

Comportements et sentiments. De l'ambiguïté dans les émotions ?

**Pierre Molette
LERASS PsyCom
Toulouse mai 2014**

www.tropes.fr
www.owledge.org
www.lerass.com

Différentes approches pour l'analyse de textes

Plusieurs méthodes utilisant des logiciels :

- Post codification (« Stabilo »). On va annoter manuellement les mots des textes.
- Search « minimaliste ». On cherche une chaîne de caractères et on lit les textes (par exemple, Google).
- Lexicométrie. On fait des statistiques poussées, sans chercher à résoudre les polysémies.
- Sémantique, supervisée. Le logiciel propose une classification, qui sera rectifiée par l'utilisateur.
- Sémantique, purement automatique. Il faut « faire confiance » au logiciel.

Chacune de ces approches présente des avantages et des inconvénients.

Pour un panorama des méthodes d'analyse de textes, lire Marchand (1998).

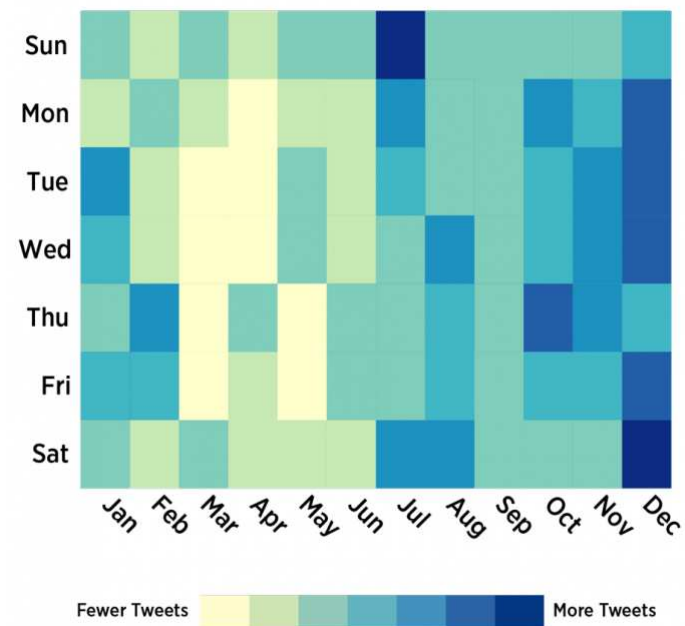
Plus la quantité de données est grande, plus il est difficile d'analyser les résultats.

Problème sémantiques sur les gros corpus

C'est pour cette raison que des analyses se limitent souvent à compter le nombre de réponses pour une simple expression. Par exemple, cette étude Twitter se limite à une recherche (« feeling sad ») :

When do you “feel sad”?

Tweets using the phrase “feeling sad” by day, 2013. Rate per million Tweets sent



Cela fonctionne sur un corpus énorme. Mais on ne peut pas considérer que c'est suffisant.

Les logiciels de classification sémantiques

Classer, c'est prendre le risque de se tromper. Mais c'est nécessaire pour comprendre.

A contrario, les moteurs de recherche internet affichent - sans les classer - d'interminables listes de résultats correspondant à ce que « tout le monde devrait chercher » (avec une forte orientation marketing). Cette approche est généralisée. Bien que scientifiquement peu satisfaisante quand on veut analyser objectivement.

Une classification sémantique automatique n'est réellement exploitable que si le taux d'erreur reste acceptable pour l'utilisateur. C'est possible sur de gros volumes de textes.

Comme dans tous les logiciels traitant le langage naturel, de nombreux termes polysémiques perturbent les résultats d'analyse. Comment améliorer les résultats ?

Partant du principe qu'il est préférable d'utiliser des termes spécifiques pour réduire le champ sémantique de certains mots, il convient de retirer (ou de spécialiser) certains concepts pour améliorer la qualité de l'analyse. En empruntant aux tests médicaux, disons qu'il faut étudier le rapport bénéfice/risque des mots ambigus pour isoler ceux qui ont trop d'effets indésirables (c'est-à-dire qui sont pratiquement toujours utilisés à contresens).

Les scénarios « Concepts » de Tropes ont été conçus pour réduire fortement les ambiguïtés. Ils s'appuient sur un réseau sémantique et une logique d'intelligence artificielle, qui évoluent depuis une vingtaine d'années.

Les logiciels Tropes et OwlEdge

Tropes (1994-2014) :

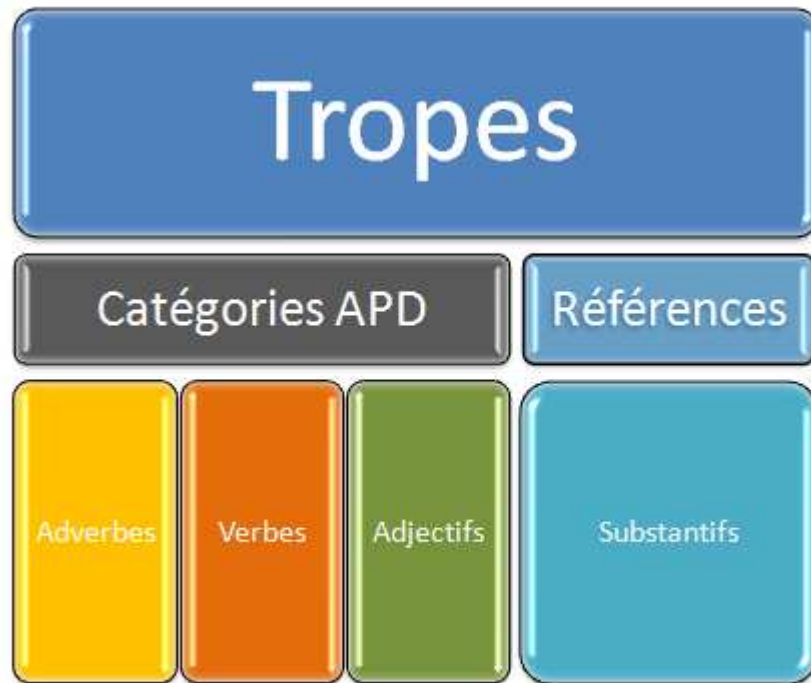
- Logiciel d'analyse sémantique de textes
- Développé par Pierre Molette et Agnès Landré, sur la base des travaux de Rodolphe Ghiglione
- Fondé sur l'Analyse propositionnelle du discours et l'Analyse cognitivo discursive (entre autres)
- www.tropes.fr

OwlEdge (2011-2014) :

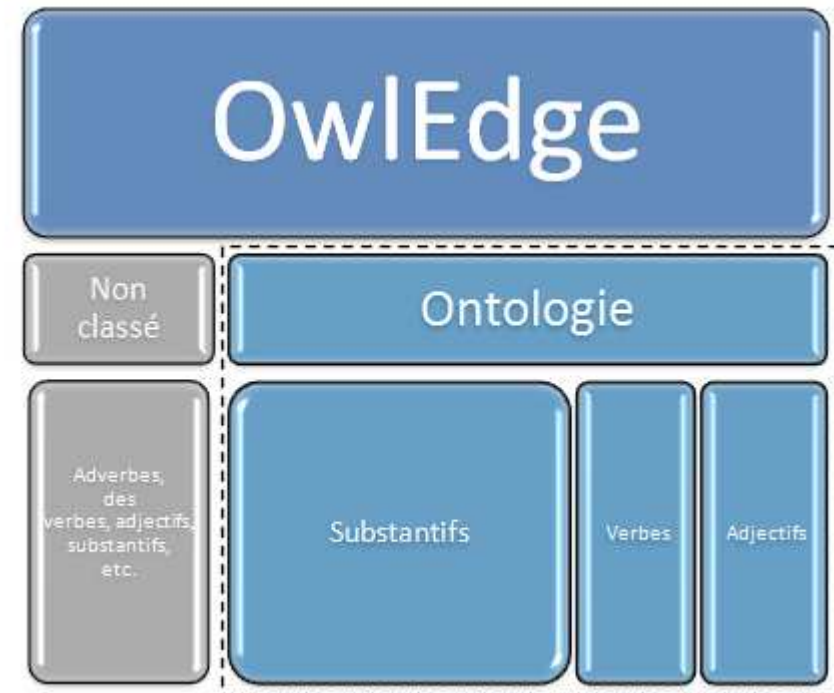
- Moteur de recherche et d'analyse par concepts
- Développé par Pierre Molette
- Exploite des ontologies sémantiques, dont certaines sont compatibles avec Tropes
- www.owlledge.org

Pour une explication du fonctionnement de Tropes, lire Ghiglione et al. (1998) ou Molette (2009).

TROPES VS OWLEDGE



Tropes, qui se fonde sur la théorie de l'APD, catégorise les mots à partir de la grammaire.



OwlEdge exploite une classification unifiée sans cloisonnement grammatical.

Ontologies et Scénarios

Les *scénarios* de Tropes sont des classifications sémantiques personnalisables (ontologies).

Il existe plusieurs scénarios de référence. Par exemple :

- Concepts Fr V8 (recouvrant 100.000 termes français) - Molette, Landré (1998-2014)
- Concepts US V8 & Natural Sciences US V8 (environ 100.000 termes anglais) - Molette et al. (200 ?)
- EMOTAIX - Lexique des émotions (4.600 termes) - Piolat, Bannour (2009)
- MeSH (Medical Subject Heading), thésaurus médical américain (NLM, MEDLINE), obsolète (2002).
- AGROVOC - Vocabulaire de l'agronomie multilingue (FAO, INRA, Cirad), 2006.
- etc.

Certains scénarios sont livrés avec le logiciel Tropes (*Concepts V8...*), d'autres sont disponibles sur demande auprès de leurs auteurs (par exemple EMOTAIX).

➔ Tous les utilisateurs du logiciel Tropes peuvent (et doivent) définir leurs propres classifications.

Comparaison de deux corpus contenant des sentiments

Etude de deux corpus, a priori peu comparables :

	Date de recueil	Nombre de documents	Format	Occurrences de mots	Taille totale du texte
Assemblée Nationale (Fr)	Janvier 2014	61 159	HTML	389 millions	5 Go
Classiques littéraires (ABU)	XXème siècle	100	Texte	5 millions	32 Mo

Le premier est une copie du site web de l'Assemblée nationale, collecté début janvier 2014. Ce corpus comprend des rapports, les comptes-rendus des débats et de multiples informations parlementaires. Ces textes ont été essentiellement rédigés durant les trente dernières années.

L'autre corpus contient cent classiques littéraires, qui ont été recueillis sur le site de l'ABU (Association de Bibliophiles Universels), <http://cedric.cnam.fr/ABU>, 1999

→ Ces deux corpus ont été analysés avec l'ontologie fournie par défaut avec Tropes (le scénario Concepts Fr V8), sans personnalisation des classifications. On y a toutefois ajouté une liste de 5000 noms propres. Le traitement est donc purement automatique.

→ Une personnalisation des classifications permettrait d'augmenter la précision des résultats.

Contenu de la branche « Comportement et sentiments »

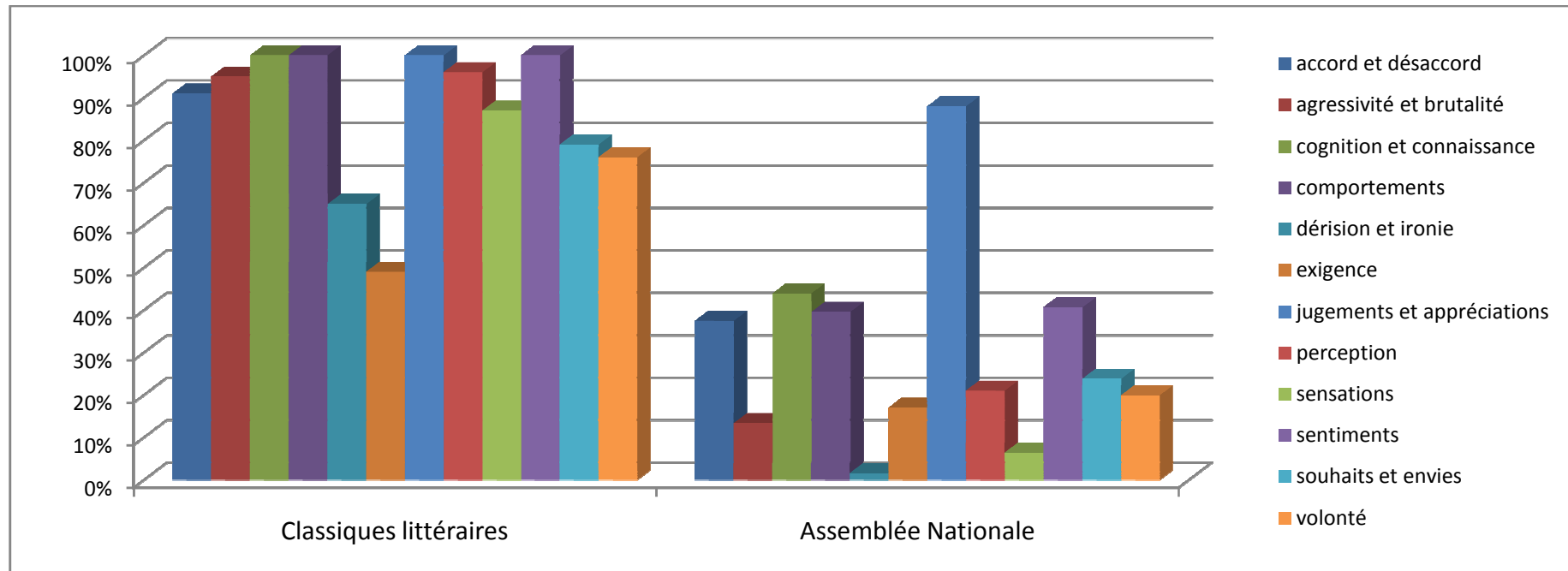
Treemap d'OwlEdge – Scénario Concepts Fr V8 de Tropes

comportements et sentiments		
comportements abstention abstinence acharnement ambition associabilité attitude	cognition et connaissance cognition compétence connaissance connexion esprit	sensations écoeurement faim fatigue et surmenage frisson
sentiments abomination amitié amour et aimer bonheur complaisance confiance et défiance	perception écoute goût odorat perception	souhaits et envies envie souhait vœux
jugements et appréciations estimation et diagnostic jugements de valeur jurons et insultes opinions et préjugés	agressivité et brutalité agressivité et violence coup inhumanité mauvais traitement	volonté intention préméditation volontaire
	accord et désaccord accords désaccords réclamation	dérision et ironie dérision ironie moquerie
		exigence

Nombre de documents. Données affichées pour les Classiques littéraires.

Comparaison de corpus asymétriques (suite)

Une approche simple consiste à comparer le nombre de documents classés, en pondérant les valeurs. Par exemple voici les résultats d'OwlEdge (sur « comportement et sentiments ») :

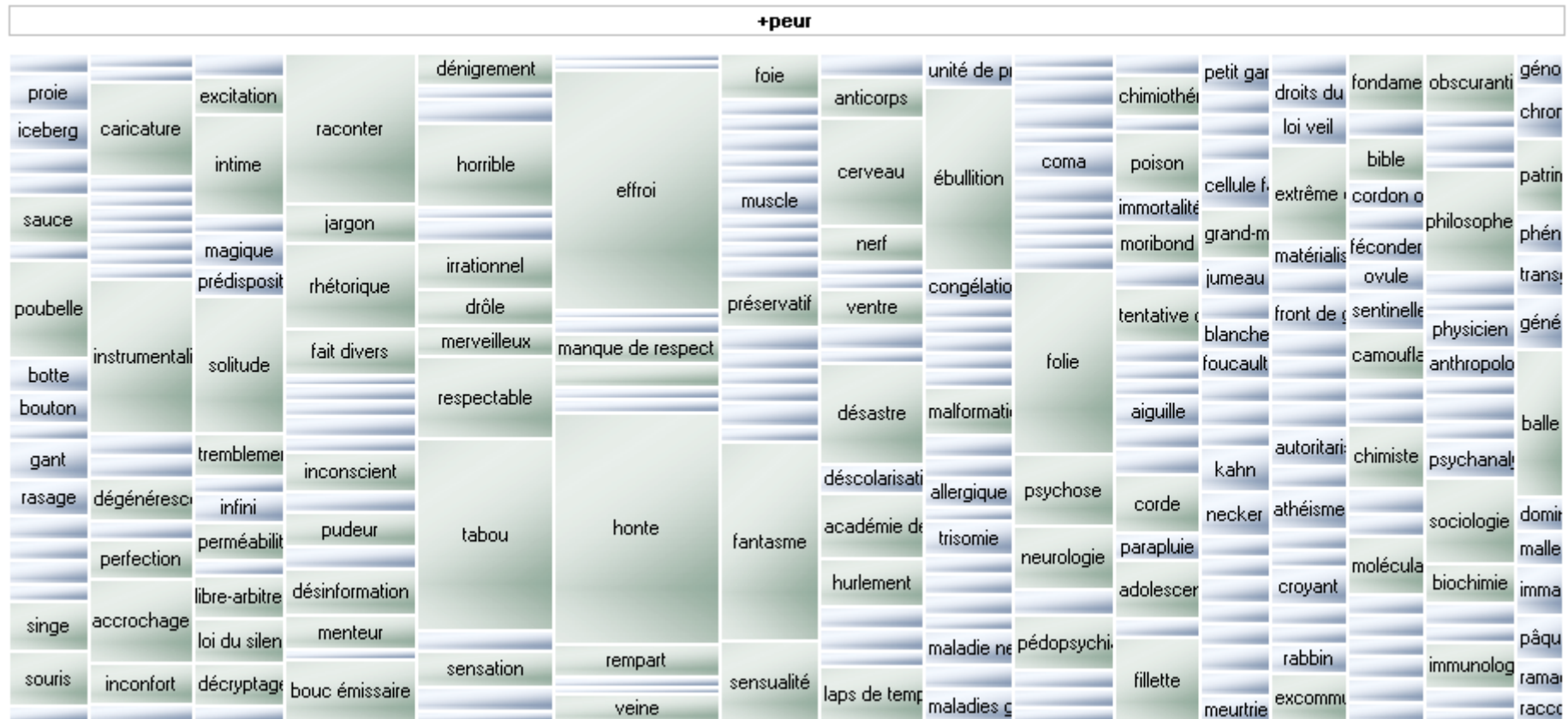


Pour les classiques littéraires, pratiquement toutes les branches de cette classification sont activées.

Ce n'est pas le cas pour l'Assemblée nationale, excepté pour les « jugements et appréciations », présents dans 88% des documents. Il ne faut pas en déduire qu'il y a peu d'émotions à l'Assemblée nationale...

Clustering d'OwlEdge sur la peur, documents

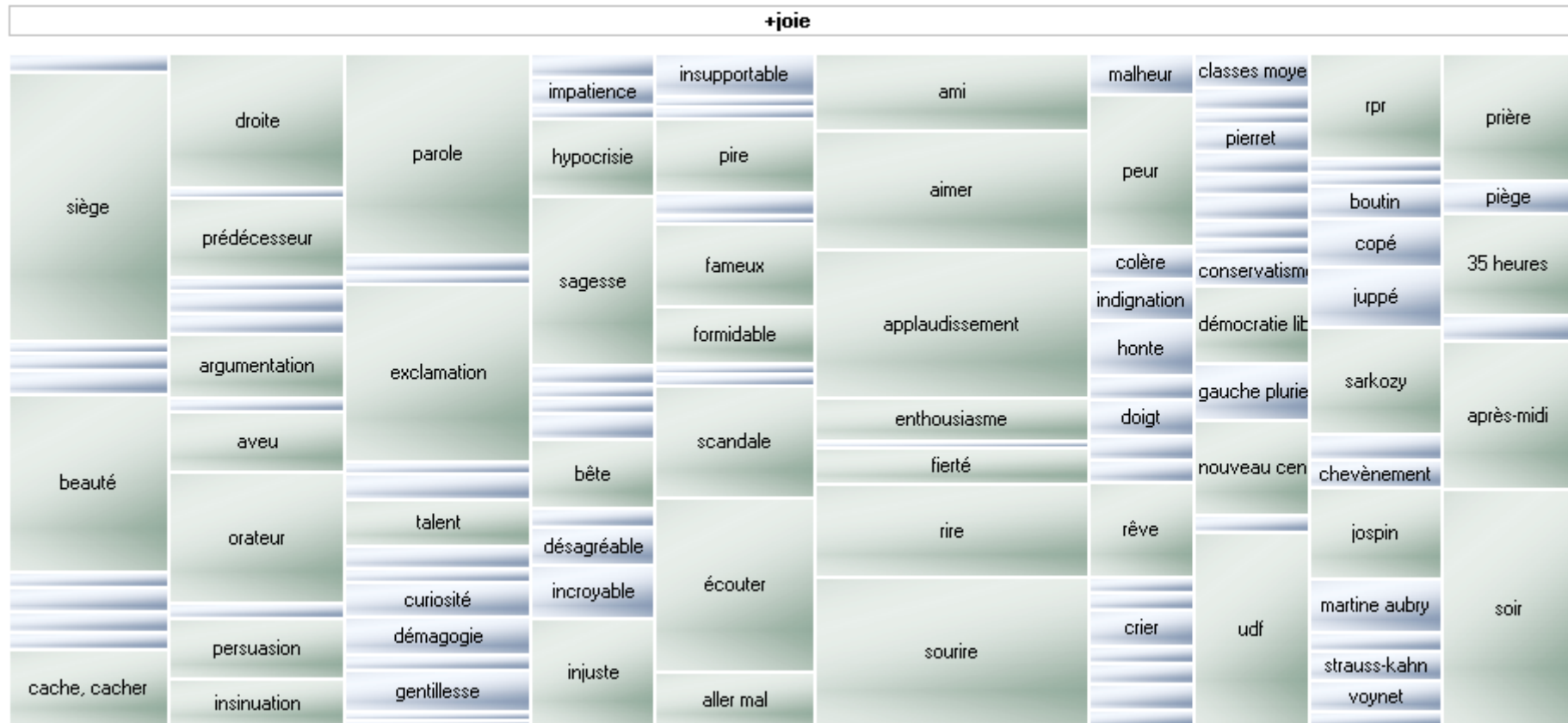
Résultat d'un clustering sur le concept de « peur », corpus Assemblée nationale :



La recherche intègre les sens associés à la peur (incluant l'angoisse et la crainte). On parle de troubles mentaux et de folie. Mais aussi d'autre chose... Est-ce que la science fait peur ?

Clustering d'OwIEdge sur la joie, documents

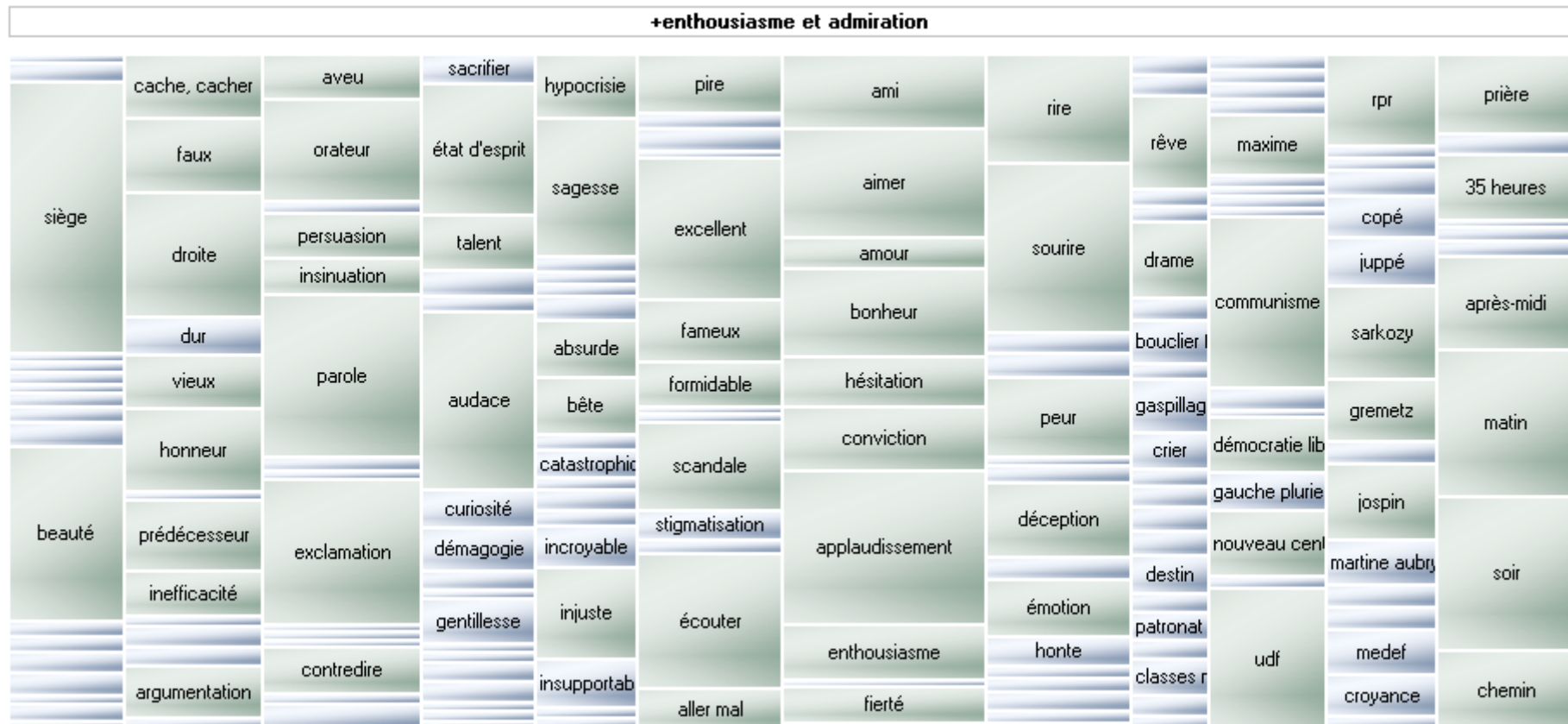
Résultat d'un clustering sur le concept de « joie », corpus Assemblée nationale :



La recherche intègre les sens associés à la joie (incluant rire et sourire). On voit apparaître des partis politiques et des noms propres, sur la partie droite du graphe.

Clustering d'OwlEdge sur l'enthousiasme, documents

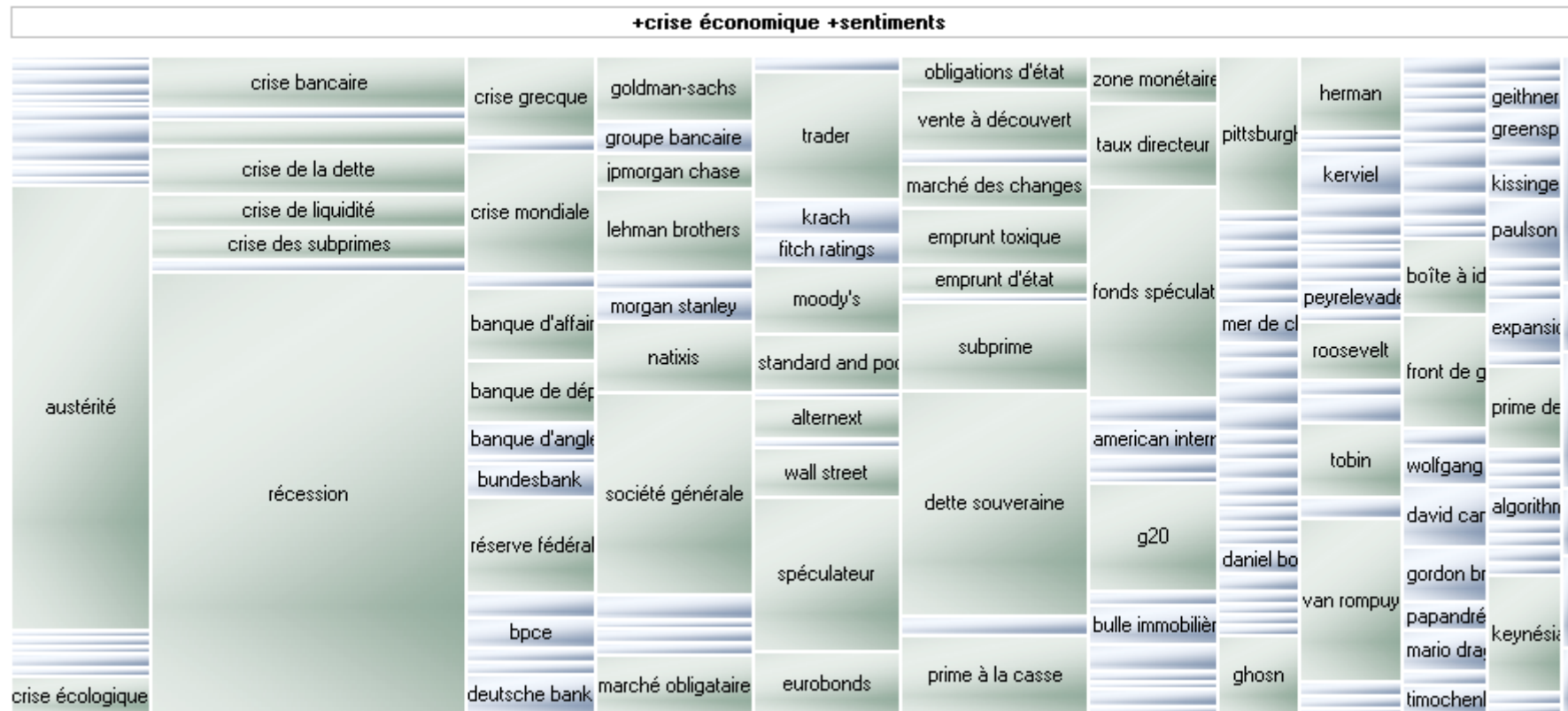
Résultat d'un clustering sur le concept « d'enthousiasme et admiration », corpus Assemblée nationale :



La recherche intègre les sens associés (incluant les applaudissements). Comme dans le cas précédent, on voit apparaître des partis politiques et des noms propres, sur la partie droite du graphe.

Clustering d'OwlEdge sur la crise économique, documents

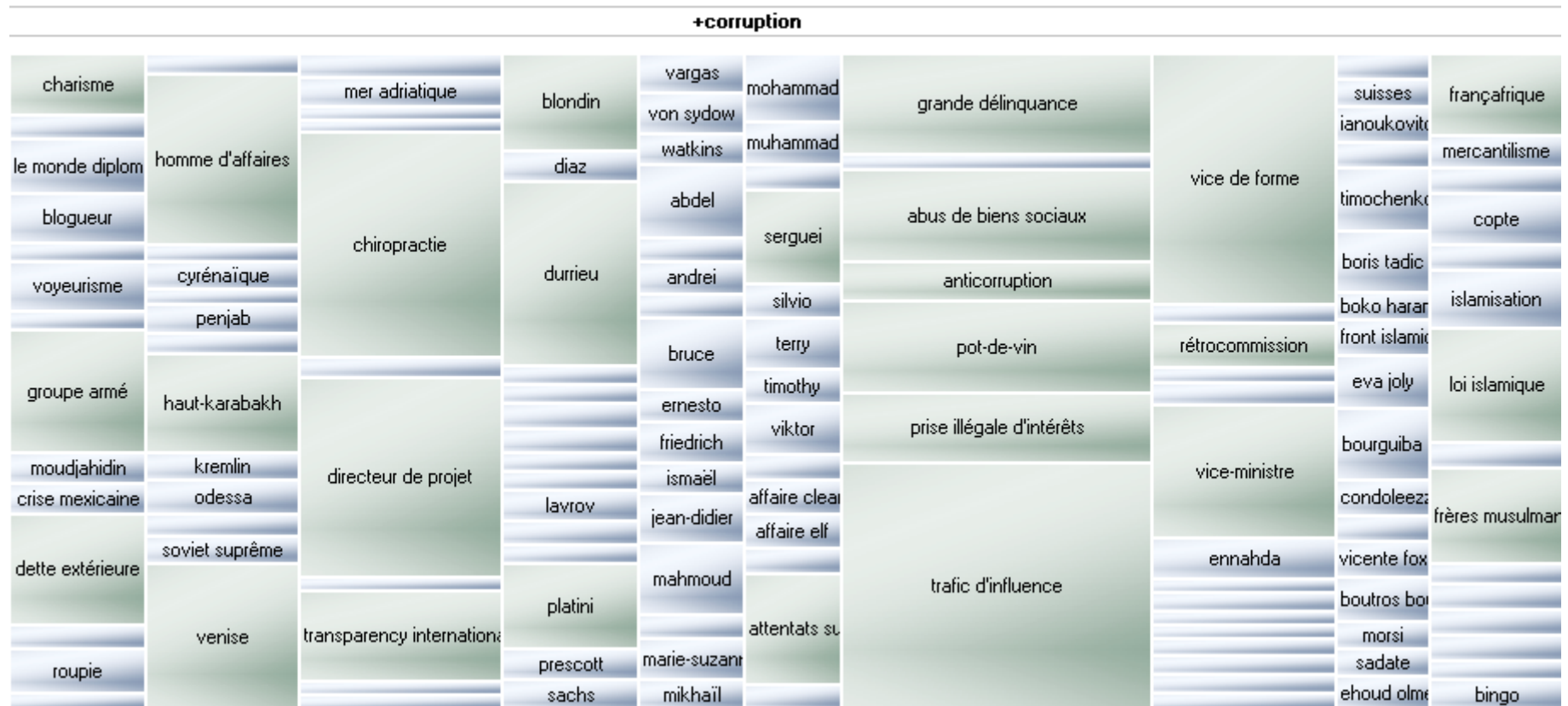
Résultat d'un clustering sur « crise économique » et « sentiments », corpus Assemblée nationale :



Les acteurs économiques et de nombreuses personnalités politiques étrangères sont présents sur ce graphe. Les comportements et sentiments de Tropes ne sont pas significatifs. Ils sont implicites ! Notez que les politiques français ne sont pas sur ce graphe, excepté le Front de gauche.

Clustering d'OwlEdge sur la corruption, documents

Résultat d'un clustering sur le concept de « corruption », corpus Assemblée nationale :



Certaines confessions religieuses, des corps de métiers et personnalités étrangères sont bien représentés sur ce graphe. Les personnalités françaises sont atypiques (ce ne sont pas celles qu'on trouve en analysant un corpus de presse avec la même requête). Que faut-il en déduire ?

BIBLIOGRAPHIE

Rodolphe GHIGLIONE et Al. *L'analyse automatique des contenus*. Paris, Dunod, 1998.

John LYONS. *Semantics*. Cambridge University Press. 1977.

Pascal MARCHAND. *L'Analyse du Discours Assistée par Ordinateur*. Paris, Armand Colin, 1998.

Pierre MOLETTE. *De l'APD à Tropes : comment un outil d'analyse de contenu peut évoluer en logiciel de classification sémantique généraliste*. Colloque Psychologie Sociale et Communication. Tarbes, 2009.

Pierre MOLETTE. *OwlEdge. Un moteur d'analyse de corpus par concepts*. Communication au séminaire du Labex SMS. Toulouse, mars 2014.

Annie PIOLAT, Rachid BANNOUR. *EMOTAIX : un scénario de Tropes pour l'identification automatisée du lexique émotionnel et affectif*. L'Année Psychologique, 109, 657-700.

TRANSPARENCY INTERNATIONAL. *Transparence de la vie publique et maintenant ?* Rapport 2013.

ANNEXES

Problèmes d'identification sémantique

L'Assemblée nationale regorge d'expressions difficiles à décoder, par exemple :

La Commission est saisie de l'amendement CL 16 du rapporteur.

M. le rapporteur. Je propose la création, s'agissant du blanchiment, de la corruption et du trafic d'influence, d'un statut de « repenti », à l'instar de ce que prévoit déjà le code pénal dans d'autres domaines. Une personne qui empêcherait la commission de l'infraction pourrait ainsi être exemptée de peine, et celle qui dénoncerait une infraction commise, bénéficier d'une réduction de peine.

De telles dispositions seraient de nature à donner à la justice un avantage important, notamment pour ce qui concerne les affaires les plus complexes.

M. Patrick Devedjian. Pourquoi ne pas également faire nôtre l'institution de la *bocca di leone*, qui encourageait les citoyens vénitiens à pratiquer la délation anonyme au nom de la protection de la république ?

M. le président Jean-Jacques Urvoas. Cela mérite un voyage d'étude ! (*Sourires.*)

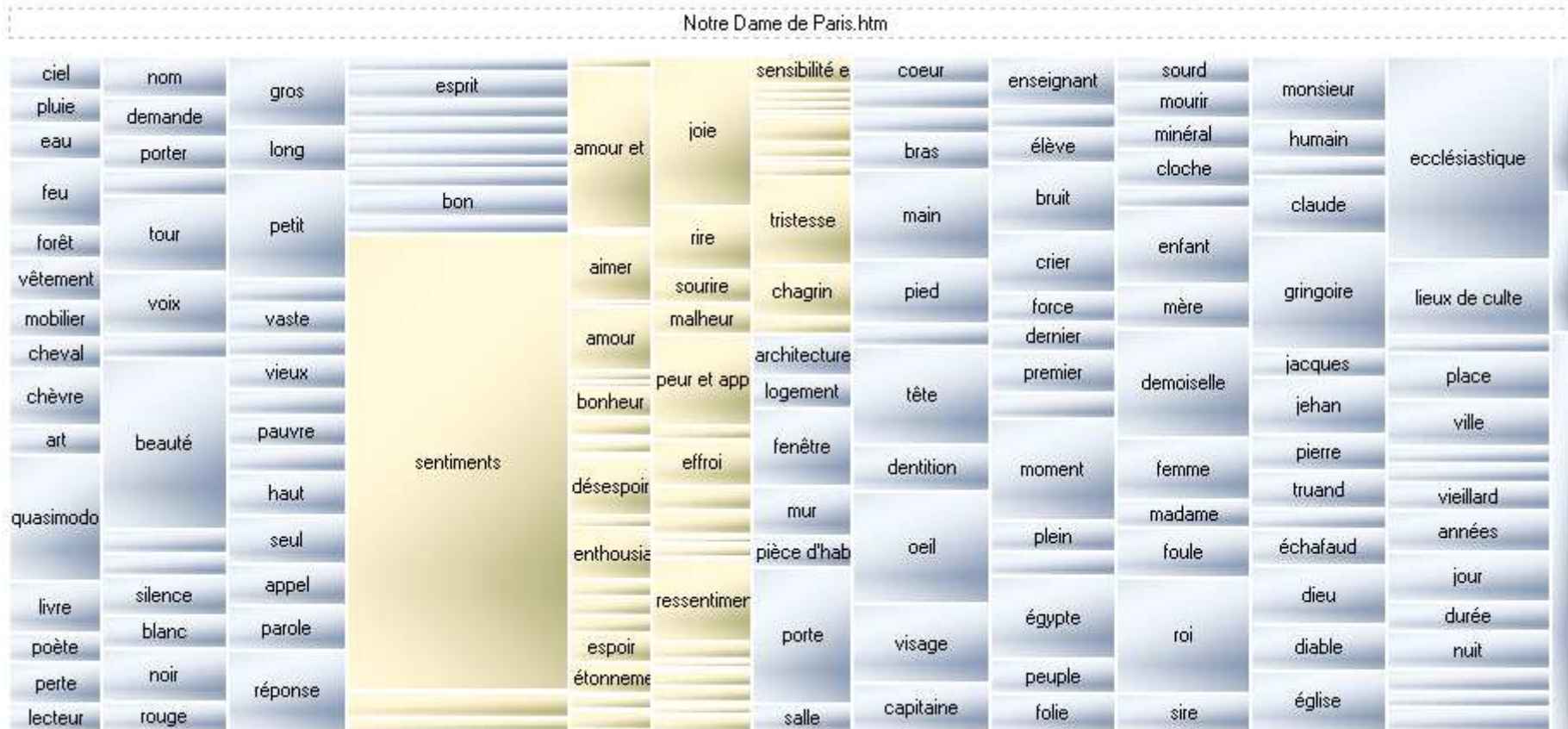
La Commission adopte l'amendement.

Commentaires sur le rapport de Yann GALUT. Nos 1130 et 1131. Enregistrement du 12 juin 2013.

Le traitement de ces informations nécessite le décodage d'implicites, complexes...

Les sentiments, dans un classique littéraire

Analyse d'un unique texte avec OwlEdge :



Nombre d'occurrences de mots. Notre Dame de Paris (Victor Hugo). Les sentiments sont fréquents.

Paires de graphes, Assemblée nationale

Suivent différents résultats de clustering des documents d'OwlEdge, sur des grappes de concepts antonymes, sur le corpus Assemblée nationale.

Tout ceci mériterait d'être complété par une étude plus longue.

