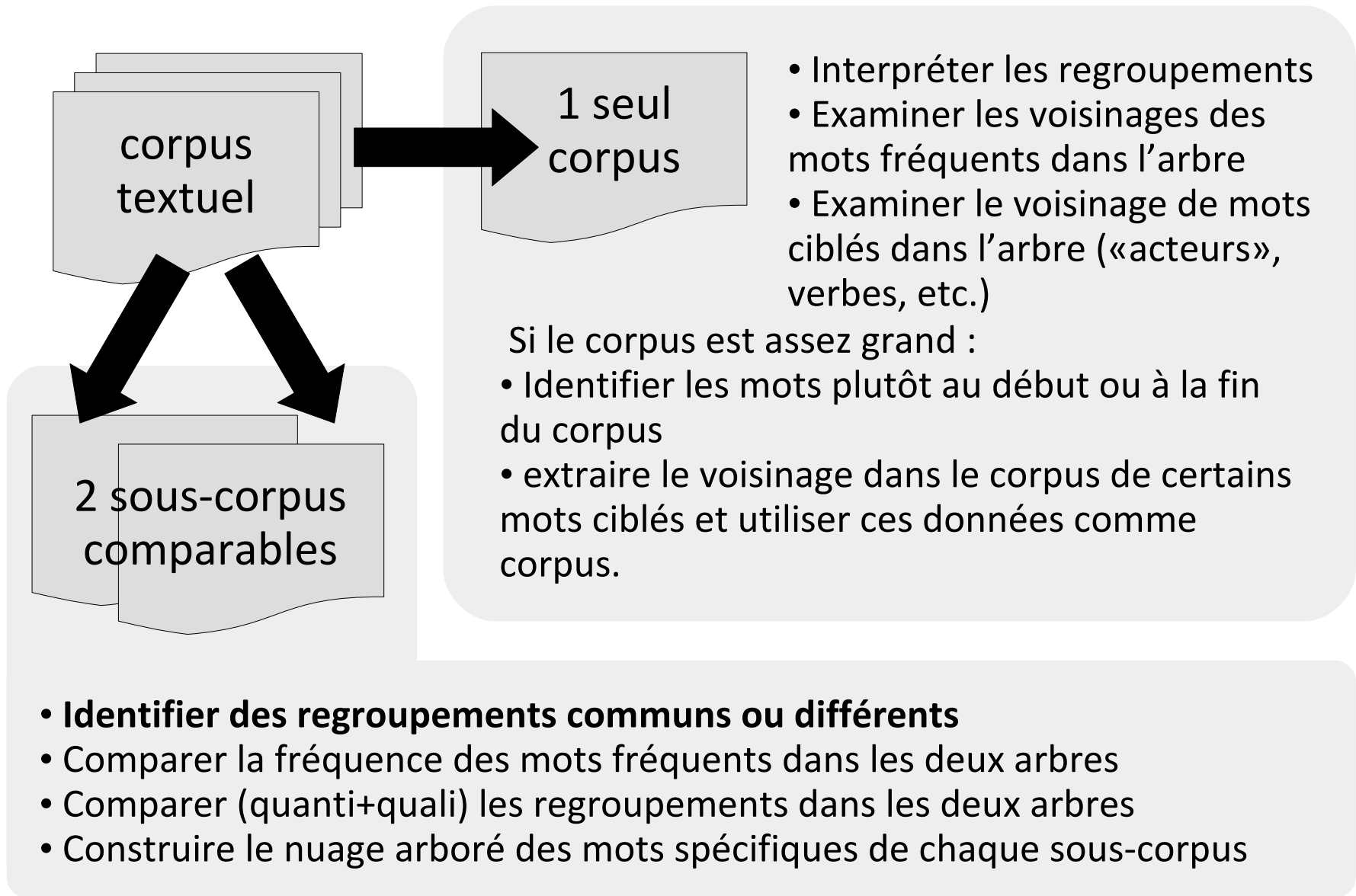


Illustration sur le corpus Mediator

Comparer les articles d'agences et articles de journalistes

Corpus : 595 articles d'agences contre 1496 articles de journalistes de 2011 évoquant l'affaire du Mediator dans la presse française.

Ensemble des articles

Gambette & Martinez,
Texto!, 2013

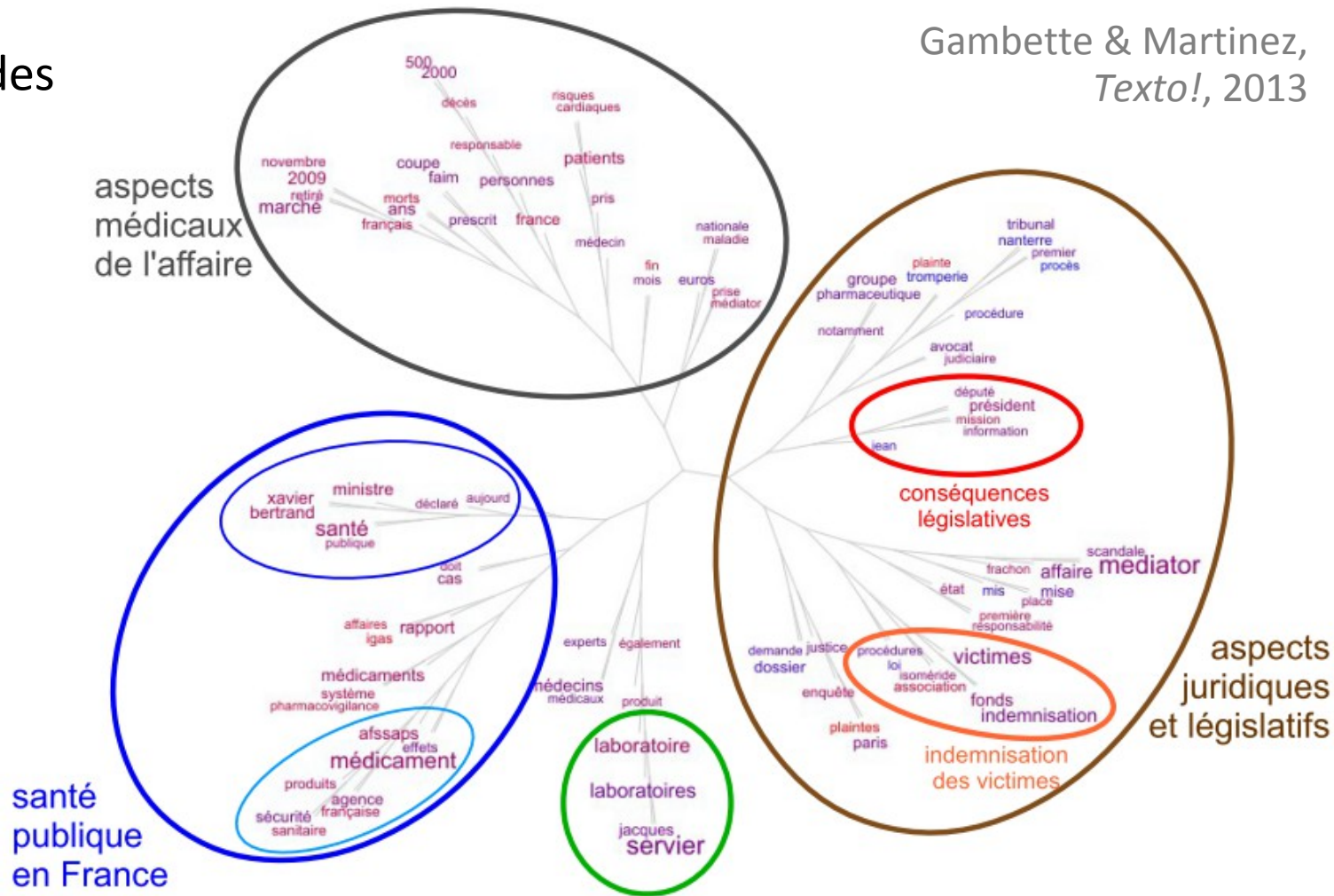
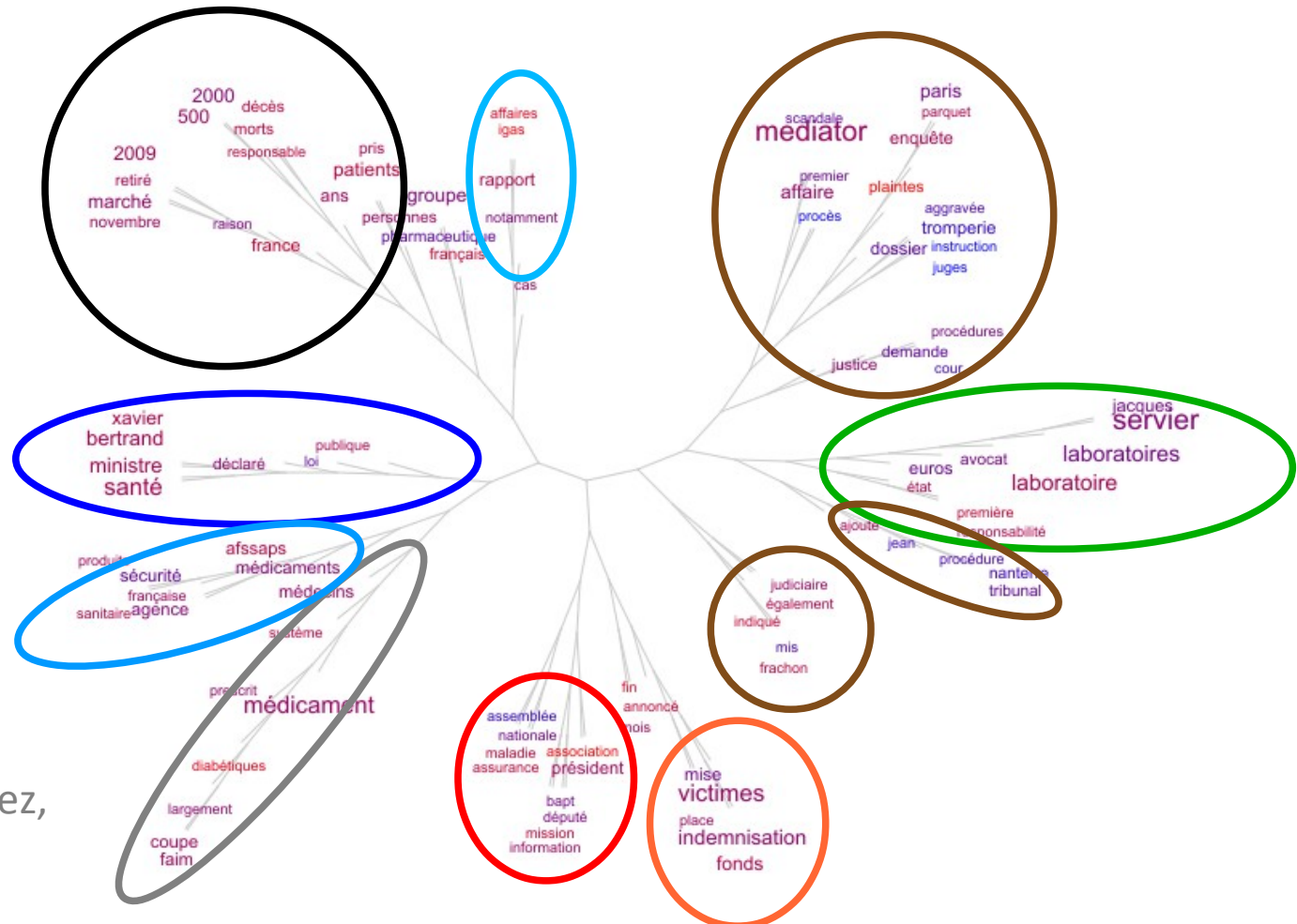


Illustration sur le corpus Mediator

Comparer les articles d'agences et articles de journalistes

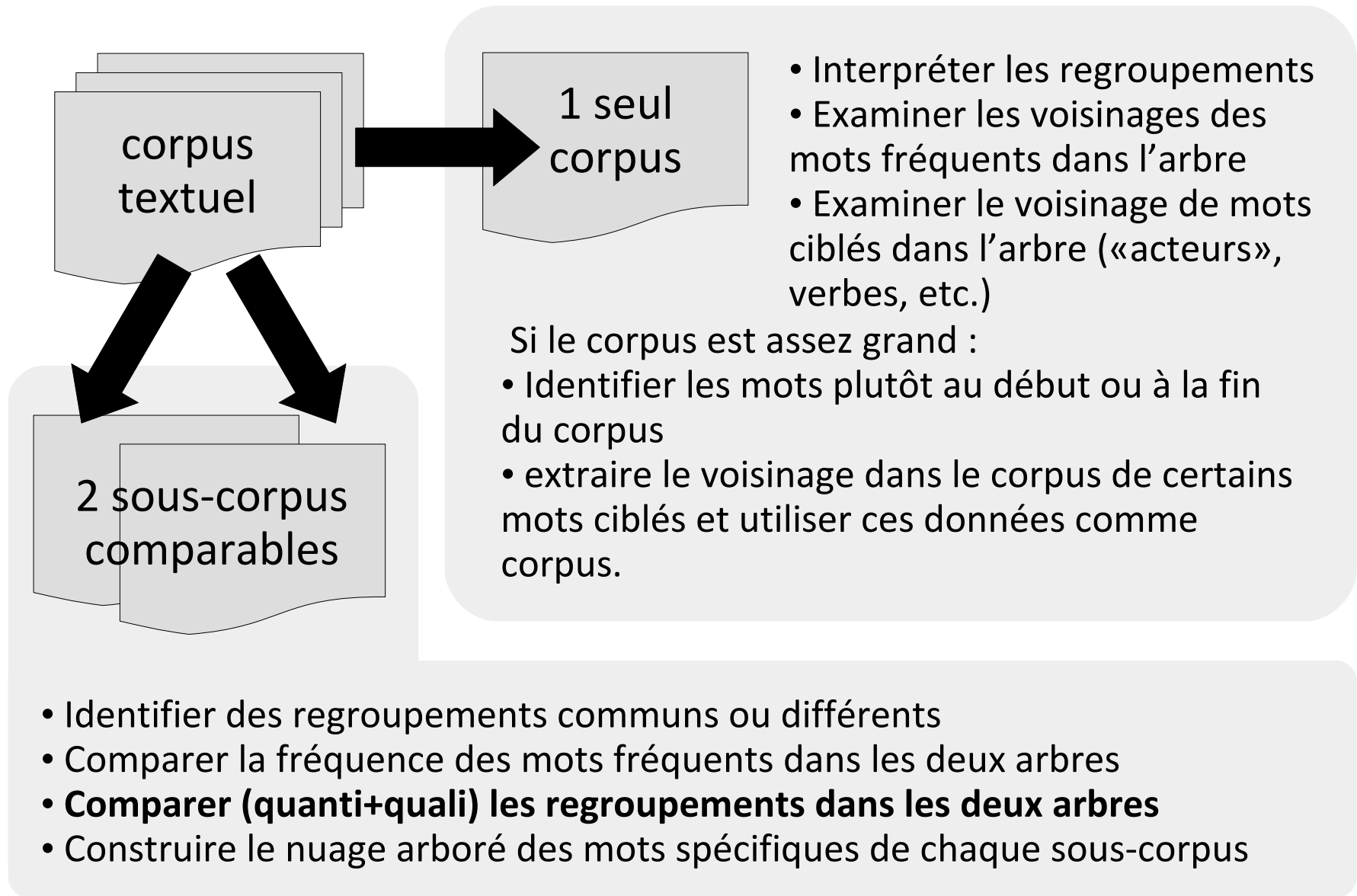
Corpus : 595 articles d'agences contre 1496 articles de journalistes de 2011 évoquant l'affaire du Mediator dans la presse française.

Articles
d'agences

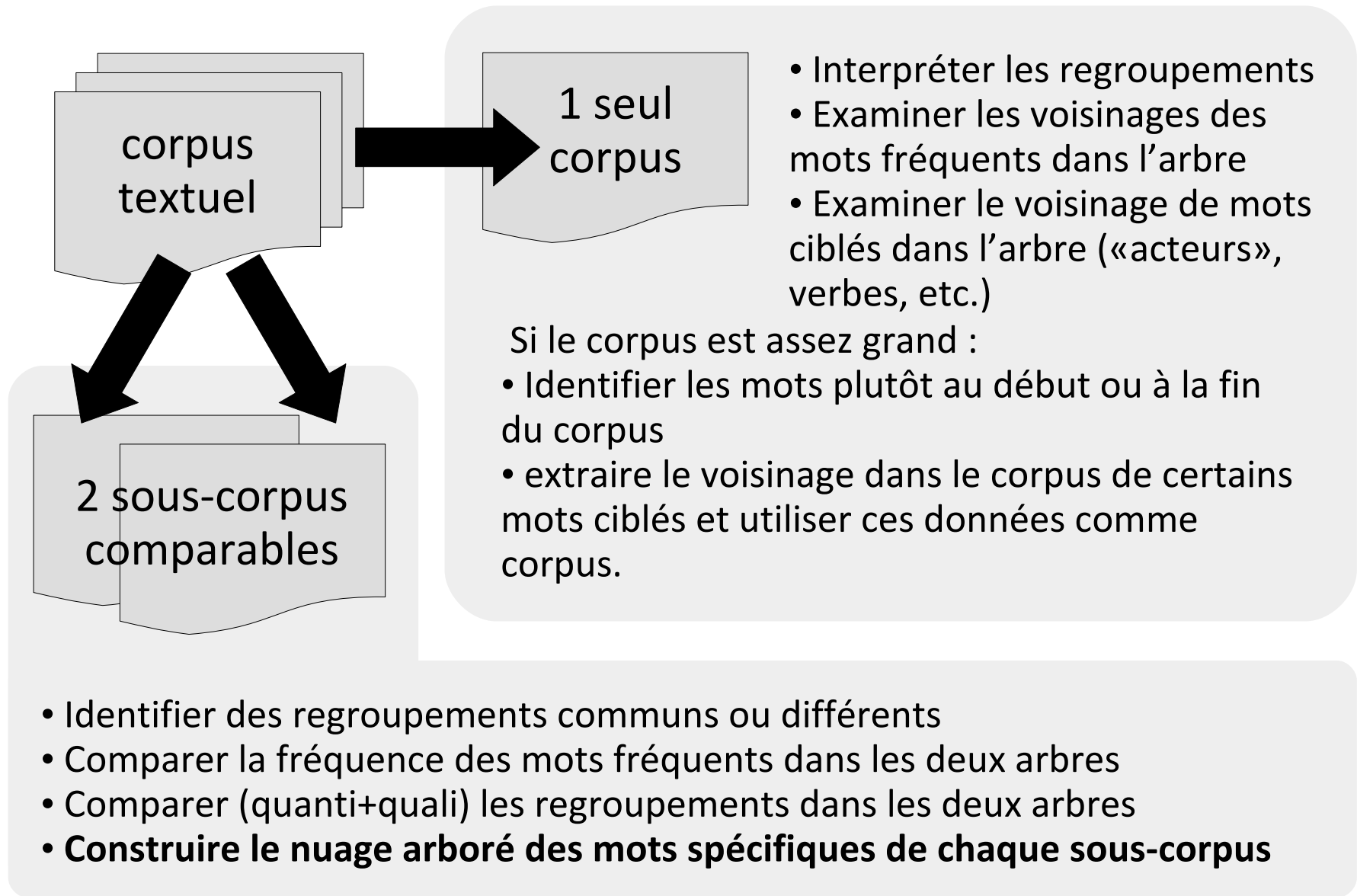


Gambette & Martinez,
Texto!, 2013

Exploration de corpus avec TreeCloud

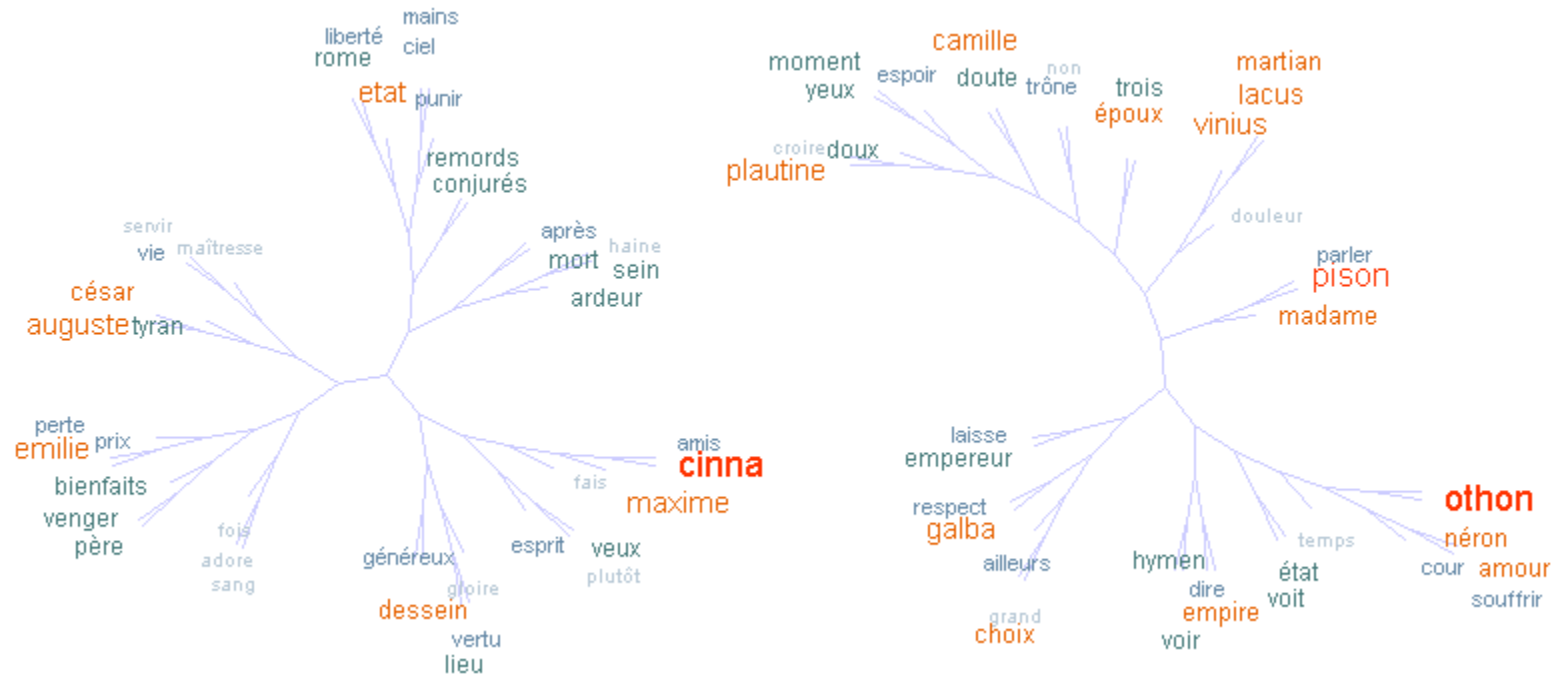


Exploration de corpus avec TreeCloud



Méthode : comparaison des spécifiques

Amstutz & Gambette,
JADT 2010



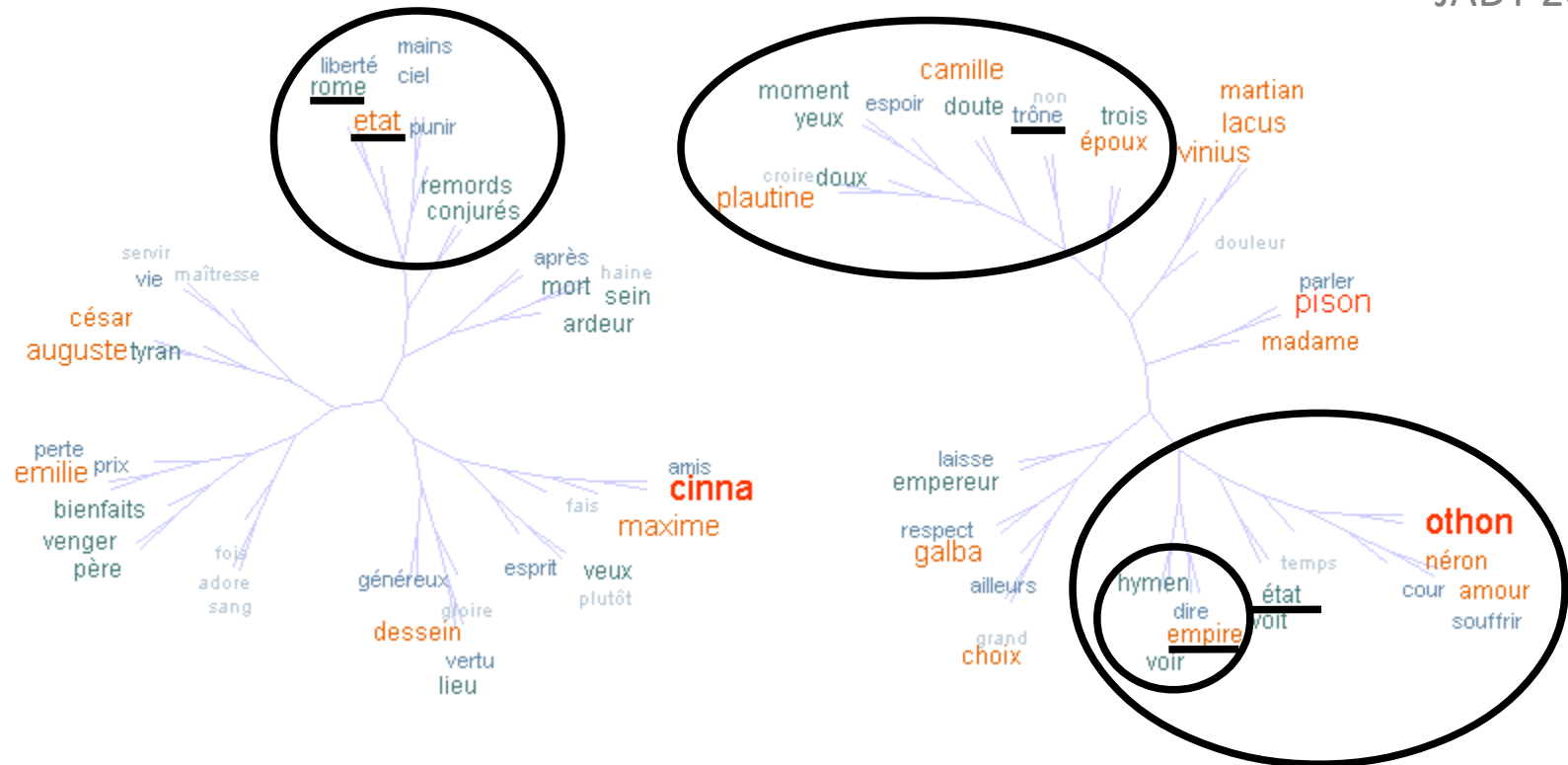
*Nuages arborés des **mots spécifiques** de Cinna et Othon, dimensionnés et colorés d'après leur spécificité calculée dans Lexico3.*



Quels moyens au service de la cause politique ?

Méthode : comparaison des spécifiques

Amstutz & Gambette,
JADT 2010

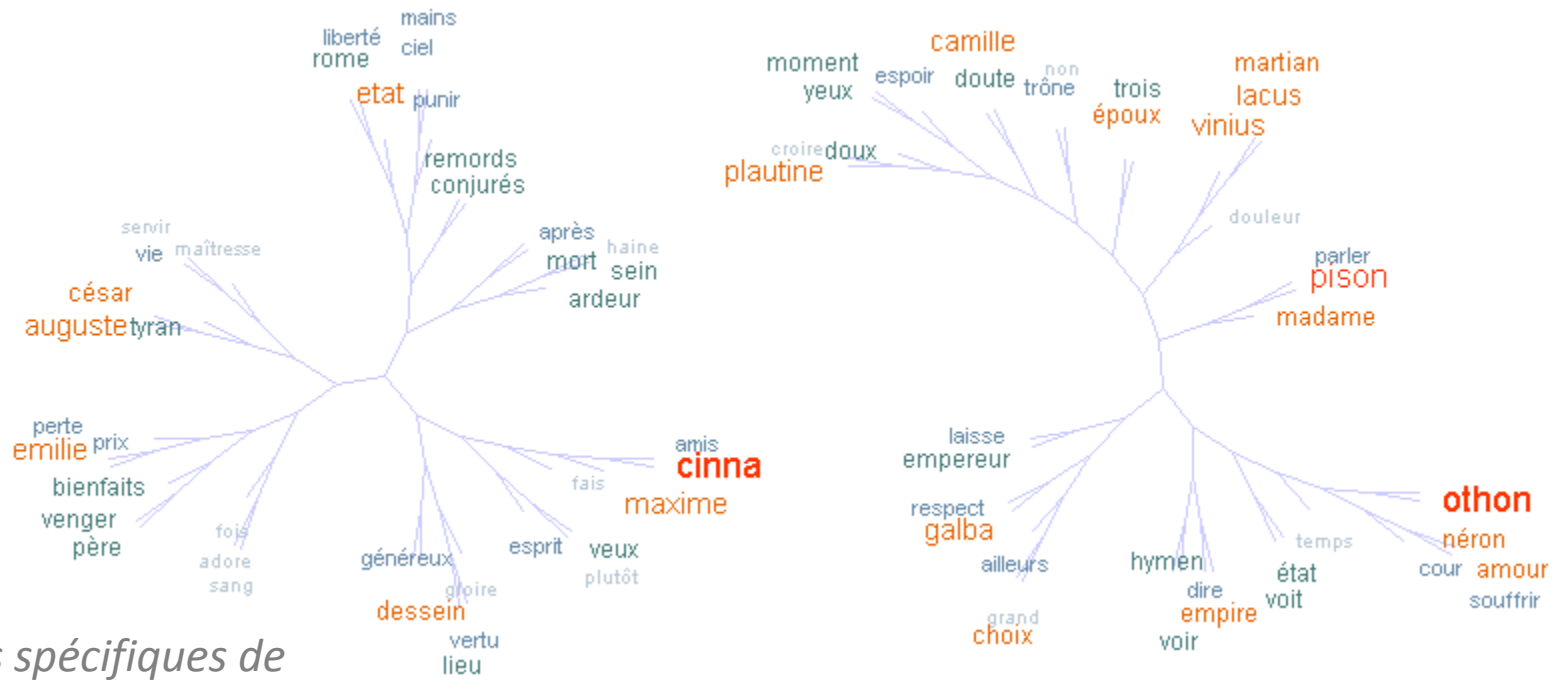


*Nuages arborés des **mots spécifiques** de Cinna et Othon, dimensionnés et colorés d'après leur spécificité calculée dans Lexico3.*



Quels moyens au service de la cause politique ?

Méthode : comparaison des spécifiques



mots spécifiques de Cinna et Othon d'après Lexico3

	<i>Cinna</i>	<i>Othon</i>
Lieu du pouvoir et objet de la confrontation entre les personnages	Rome (« liberté »)	Empire (« trône »)
Souverain en place	tyran	Empereur
Membres du corps politique	amis	maîtres / seigneurs
Moyens au service de la cause politique	gloire	amour matrimonial (« amour », « hymen », « choix »)
Caractérisation de la pièce	Pièce de FONDATION	Pièce de SUCCESSION DYNASTIQUE

Références

Disponibles sur TreeCloud.org :

Philippe Gambette, Jean Véronis (2009)

Visualising a Text with a Tree Cloud,

IFCS'09, Studies in Classification, Data Analysis, and Knowledge Organization 40, p. 561-570

<http://www.slideshare.net/PhilippeGambette/visualising-a-text-with-a-tree-cloud>

Delphine Amstutz & Philippe Gambette (2010)

Utilisation de la visualisation en nuage arboré pour l'analyse littéraire,

JADT'10 (Proceedings of the 10th International Conference on statistical analysis of textual data),

Statistical Analysis of Textual Data, p. 227-238

<http://www.slideshare.net/PhilippeGambette/utilisation-de-la-visualisation-en-nuage-arbor-pour-lanalyse-littraire>

Philippe Gambette, Nuria Gala & Alexis Nasr (2012)

Longueur de branches et arbres de mots,

Corpus 11:129-146

<http://www.slideshare.net/PhilippeGambette/longueur-de-branches-et-arbres-de-mots>

William Martinez & Philippe Gambette (2013)

L'affaire du Médiateur au prisme de la textométrie,

Texto! XVIII(4)

<http://www.revue-texto.net/index.php?id=3318>

Co-auteurs des travaux en cours :

- Edna Hernandez : méthodologie d'utilisation de TreeCloud pour les analyses exploratoires
- Claude Martineau : intégration de prétraitements Unitex dans TreeCloud
- Xavier Le Roux, Hilde Eggermont : analyse du corpus de projets sur la biodiversité