
La réduction de termes complexes dans les langues de spécialité

Yannis Haralambous*,** — Elisa Lavagnino*,***

* *Télécom Bretagne, Technopôle Brest Iroise, CS 83818, 29238 Brest Cedex 3, France
Yannis.Haralambous@telecom-bretagne.eu*

** *UMR CNRS 6285 Lab-STICC*

*** *Université de Gênes, Sezione di Francesistica, Piazza santa Sabina, 2, 16124 Genova, Italie. Elisa.Lavagnino@telecom-bretagne.eu*

RÉSUMÉ. Nous étudions, par des méthodes statistiques sur des corpus français et italiens, le phénomène de réduction des termes complexes dans les langues de spécialité. Il existe deux types de réductions : anaphorique et lexicale. Nous montrons que la réduction anaphorique dépend du type de discours (de vulgarisation, pédagogique, spécialisé) mais ne dépend ni du domaine, ni de la langue, alors que la réduction lexicale dépend du domaine et est plus fréquente dans les domaines techniques à évolution rapide. D'autre part, nous montrons que la réduction anaphorique a tendance à suivre la forme pleine du terme, nous définissons une notion d'arbre anaphorique de terme et nous étudions ses propriétés. Concernant la réduction lexicale, nous tentons de démontrer statistiquement qu'il existe une notion de cycle de vie de terme, où la forme pleine est progressivement remplacée par une réduction lexicale.

ABSTRACT. Our study applies statistical methods to French and Italian corpora to examine the phenomenon of multi-word term reduction in specialty languages. There are two kinds of reduction: anaphoric and lexical. We show that anaphoric reduction depends on the discourse type (vulgarization, pedagogical, specialized) but is independent of both domain and language; that lexical reduction depends on domain and is more frequent in technical, rapidly evolving domains; and that anaphoric reductions tend to follow full terms rather than precede them. We define the notion of the anaphoric tree of the term and study its properties. Concerning lexical reduction, we attempt to prove statistically that there is a notion of term lifecycle, where the full form is progressively replaced by a lexical reduction.

MOTS-CLÉS : Terminologie, langues de spécialité, extraction automatique de termes, termes complexes, réduction anaphorique, réduction lexicale.

KEYWORDS: Terminology, specialty languages, automatic term extraction, multi-word lexical units, anaphoric reduction, lexical reduction.

1. Introduction

En terminologie, nombreuses ont été les études portant sur les termes complexes (Collet, 2000 ; Jacques, 2003 ; Portelance, 1991), leurs variantes réduites (Freixa, 2002 ; Jacques, 2003 ; Møller, 1998), les textes et le discours spécialisés et le mécanisme de la réduction dans les langues de spécialité (Jacques, 1996 ; Lavagnino, 2011).

Pour justifier l'importance d'étudier le comportement des unités complexes dans le cadre des langues de spécialité, il suffit de dire qu'elles composent 80 % des éléments constituant les textes de spécialité (Jacques, 2003). Il est donc clair que leur rôle mérite d'être analysé et étudié dans tous ses aspects.

Notre contribution s'inscrit dans le cadre théorique de la terminologie textuelle et computationnelle (Bourigault et Slodzian, 1999 ; Conceição, 2005) appliquée aux discours spécialisés. Dans notre cas, le fait d'aborder les problèmes de la variation des termes complexes en terminologie, c'est prendre en compte au moins trois types de faits, à savoir :

- les caractéristiques de la structure interne d'un terme complexe (Portelance, 1996 ; Collet, 2000), son évolution (Møller, 1998) et les effets que cette évolution peut avoir sur la langue de spécialité (Jacques, 1996) ;

- les caractéristiques des environnements textuels étudiés (Desmet, 2006 ; Alexandru et Gaudin, 2005) l'usage et le comportement des termes dans le discours (Jacques, 2000) et les effets sur la communication, les usagers et la langue de spécialité (Lerat, 2009 ; de Vecchi et Estachy, 2008) ;

- le traitement automatique des terminologies (Daille *et al.*, 1998) et la terminologie computationnelle (Savary et Jacquemin, 2003).

La nature du terme et surtout son rôle dans le discours sont d'une importance fondamentale dans notre recherche. Ce n'est que récemment que l'intérêt pour la variation des termes a été l'objet de débats et d'analyses (Soglia, 2002). La fonction des termes n'est plus uniquement celle de dénommer une entité dans un système conceptuel figé ; l'évolution continue de la science et des techniques demande un renouvellement constant des terminologies et leur adaptation aux nécessités des usagers et des situations communicatives.

2. Autour de la définition du terme complexe

Le terme complexe est une unité monoréférentielle (critère de monoréférentialité) qui appartient à une langue de spécialité (critère terminologique), et qui représente une notion univoque (critère notionnel) (Collet, 2000). La structure d'un terme complexe est binaire par définition (critère de binarité : tête/expansion, déterminant/déterminé, hyper/hyponyme, *cf.* ci-dessous).

Il changement climatique et le **degrado ambientale** sont susceptibles de provoquer l'augmentation de la migration de l'Afrique subsaharienne, avec des effets potentiellement dévastateurs pour des centaines de millions de personnes, surtout des pauvres, qui y vivent. [...] Aujourd'hui, le **degrado** est un problème sérieux pour 32 pays d'Afrique, et plus de trois cents millions de personnes qui affrontent déjà la rareté de l'eau. [...] En outre, au Sénégal, les migrations internes et internationales ont été provoquées par les changements environnementaux qui ont réduit les opportunités de travail dans l'agriculture, qui est diminuée avec l'augmentation du **degrado ambientale**. [...] (<http://www.ecologiae.com/riscaldamento-globale-rifugiati-climatici-aumentare/14240/>)

Les terres vitées destinées à la production de vins à **Denominazione di Origine Controllata** (D.O.C.) doivent être inscrites – sur demande des producteurs intéressés par le biais du territoire compétent – en vue de la vérification par l'Inspection provinciale de l'Agriculture – en apposant l'Albo pubblico institué auprès de chaque Chambre de Commerce. [...] L'inscription à l'Albo sert à pouvoir effectuer la « demande de vins » – de la part des producteurs de terres déjà inscrites – destinées à la production de vin à **Denominazione di Origine Controllata**, au territoire compétent. Ceci afin de commercialiser le produit avec la respectueuse **denominazione**. [...] (http://www.cameradicommercio.ag.it/index.php?option=com_content&task=view&id=47&Itemid=110)

La **denominazione di origine controllata** « Soave » et « Soave classico » est réservée aux vins « Soave » (avec la spécification de la sous-zone Colli Scaligeri), « Soave » effervescent et « Soave classico », qui répondent aux conditions et aux exigences établies par le présent règlement de production. [...] Les conditions environnementales et de culture des vignes destinées à la production de vins à **denominazione controllata** « Soave » et « Soave » Classico doivent être celles traditionnelles de la zone et, de plus, aptes à conférer aux vins et au vin dérivé les caractéristiques spécifiques. [...] (<http://www.ilsoave.com/disciplinare.php>)

Figure 1. Exemples de termes complexes dans leur contexte

Parmi les critères structurels définissant les termes complexes, la binarité occupe sans doute une place fondamentale puisque tous les termes complexes ont une structure binaire¹, composée de deux éléments qui peuvent être simples ou composés (Collet, 2000).

Par exemple (cf. figure 1) :

– [*degrado*] [*ambientale*] : le terme complexe est composé de deux constituants simples. Ce facteur influence le mécanisme de la réduction qui, dans la plupart des termes à deux composants, ne s'avère qu'en donnant lieu à des variantes anaphoriques, comme *degrado*, ci-dessus ;

– [*denominazione*] [*di*] [*origine controllata*] : le terme complexe est composé de deux constituants eux-mêmes composés. Cette structure permet la réalisation de deux typologies de réduction (variantes : *denominazione* et *denominazione controllata*, ou *DOC*). Rappelons ici le débat concernant la relation entre le mécanisme de la siglaison et la réduction qui reste toujours ouvert. Dans notre étude, nous avons traité les deux mécanismes comme non équivalents à cause de la motivation qui les justifie : la

1. La binarité, en tant que caractéristique du terme complexe, nous permet de considérer sa structure divisée en deux parties permettant une meilleure identification des rôles des éléments : la tête, responsable de la caractérisation syntaxique du TC et les constituants qui en définissent les changements au niveau sémantiques. Pour l'instant, nous n'avons pas considéré les relations hiérarchiques entre les éléments, qui pourraient nous permettre d'envisager les rapports reliant les différentes composantes du terme complexe (Collet, 2000).

siglaison n'est pas un mécanisme spontané répondant à des exigences textuelles, mais plutôt un mécanisme émanant des experts d'un certain domaine (Abreu-Garcia, 1992).

La binarité peut être analysée selon les trois niveaux structuraux des termes complexes :

- une structure syntagmatique : tête et composants ;
- une structure sémantique : déterminant/déterminé ;
- une structure onomasiologique : hyperonyme/hyponyme.

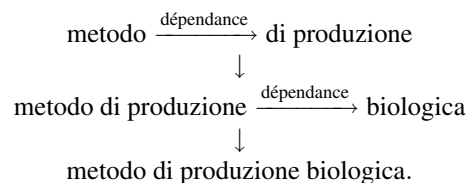
Au niveau syntagmatique, le terme complexe est une unité composée linéaire qui est formée de composants considérés soit comme des éléments forts (noms, verbes, adjectifs et adverbes) ou bien lexicaux, soit comme des éléments faibles (prépositions, articles et conjonctions) ou bien grammaticaux (Collet, 2000 ; Delmonte, 2004). Ces éléments sont en relation d'interdépendance grammaticale entre eux. La structure syntagmatique du terme complexe peut être schématiquement représentée comme suit (Collet, 2000) :

$$\text{terme complexe} = (\text{tête}) + (\text{composant}_1 + \text{composant}_2 + \text{composant}_3).$$

Lorsque nous nous situons au niveau sémantique, nous retrouvons la tête (déterminé) et les composants (déterminants). Les déterminants caractérisent au niveau sémantique le déterminé, qui donc change son référent notionnel.

Au niveau onomasiologique, les composants sont étudiés selon la typologie des relations qui s'instaurent entre eux, par exemple l'analyse des rapports d'hyperonymie/hyponymie.

Pour représenter graphiquement cette relation, prenons l'exemple *metodo di produzione biologica*. Nous constatons que la tête *metodo* instaure une relation d'hyperonymie/hyponymie avec le terme complexe en forme pleine ; d'ailleurs nous distinguons plusieurs degrés de dépendance :



La binarité d'un terme complexe est à la base de son instabilité et des changements que sa linéarité peut subir aux trois niveaux structuraux. Ces modifications sont déterminées par les déplacements des composants ou bien leur chute, comme dans le cas de la réduction.

Pour l'exemple *agricoltura biodinamica* nous avons :

- plan syntagmatique : [agricoltura](tête) + [biodinamica](composant) ;
- plan sémantique : [agricoltura](déterminé) + [biodinamica](déterminant) ;

– plan onomasiologique : [agricoltura](hyperonyme) + [biodinamica](hyponyme).

Si le composant *biodinamica* tombait, des changements seraient évidents à tous les niveaux. Il est donc fondamental de souligner que même si la réduction affecte tous les niveaux structuraux d'un terme complexe, les effets de ce mécanisme ne se reflètent pas sur la fonction dénomminative du terme complexe plein.

3. Autour de la définition de la réduction

Plusieurs études se sont intéressées à la question fondamentale suivante : « comment l'immersion dans un discours permet-elle d'omettre une partie *a priori* essentielle d'un terme complexe, c'est-à-dire formé par plusieurs mots, comme par exemple 'contrôleur de carrefour', 'équipement de terrain', 'effort à la commande', 'réseau routier national', etc. ? » (Jacques, 1996).

La réduction est un mécanisme discursif qui, à travers l'éllision d'au moins un constituant, transforme globalement un terme complexe en maintenant ses noyaux référentiel et notionnel (Collet, 1997).

Sur la base des recherches de Guilbert (Guilbert, 1975) et Portelance (Portelance, 1989), on peut affirmer que la réduction coïncide avec la suppression d'une information non différentielle, c'est-à-dire l'élément ou les éléments non fondamentaux pour la catégorisation du terme. Toutefois, cette acception limite ce mécanisme aux seuls phénomènes sociolinguistiques, sans tenir compte du fait que la réduction peut répondre également des exigences contextuelles reliées au contexte et au cotexte (Alexandru et Gaudin, 2005) ; cette acception diminue, en outre, la valeur de la réduction et exclut la possibilité que ce mécanisme puisse former des alternatives au syntagme plein. La réduction représente en revanche, comme nous le verrons, un élément important de cohésion textuelle.

De manière générale, ce phénomène a été traité soit comme un simple mécanisme d'anaphore qui permet par exemple de reprendre *le parc naturel* par *le/ce parc*, soit comme un processus de lexicalisation qui, par exemple, a transformé *voiture automobile* en *automobile* ou *téléphone portable* en *portable*. Schématiquement, l'effacement de l'expansion du terme complexe resterait, en tant qu'anaphore, étroitement dépendant du contexte et du cotexte, tandis que l'effacement de sa tête produirait une nouvelle unité lexicale. Le discours instaure ainsi un jeu entre les phénomènes de réduction participant à sa propre cohésion et ceux participant à la construction de nouvelles unités lexicales d'un domaine.

En outre, la variation par réduction satisfait le besoin d'économie de la langue (Tran *et al.*, 2008). En effet, la répétition de la forme pleine du syntagme devient lourde et non économique, alors que l'emploi d'une variante réduite conserve le noyau notionnel et référentiel, tout en évitant l'ambiguïté. Collet (2004) ajoute que la réduction, en ce qu'elle est précisément due à la réitération du terme complexe, d'une part, et en ce qu'elle constitue une forme différente de sa source, d'autre part, est à ranger

parmi les moyens de cohésion lexicale. La finalité de la réduction est donc double : sur le plan de l'encodage (Jacques, 2003), elle permet une économie d'énergie, sur le plan du décodage (Collet, 2000, p. 528–529), elle contribue à la cohésion textuelle en ce qu'elle constitue une forme de répétition d'un item lexical. Les caractéristiques de ces microstructures (ce que Collet appelle les « contextes réductionnels ») permettent la réduction et autorisent les locuteurs à se dispenser d'une partie *a priori* essentielle du terme. Ce mécanisme est donc en relation avec l'organisation plus globale du texte, témoignage de l'interaction entre le lexique et le discours.

3.1. La réduction à la base du polymorphisme

À travers le mécanisme de la réduction, une notion représentée par un terme complexe se trouve être dénommée par plusieurs unités qui partagent le même référent dans un domaine de spécialité. Dans les textes examinés, nous avons trouvé diverses formes du terme complexe comme dans l'exemple qui suit :

- terme plein : *mode de production biologique* ;
- variante 1 : *mode* ;
- variante 2 : *mode de production* ;
- variante 3 : *production biologique*.

Ainsi, certains termes complexes sont réalisés dans le même texte, sous trois formes différentes :

- une forme que nous appellerons « pleine », par exemple *mode de production biologique* ;
- une forme réduite à la tête du terme, *mode de production* ou *mode* ;
- pour certains, une forme réduite à l'expansion du terme, *e.g. production biologique*.

La polymorphie désigne la coexistence de plusieurs formes possibles d'expression pour certains termes complexes. S'il y a réduction, c'est parce qu'il y a adaptation du producteur du texte aux nécessités du discours. Certaines de ces réductions procurent un réel bénéfice sur le plan sémantique, en permettant de condenser certaines informations sur une seule occurrence. Par exemple le terme *produit issu de l'agriculture biologique* étant très long et peu économique, les usagers se sont tournés vers une forme abrégée qui puisse substituer le terme complexe dans le texte : *produit biologique*.

Telle qu'illustrée ci-dessus, la coexistence de diverses formes d'un terme complexe – une forme pleine et une ou plusieurs formes réduites – suscite deux questions essentielles. La première concerne la compréhension des facteurs de réduction et du mécanisme à l'œuvre dans l'effacement des constituants majeurs d'un terme complexe ; il s'agit de comprendre pourquoi et comment, en un point donné du discours, l'auteur choisit plutôt une forme qu'une autre, et comprendre ce qui, dans le discours, l'auto-

rise à user d'une forme dans laquelle toute l'information attachée au terme complexe n'est pas exprimée. Cela permet également d'expliquer ce qui rend possible et favorise l'effacement de constituants. La seconde question concerne les conséquences sémantiques de l'effacement de constituants. L'enjeu est de dégager les facteurs propres au discours qui permettent de mieux comprendre la polymorphie des termes complexes.

3.2. La valeur des variantes réduites

Le terme complexe et les variantes correspondantes donnent vie à des couples réductionnels. Les variantes peuvent avoir différents degrés d'autonomie référentielle, selon leur charge sémantique par rapport au terme complexe en forme pleine. La charge sémantique d'une composante représente sa valeur informative à l'intérieur du terme : plus une composante est fondamentale pour la transmission du sens du terme complexe, moins elle subira l'effet de la réduction. La valeur informative d'une composante ne constitue pas le seul facteur qui influence sa chute potentielle. L'usage et son figement à l'intérieur d'une terminologie représentent d'autres facteurs qui peuvent être cause de réduction. Comme nous l'avons déjà anticipé, même les exigences textuelles peuvent être un facteur causant la réduction.

Selon Guilbert (1975) et Portelance (1989), nous observons que la réduction est un mécanisme qui supprime l'information non différentielle, c'est-à-dire l'élément ou les éléments qui ne caractérisent plus le terme complexe.

Le mécanisme de la réduction peut également être justifié par des raisons pragmatiques : le facteur le plus incisif dans le cas de la réduction anaphorique est l'organisation du discours ; contrairement à cela, la réduction lexicale est influencée par des procédures lexicales qui donnent lieu à des unités terminologiques indépendantes.

Nous pouvons déjà en déduire qu'il existe des conditions internes et externes aux termes complexes, et que celles-ci peuvent influencer leur évolution. Les conditions internes sont plutôt reliées aux caractéristiques du terme complexe et à sa structure, tandis que les conditions externes subissent l'influence de l'environnement textuel, des caractéristiques de la langue de spécialité considérée et du degré de figement des termes complexes.

Dans la section suivante, nous introduisons les deux types de réductions analysés dans notre étude.

3.3. Les réductions anaphorique et lexicale

En général, nous pouvons identifier différents types de processus de réduction qui peuvent être classés sur un continuum. Celui-ci est organisé selon les relations qui s'instaurent parmi les termes complexes en forme pleine et les variantes réduites, les composantes qui chutent et la valeur de la variante qui se forme. De toute façon, dans toutes les études (Collet, 2000 ; Jacques, 2003 ; Adelstein, 2002 ; Cardero García,

[...] Au niveau européen, les règlements (CEE) n° 2092/91 du 24 juin 1991 et (CE) n° 1804/99 du 19 juillet 1999 définissent les règles du mode de production biologique et du contrôle des produits, le premier pour les végétaux et le second pour les produits animaux. Les organismes génétiquement modifiés (OGM) et produits dérivés sont exclus des **modes de production biologique**.

Fin 2004, les surfaces en **production biologique** (déjà certifiées ou en conversion) s'élèvent à 535 000 hectares, soit 1,9 % de la surface agricole utilisée (SAU) française. [...] Ce **mode de production** touche environ 11 000 exploitations, ce qui correspond à la moitié des 25 000 exploitations et 1 million d'hectares visé par le plan pluriannuel de développement de l'agriculture biologique, présenté en décembre 1997 par le ministre en charge de l'Agriculture. [...] L'implantation de ce mode d'exploitation est faible dans les zones de grandes cultures intensives du Bassin parisien. (<http://www.stats.environnement.developpement-durable.gouv.fr/donnees-essentielles/activites-humaines/agriculture-et-environnement/1-agriculture-biologique.html>)

Le vin bio c'est le sang de la terre, différent d'un terroir à l'autre, sain, authentique, sans artifices. C'est un nectar qui célèbre la vie, c'est l'invité incontournable des dîners festifs entre amis. Pour obtenir du vin, le viticulteur qu'il soit en mode conventionnel ou en mode biologique doit cultiver sa vigne avant de vinifier son raisin. Nous allons vous exposer la spécificité du travail du viticulteur en mode viticulture biologique. [...] (<http://www.terroirselect.info/territoires/Provence/cultiver-vigne-mode-bio.html>)

[...] Il **metodo di produzione biologico** è disciplinato a livello comunitario dai reg. CE 2092/91 (normativa base) e 1804/99 (disposizioni per le produzioni animali). L'Italia è il primo paese in Europa per numero di aziende che applicano il metodo di produzione biologico. [...] (<http://www.coldiretti.it/anagribios/agricoltura.htm>)

[...] L'agricoltura biologica è un **metodo di produzione** definito dal punto di vista legislativo a livello comunitario con un regolamento, il Regolamento CEE 2092/91, e a livello nazionale con il D.M. 220/95. [...] (http://www.aiab.it/index.php?option=com_content&view=article&id=112&Itemid=136)

[...] Parla di dati incoraggianti anche il sottosegretario alle Politiche agricole alimentari e forestali, Stefano Boco, che « dimostrano una significativa attenzione degli agricoltori verso il **metodo biologico** ». [...] (http://www.aiol.it/last_node/articolo?page=863)

Figure 2. *Exemples de réductions dans leur contexte*

2003) de ce mécanisme, ce sont les réductions anaphorique et lexicale qui retiennent l'attention de la plupart des auteurs.

La réduction anaphorique est un processus discursif et textuel, tandis que la réduction lexicale est générée par des conditions internes au syntagme plein (caractéristiques morphosyntaxiques, notionnelles, statut terminologique des constituants) ou par des conditions externes (niveau de spécialité du texte, typologie textuelle). Cette distinction est globalement acceptée, mais certains spécialistes utilisent des dénominations différentes pour les mêmes notions. Par exemple Alber-DeWolf (1984) parle d'« ellipses contextuelle et néonymique », Kocourek (1991) d'« ellipses contextuelle et lexicale », Jacques (1996) définit les variantes comme « reprise anaphorique » et « terme réduit sans antécédents textuels ». Dans cet article (ainsi que dans (Lavagnino, 2011)) nous avons choisi les termes de « réductions lexicale et anaphorique », qui nous semblent les plus transparents.

Au niveau syntaxique phrastique et interphrastique, les variantes ont la même valeur. Nous détectons une différence au niveau interne de la structure du terme complexe : une variante lexicale peut s'avérer sous la forme d'un changement non linéaire de la structure de terme complexe (par exemple : *mode de production biologique* et *mode biologique*), tandis que la réduction anaphorique entraîne un changement qui est toujours linéaire et qui détermine la chute des composants qui suivent la tête du terme complexe (par exemple *degrado ambientale* et *degrado*).

Afin de différencier les deux formes de réduction, Jacques (2000) affirme que :

- la réduction lexicale n'est pas liée au milieu contextuel immédiat, qu'elle est durable, qu'elle crée des variantes susceptibles de devenir des membres permanents de la terminologie du domaine, et qu'au niveau onomasiologique le syntagme plein est conservé même hors contexte ;

- la réduction anaphorique a une valeur contextuelle, qu'elle se déroule au noyau même du terme complexe, qu'elle a une valeur cohésive parmi les phrases d'un même texte et qu'au niveau onomasiologique, le syntagme plein est conservé uniquement en contexte.

Pour mieux expliquer notre propos, voici quelques exemples (*cf.* figure 2) :

	Exemple italien	Exemple français
Terme	<i>metodo di produzione biologica</i>	<i>mode de production biologique</i>
Syntagme nominal plein	[metodo] subs. + [di] prép. + [produzione] subs. + [biologica] adj.	[mode] subs. + [de] prép. + [production] subs. + [biologique] adj.
Tête	metodo	mode
Constituants	di produzione biologica	de production biologique
Réduction lexicale	metodo biologico	mode biologique
Réduction anaphorique	metodo di produzione metodo	mode de production mode

Dans l'exemple italien, la tête du terme complexe reste liée à ses constituants dans le cas de la réduction lexicale ; en revanche, pour la réduction anaphorique la variante a seulement une valeur de cohésion anaphorique, alors que la tête détient une charge sémantique supérieure.

À travers ce type d'analyse, on définit un critère fonctionnel : la variante lexicale réduite sert également comme connecteur textuel, mais elle ne se limite pas à cette fonction.

Du point de vue onomasiologique, la différence entre les deux formes de réduction réside dans le sens transmis aux variantes : dans le cas de la réduction lexicale, la variante n'est pas un hyperonyme pur du terme complexe, mais elle en conserve la valeur onomasiologique ; au contraire, la variante anaphorique représente l'hyperonyme du terme complexe, donc sur le plan onomasiologique elle a moins d'intention et surtout est étroitement liée au contexte. Exemple :

Terme	<i>logo comunitario di controllo CEE</i>
Syntagme nominal plein	[logo] subj. + [comunitario] adj. + [di] prép. + [controllo] subj. + [CEE] acronyme/subst.
Tête	logo
Constituants	comunitario di controllo CEE
Réduction lexicale	logo di controllo CEE, ou logo CEE (relation non hyperonymique)
Réduction anaphorique	logo (hyperonyme)

Sur le plan formel, on peut dire que la réduction anaphorique permet seulement l'élision des constituants – en revanche, la réduction lexicale permet la suppression de la tête, des constituants, des éléments forts et faibles. Exemple :

Terme	<i>agricoltura biologica</i>
Syntagme nominal plein	[agricoltura] subst. + [biologica] adj.
Tête	agricoltura
Constituants	biologica
Réduction lexicale	il biologico
Réduction anaphorique	agricoltura

On a donc enregistré une élision du constituant, dans le premier cas, et de la tête, dans le deuxième exemple. La variante lexicale, qui détient une charge sémantique supérieure par rapport à la variante anaphorique, peut substituer le terme complexe dans les textes, avec un risque très réduit d'ambiguïté

Sur le plan sémantique, il y a également des différences. La variante lexicale est caractérisée par une cohésion intérieure plus forte par rapport au syntagme plein.

Par exemple, le syntagme *audit environnemental* se transforme en *audit* : il y a donc suppression de la cohésion interne du terme complexe. Mais, en cas de réduction lexicale, le comportement est à l'opposé : à cause de l'élimination des constituants, la variante a une cohésion interne réduite, voire parfois annulée.

Ainsi, par exemple, la variante lexicale *metodo biologico* du terme complexe *metodo di produzione biologico* a une cohésion interne plus forte, le signifié dans la forme réduite étant distribué parmi les constituants de la variante, qui sont donc plus cohésifs et ont plus de valeur. Le signifié de la forme pleine est distribué dans tous les éléments du syntagme.

En général, on peut affirmer qu'entre les variantes anaphoriques et les termes complexes, il s'instaure des relations de type hyperonymique, et entre les variantes lexicales et les termes complexes des relations synonymiques, bien que ces variantes ne soient pas de véritables synonymes.

3.4. Réductions anaphorique et lexicale dans le discours

Selon le domaine, le mécanisme de réduction peut se manifester de différentes ma-

nières. La réduction anaphorique ne subit pas l'influence du domaine, vu qu'il s'agit d'un phénomène intratextuel lié aux caractéristiques du texte. En revanche, la réduction lexicale est influencée par le domaine. Plus un domaine est technique et plus il sera caractérisé par une évolution rapide de la terminologie qui peut produire (Dury et Drouin, 2010) :

- 1) l'effacement de certains éléments n'ayant plus une fonction différentielle dans un terme complexe ;
- 2) l'ajout d'éléments aux termes déjà existants ;
- 3) la disparition des termes.

Dans le premier cas, nous retrouvons le mécanisme de la réduction. En revanche, l'expansion des termes complexes peut, elle aussi, et dans un deuxième temps, déterminer la réduction, si la linéarité du terme complexe change pour des raisons d'économicité. Cette dernière affirmation est reliée au concept de série syntagmatique : pour chaque terme de la série qui enregistre un ajout d'information à travers l'expansion de sa structure, nous pouvons constater l'effacement d'un composant devenant implicite et d'apport non différentiel.

Ce mécanisme de réduction assez fréquent s'appuie sur le principe de redondance et s'apparente à l'apocope en langue parlée (Portelance, 1991).

En plus, un domaine à haut degré de réduction est normalement caractérisé par des lexiques composés surtout par des termes complexes, et donc plus sensibles à subir des mécanismes comme la réduction qui sont enchaînés dans des séries. Dans ses travaux, Portelance (Portelance, 1989 ; Portelance, 1991 ; Portelance, 1996) décrit cette tendance des terminologies dans les domaines techniques, en justifiant la chute des éléments qui ont cessé d'être différentiels pour éviter leur redondance. Ce type de chute est à la base de la réduction lexicale dans le domaine technique. La réduction anaphorique dans les domaines techniques ne permet normalement pas l'évolution d'une langue de spécialité puisque les variantes anaphoriques ne sont pas indépendantes.

4. Autour de la définition des corpus

4.1. *Les langues de spécialité et les discours de spécialité*

Si la langue de spécialité est perçue comme une variété de la langue générale (Prandi, 2006), la variation des langues de spécialité est soumise aux mêmes critères que la langue générale (variations diachroniques, diatopiques et diastratiques). Nous pouvons retrouver d'autres facteurs de variation qui influencent également la communication spécialisée, comme la situation de communication, les intentions et les buts de la communication. Ceux-ci conditionnent les ressources syntaxiques, morphologiques et textuelles utilisées dans les textes spécialisés (types de phrases, formes verbales, articulateurs du discours, etc.) (Desmet, 2001).

Quand nous nous référons aux langues de spécialité, il est important de ne pas confondre le concept de discours avec celui de texte ou énoncé. Le discours se définit comme le produit des multiples pratiques discursives à l'œuvre dans la vie sociale (Desmet, 2001). Il ne peut pas être dissocié du contexte socioculturel auquel il s'insère. L'énoncé représente sa manifestation ponctuelle, un objet concret et observable. Le texte enfin se réfère au modèle abstrait à partir duquel s'organisent les énoncés.

Les textes peuvent être classés selon des critères permettant d'identifier les niveaux suivants :

- 1) niveau fonctionnel (fonctions textuelles) ;
- 2) niveau situationnel (contexte social des activités communicatives, lieu, temps, nombre, rôle et relations entre les locuteurs) ;
- 3) niveau du contenu sémantique (sujet d'un texte, différentes perspectives et développement thématique) ;
- 4) niveau formel et grammatical (formes linguistiques et non linguistiques, aspects grammaticaux, ressources syntaxiques et lexicales) (Desmet, 2001).

Les critères que nous venons de citer se distinguent en :

- critères internes, qui concernent directement la structure du texte ;
- critères externes, qui considèrent le contexte où se situe le texte.

Le contexte est caractérisé par la situation communicative, l'émetteur, le destinataire et les typologies textuelles.

Un autre concept important à définir est celui de domaine, seule façon de délimiter, de dénommer une structure cognitive, conceptuelle (de Bessé, 2000).

Notre projet se fonde sur deux approches différentes des corpus : d'une part, l'analyse linguistique du phénomène de la réduction s'appuie sur une étude de corpus qui ont été constitués *ad hoc*, donc selon une approche *corpus-based*, d'autre part, la validation expérimentale est fondée sur des corpus ayant été constitués en vue, non pas de l'analyse de la réduction, mais de la validation des conclusions tirées de l'étude linguistique, donc selon une approche *corpus-driven* (Condamines, 2005).

En parlant de corpus, nous désignons l'aspect normatif de la langue, notamment sa structure et son code. Le corpus regroupe un ensemble de textes ayant une visée commune. Les catégories des textes contenus dans un corpus peuvent être différentes, mais normalement elles partagent un objectif commun. Les critères que nous considérons comme les plus pertinents pour classer les textes dans un corpus sont : le sujet, la perspective d'énonciation, le niveau de spécialisation, les sources, la typologie textuelle, la langue.

4.2. Les typologies textuelles

Le mécanisme de réduction anaphorique dépend des caractéristiques du texte dans lequel le terme est inséré. En général, nous pouvons relier ce mécanisme au classement des typologies textuelles de Sabatini (1999), (Lavagnino, 2011). Plus un texte est contraignant et plus évidente sera la tendance de la terminologie à rester figée, comme nous le montrerons en section 5.

Cette classification est inspirée de celle de Desmet (2006) et vise à créer des typologies textuelles homogènes qui puissent être utilisées dans les mêmes situations communicatives. Celles-ci, étant ciblées en langue de spécialité, imposent des contraintes qui peuvent concerner le niveau situationnel influencé par les activités communicatives, lieu, temps, nombre, rôle des locuteurs et relations entre eux.

Dans notre cas, les textes ont été classés selon des contraintes prenant en compte le niveau situationnel de la communication mais aussi le contexte linguistique et le microcontexte (densité de termes, structure du texte).

Dans notre étude, nous avons traité les domaines indiqués dans le tableau 1. :

Domaine	Réduction lexicale	Réduction anaphorique
Espaces naturels	oui	oui
Médecine vétérinaire	non	oui
Cancer	non	oui
Emballages	oui	non
Philosophie	oui	non

Tableau 1. Les domaines traités et les types de réductions étudiés

Pour chaque domaine, nous avons créé des corpus de référence, à partir du contenu textuel de revues spécialisées du domaine.

En ce qui concerne la réduction anaphorique, nous avons fondé notre étude contrastive sur des textes tirés du Web et subdivisés dans des catégories textuelles qui se fondent sur la classification suivante :

- catégorie 1 : discours de vulgarisation et de semi-vulgarisation scientifique ou technique. Par exemple, articles tirés des journaux généraux, brochures, sites Web non spécialisés ;

- catégorie 2 : discours scientifique ou technique à des fins pédagogiques. Par exemples textes universitaires, textes destinés aux experts des domaines pour la veille technologique et scientifique ;

- catégorie 3 : discours scientifique ou technique spécialisé et/ou officiel, discours législatif. Par exemple, lois qui règlent les domaines, articles scientifiques.

4.3. *Statistique descriptive du comportement des termes complexes*

Plusieurs auteurs se sont intéressés à l'extraction de termes complexes dans les langues de spécialité (Daille, 1996 ; Frantzi *et al.*, 1998 ; Jacquemin, 2001 ; SanJuan et Ibekwe-SanJuan, 2002 ; Ngonga Ngomo, 2008). La réduction, en tant que cas particulier de la variation, a également été étudiée (Jacquemin, 1999 ; Daille, 2003 ; Nenadić *et al.*, 2004).

En marge de ces importants travaux, souvent liés à des développements d'outils d'extraction, il y a eu quelques études statistiques du comportement des termes, portant surtout sur la performance des outils (Daille *et al.*, 1998 ; Paziienza *et al.*, 2005).

Dans cet article nous nous intéressons spécifiquement au phénomène de réduction des termes complexes, ainsi, nous visons à contribuer, par une meilleure compréhension du comportement des termes complexes, à l'optimisation des systèmes d'extraction de termes, de l'analyse sémantique et de la traduction automatique.

Notamment, nous avons distingué :

- 1) une phase concernant l'extraction automatique des termes complexes des textes de spécialité ;
- 2) une phase concernant la validation expérimentale des hypothèses qui sont décrites dans le chapitre précédent.

Les deux approches utilisent les mêmes protocoles qui visent des objectifs communs.

Pour ce qui concerne l'extraction automatique des termes, nous nous sommes servis du logiciel *Acabit*. Cette validation empirique ne pouvait s'appuyer que sur des corpus créés *ad hoc*, au risque de montrer les hypothèses de départ sans mettre en évidence d'éventuelles contre-hypothèses. Pour éviter cela, et afin de vérifier ultérieurement les axiomes de départ, nous avons également décidé d'élaborer d'autres corpus qui se réfèrent à d'autres langues de spécialité.

Dans le cas de la réduction anaphorique, nous avons constitué des corpus d'apprentissage d'où nous avons extrait les termes à travers *Acabit*. Après avoir obtenu les listes de termes complexes associés à leurs variantes éventuelles, nous avons lancé des requêtes sur Internet pour retrouver d'autres textes. À ce stade, nous avons contacté des experts afin de procéder à la validation des listes des termes complexes et des variantes. La collaboration avec les experts s'est avérée fondamentale pour arriver à des conclusions sur les cas ambigus où les variantes anaphoriques pouvaient être confondues avec des hyperonymes du terme complexe.

En outre, pendant que les experts se concentraient sur la désambiguïsation des variantes, nous avons créé un instrument pour la catégorisation des textes. Cette plateforme informatique présente les textes classés selon le terme complexe détecté par *Acabit* suivi d'une série de variantes éventuelles. Dans le cas de la réduction lexicale, le premier problème que nous avons rencontré concernait les typologies de termes complexes qu'il fallait traiter pour la validation des résultats. Nous avons décidé

de traiter des termes complexes composés d'au moins trois composants, considérés comme pertinents dans le cadre de notre étude.

Dans ce cas, les experts ont été contactés, en premier lieu, pour la confirmation des relations réductionnelles entre les termes complexes et les variantes détectées, comme pour la réduction anaphorique. Ensuite, leur avis a été nécessaire pour une évaluation des composants qui subissaient le mécanisme de réduction. Cette deuxième problématique dérivait de l'analyse des langues de spécialité moins techniques, notamment par les composants adverbiaux. Leur effacement ne détermine pas de réduction, puisque la présence des adverbes détient une fonction de modulation de la valeur sémantique exprimée. Leur chute change l'intensité du concept exprimé, par exemple *dimensione propriamente etica* et *dimensione etica*.

Concernant le choix des experts des domaines, nous nous sommes adressés à des experts différents selon la spécialité. Pour ce qui concerne le domaine des espaces naturels, nous avons contacté des organismes de recherche responsables de la base de données multilingue concernant les réseaux naturels transalpins. Pour le domaine vétérinaire, nous nous sommes adressés aux organismes de contrôles des activités agricoles biologiques, notamment ceux qui évaluent les élevages. Ces organismes avaient déjà été consultés dans le cadre d'un projet concernant un glossaire multilingue sur les activités biologiques. Ils ont cumulé une expérience décennale dans l'évaluation des entreprises du domaine au niveau européen.

Pour ce qui concerne l'évaluation des termes appartenant au domaine de la philosophie, nous avons contacté des chercheurs universitaires qui participent activement à la recherche dans les domaines spécialisés en philosophie et sciences humaines.

Enfin, dans le cadre de la médecine, nous avons collaboré avec les associations de divulgation de l'information en médecine citées dans les sources consultées pour la création de nos corpus. D'autre part, dans le cas des emballages, nous avons contacté les entreprises citées dans les revues prises en considération dans nos corpus pour avoir leur avis sur la terminologie utilisée.

5. La réduction anaphorique

5.1. Hypothèses

Les hypothèses que nous avons décidé de valider au niveau informatique sont le résultat de réflexions concernant les facteurs qui nous venons de décrire au niveau linguistique. Ainsi, nous formulons les hypothèses suivantes.

(1) La réduction anaphorique est corrélée avec le type de discours, selon la classification donnée ci-dessus. En particulier, elle est plus présente dans les textes de catégorie 1, moins présente dans ceux de catégorie 2, quasi absente dans ceux de catégorie 3.

(2) La propriété (1) est indépendante du domaine et de la langue.

(3) S'agissant d'un phénomène anaphorique, dans un document les formes pleines ont tendance à apparaître avant les formes réduites (qui se réfèrent à elles).

5.2. Validation expérimentale

5.2.1. Protocole utilisé

La première hypothèse est liée à des propriétés internes au texte. Pour la valider nous avons étudié les occurrences d'un certain nombre de termes complexes et de leurs réductions anaphoriques dans un corpus composé de documents tirés du Web.

Voici les cinq étapes du protocole utilisé (cf. figure 3) :

- 1) choix de termes complexes dans un domaine précis de langue de spécialité (cf. § 5.2.2);
- 2) récupération des documents contenant la forme pleine de chaque terme (cf. § 5.2.3);
- 3) extraction des formes pleines ainsi que réduites (potentielles) des termes dans ces documents ;
- 4) validation par un expert de l'appartenance de chaque document au domaine, et balisage des formes réduites relevant de la réduction anaphorique ;
- 5) calculs.

Les experts ont été sélectionnés selon les domaines de spécialité. Leur contribution a été assistée par des terminologies qui pouvaient résoudre les problèmes au niveau épistémologique.

5.2.2. Choix des termes complexes

Pour obtenir une liste pertinente de termes complexes dans des domaines donnés (point 1 du protocole), nous avons choisi comme point de départ le contenu textuel de quatre revues spécialisées², (tableau 2). Comme on peut le voir sur la figure 3, l'extraction de termes complexes a été faite par le logiciel Acabit (Daille, 1996) après balisage POS par TreeTagger et lemmatisation par Flemm (pour le français uniquement). Ces deux logiciels ne fournissant pas l'information du genre des noms, nous avons complété la chaîne de traitement en nous servant de ressources lexicales³ pour introduire cette information, qui améliore l'extraction de termes complexes en permettant la vérification de l'accord en genre.

2. *Bulletin du cancer* <http://www.john-libbey-eurotext.fr/fr/revues/medecine/bdc/sommaire.md>, ISSN 1769-6917; *Vet.journal* <http://www.evsrl.it/vet.journal/>, sans ISSN; *Espaces naturels* http://www.espaces-naturels.fr/a_la_une/la_revue_espaces_naturels, ISSN 1637-9896; *Parchi* <http://www.parks.it/federparchi/rivista/>, sans ISSN.

3. Lexique 3 pour le français (<http://www.lexique.org/>), Morph-it ! pour l'italien (<http://dev.sslmit.unibo.it/linguistics/morph-it.php>).

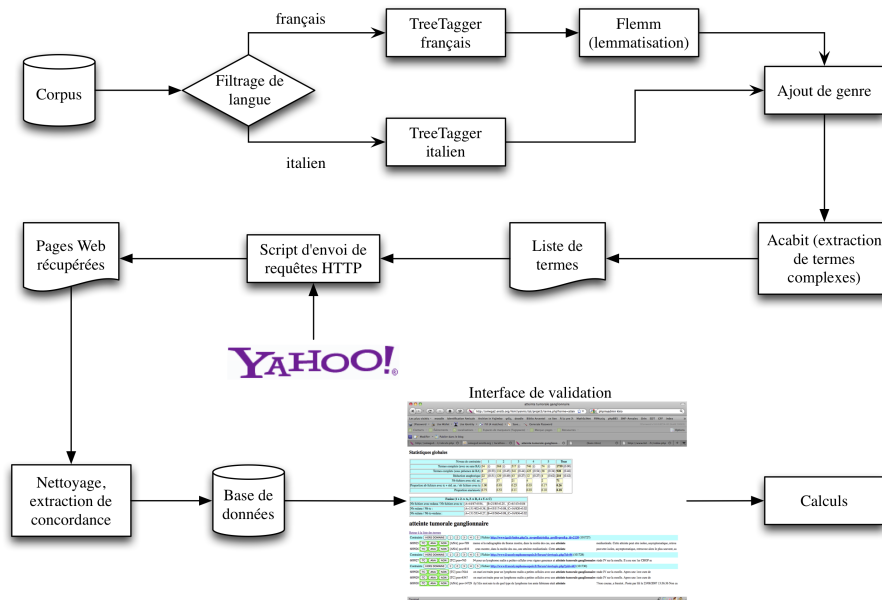


Figure 3. Le protocole utilisé pour la validation des hypothèses de réduction anaphorique

Corpus	Langue	Années	Taille en Mo
<i>Bulletin du cancer</i>	français	1997–2010	43,02 Mo
<i>Journal vétérinaire vet.journal</i>	italien	2003–2010	11,41 Mo
<i>Espaces naturels</i>	français	2003–2009	5,49 Mo
<i>Parchi</i>	italien	1990–2009	18 Mo

Tableau 2. Les revues spécialisées utilisées comme point de départ pour obtenir des listes de termes complexes pour chaque domaine

5.2.3. Récupération des documents Web contenant les termes

Le moteur de recherche Google n’autorisant pas la récupération de plus de soixante-quatre URL par requête, nous nous sommes tournés vers la plate-forme Yahoo BOSS, qui ne pose pas de limite sur le nombre de requêtes. Pour chaque requête Web, Yahoo retourne les mille premières URL, par ordre de pertinence. Nous avons récupéré tous les documents proposés par Yahoo, qui étaient récupérables et de format HTML, PDF ou texte brut, cf. tableau 3. Dans les chiffres ci-dessous, nous entendons par « taille totale » celle des contenus textuels bruts obtenus après conversion et/ou nettoyage.

Domaine	Langue	Nb doc., dont :	HTML	PDF	TXT	Taille totale
Vétérinaire	italien	14 183	10 420	3 743	20	387 Mo
Cancer	français	20 790	13 627	7 124	39	647 Mo
Parcs nationaux	italien	38 502	28 253	10 226	23	1,15 Go
Parcs nationaux	français	18 200	13 840	4 328	32	849 Mo

Tableau 3. Les documents récupérés sur le Web

5.3. Extraction des termes et validation par un expert

Par le biais d'une interface Web spéciale, reliée à une base MySQL, des experts ont pu (a) valider l'appartenance effective de chaque document au domaine donné, (b) indiquer le type du document (1 = vulgarisation, 2 = pédagogique, 3 = spécialisé), (c) vérifier les conditions de forme pleine ou de forme réduite anaphorique de chaque terme complexe observé.

5.3.1. Arbre anaphorique d'un terme complexe

Soit $T = (c_*)$ un texte (où c_* sont des caractères), t_* les occurrences de la forme pleine d'un terme complexe dans T , et r_* celles des différentes formes réduites du même terme. L'ordre linéaire des caractères du texte induit un ordre des t_* et des r_* .

Notons $r_{i,1}, \dots, r_{i,n_i}$ les formes réduites placées entre les formes pleines t_i et t_{i+1} (ou entre t_i et la fin du document). Dans le cas où il n'y a pas de tel r_* , on dira que $n_i = 0$. Soit $\text{pos} : T \rightarrow \mathbb{N}$ la fonction qui associe à chaque mot sa « position » dans le texte (on compte le nombre de caractères depuis le début du fichier).

Dans T il peut y avoir également des formes réduites r'_* placées *avant* la première forme pleine (c'est-à-dire des formes telles que $\text{pos}(r'_*) < \text{pos}(t_1)$). Celles-ci sont appelées réductions cataphoriques et on notera r'_j la j -ième réduction cataphorique de la forme pleine t_1 (pour les autres t_i on considère, dans ce modèle⁴, que l'on n'a que des réductions anaphoriques).

Les occurrences t_* , $r_{*,*}$ et r'_* forment une structure d'arbre ordonné et pondéré par la fonction pos (cf. figure 4), que nous appelons *arbre anaphorique du terme t* .

Pour étudier cette structure et en tirer des renseignements sur le comportement des réductions anaphoriques de T , nous allons nous intéresser à deux types de quantités :

- 1) des quantités relatives à la structure de l'arbre :

4. On peut imaginer d'autres modélisations de ce phénomène. Ainsi, par exemple, on pourrait « attacher » chaque réduction à la forme pleine la plus proche, que ce soit avant ou après elle. On aurait alors des réductions anaphoriques *et* cataphoriques dans tout le document. Dans cet article, nous avons choisi de n'avoir de réduction cataphorique qu'avant la première forme pleine. Cf. aussi § 7, piste 4.

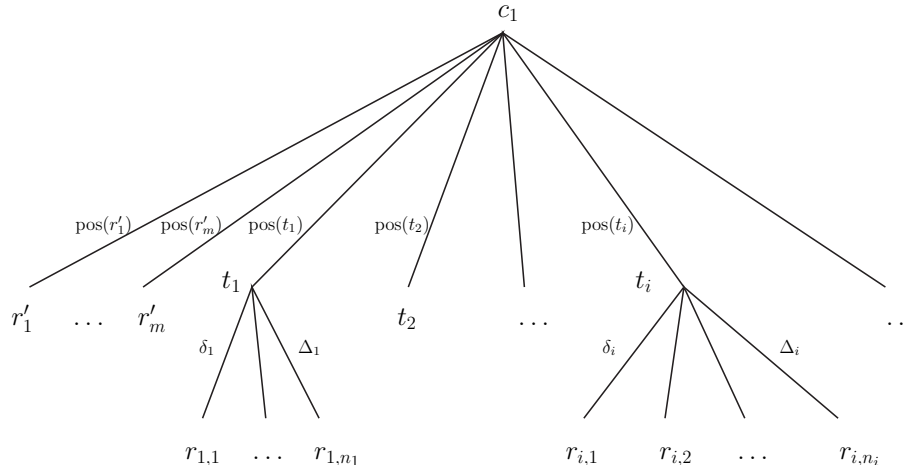


Figure 4. L'arbre des réductions anaphoriques et cataphoriques d'un terme t dont les t_* (resp. $r_{*,*}, r'_*$) sont des occurrences de forme pleine (resp. de forme réduite)

- a) d_m , le degré moyen des nœuds t_i ,
- b) d_- , le nombre de réductions cataphoriques,
- c) f , la moyenne des $f_i = \#\{t_j, t_{j+1}, \dots, t_k \mid d(t_{j-1}) > 0, d(t_\ell) = 0 \text{ pour } i \leq \ell < k \text{ et } d(t_k) > 0\}$ (où d est le degré), c'est-à-dire le nombre de formes pleines consécutives entre deux formes réduites ;

2) des quantités relatives à la pondération de l'arbre :

- a) δ , la moyenne des $\delta_i = \text{pos}(r_{i,1}) - \text{pos}(t_i)$, c'est-à-dire la distance moyenne entre une forme pleine et sa première réduction anaphorique,
- b) Δ , la moyenne des $\Delta_i = \text{pos}(r_{i,n_i}) - \text{pos}(t_i)$, c'est-à-dire la distance moyenne entre une forme pleine et la dernière réduction anaphorique avant la prochaine forme pleine ou avant la fin du fichier,
- c) $\delta_- = \text{pos}(t_1) - \text{pos}(r'_m)$, c'est-à-dire la distance entre la dernière réduction cataphorique et la première forme pleine,
- d) $\Delta_- = \text{pos}(t_1) - \text{pos}(r'_1)$, c'est-à-dire la distance entre la première réduction cataphorique et la première forme pleine.

Le scénario suivant fournit une interprétation possible de ces quantités : après avoir écrit (en moyenne) f formes pleines, l'auteur se sert (en moyenne) de d_m formes réduites, dont la première arrive (en moyenne) à δ caractères après la dernière forme pleine. Éloigné de Δ caractères (en moyenne) de la forme pleine, il considère que l'ambiguïté est devenue trop forte, il revient à la forme pleine, et le cycle reprend.

	Vétérinaire (italien)			Cancer (français)		
	vulg.	pédag.	spécial.	vulg.	pédag.	spécial.
FP	3,20	4,09	5,16	4,97	10,96	11,66
ANA/FP	0,8	0,56	0,06	0,79	0,7	0,25
CATA/FP	0,13	0,01	0,00	0,12	0,02	0,00
$\bar{\delta}$	3 297,16	2 074,05	3 762,68	1 917,76	2 158,51	8 118,12
$\bar{\Delta}$	6 579,31	6 797,56	4 673,45	4 485,39	1 0484,02	2 2292,91
\bar{d}_m	2,58	2,89	1,57	2,23	3,17	5,85
$\bar{\delta}_-$	2 129,82	688,67	NA	3 043,83	1 022,45	NA
$\bar{\Delta}_-$	2 805,29	688,67	NA	6 673,67	4 907,64	NA
\bar{d}_-	1,75	1,00	NA	1,50	1,73	NA
\bar{f}	1,75	2,42	3,39	1,72	1,99	7,38

	Parcs (italien)			Parcs (français)		
	vulg.	pédag.	spécial.	vulg.	pédag.	spécial.
FP	4,58	10,62	12,04	17,89	35,75	52,81
ANA/FP	0,73	0,55	0,19	1,03	0,50	0,17
CATA/FP	1,05	0,34	0,04	3,56	0,98	0,13
$\bar{\delta}$	5 233,06	2 293,71	849,67	1 773,01	972,44	993,31
$\bar{\Delta}$	7 499,35	2 832,36	1 0789,33	5 964,44	1 553,90	3 817,10
\bar{d}_m	2,42	1,93	3,33	3,78	1,24	2,33
$\bar{\delta}_-$	1 896,54	9 931,29	2 433,00	3 340,20	NA	NA
$\bar{\Delta}_-$	5 216,58	1 3449,50	4 811,00	1 3454,80	NA	NA
\bar{d}_-	3,46	1,71	3,00	5,40	NA	NA
\bar{f}	2,41	3,81	3,17	3,35	3,47	7,25

Tableau 4. Résultats de l'analyse des documents récupérés sur le Web. FP, ANA et CATA sont les nombres moyens de formes pleines (resp. anaphoriques, cataphoriques) par 100 Ko de texte, et ANA/FP et CATA/FP les ratios entre formes anaphoriques (resp. cataphoriques) et formes pleines. NA signifie que la donnée n'est pas calculable, faute de cas attestés dans les corpus.

De même, pour le cas cataphorique, l'auteur écrit en moyenne \bar{d}_- formes réduites cataphoriques avant la première forme pleine, qui est à une distance de $\bar{\Delta}_-$ caractères de la première et de $\bar{\delta}_-$ caractères de la dernière forme cataphorique.

5.3.2. Première et deuxième hypothèses

Notre première hypothèse stipule que la quantité et le comportement des réductions anaphoriques et cataphoriques dépendent de la typologie de texte. Nous avons classé les documents en trois catégories (vulgarisation, pédagogique, spécialisé). Le lecteur trouvera dans le tableau 4 les moyennes $\bar{\delta}$, $\bar{\Delta}$, \bar{d}_m , $\bar{\delta}_-$, $\bar{\Delta}_-$, \bar{d}_- , \bar{f} de δ , Δ , d_m , δ_- , Δ_- , d_- et f pour chacun des quatre corpus, ventilés par typologie de texte.

De ces tableaux nous tirons les conclusions suivantes :

1) comme le montre la ligne FP, la densité de formes pleines est croissante lorsque le niveau de contrainte augmente (en allant de la vulgarisation au texte scientifique spécialisé). En effet, un texte scientifique évitera l'ambiguïté en favorisant les formes pleines, au détriment de l'économie ;

2) les lignes ANA/FP et CATA/FP montrent clairement que le nombre de réductions baisse lorsque l'on passe de la catégorie 1 aux catégories 2 et 3. C'est ainsi qu'est démontrée la première hypothèse ;

3) ce comportement est similaire dans les quatre corpus, appartenant à des domaines et à des langues différents, ce qui confirme la deuxième hypothèse ;

4) en ce qui concerne les autres paramètres, on ne peut en tirer aucune conclusion, puisqu'on ne constate aucune régularité significative.

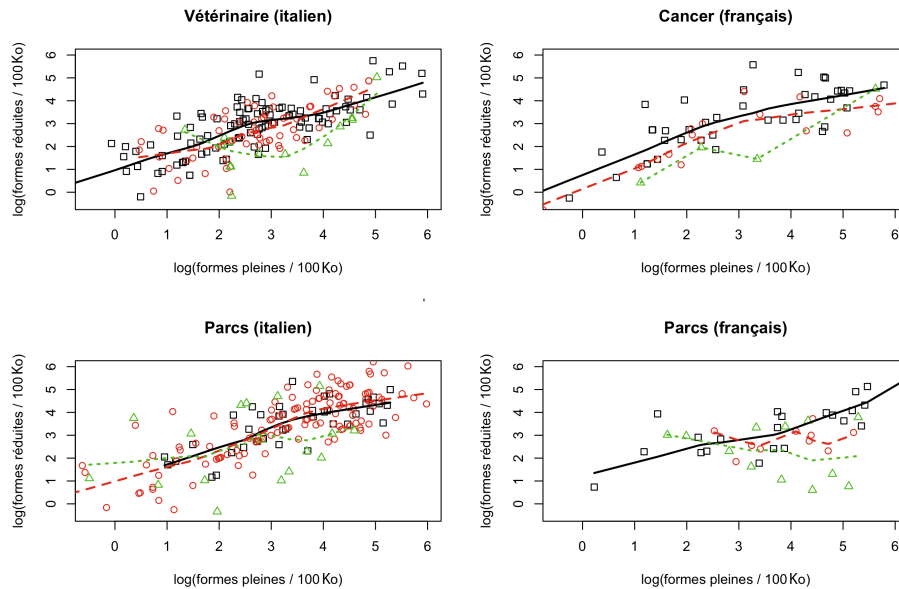


Figure 5. Distribution des documents selon le rapport formes pleines / formes réduites. Les symboles dénotent les trois catégories de texte : carré = cat. 1, cercle = cat. 2, triangle = cat. 3. Les courbes sont des régressions polynomiales LOWESS : courbe continue = cat. 1, tiretée = cat. 2, en pointillé = cat. 3.

Dans la suite, nous allons raisonner au niveau des documents.

La figure 5 montre la distribution des documents selon le rapport formes réduites / formes pleines. Les courbes tracées sont des régressions polynomiales des points, selon l'algorithme LOWESS (Cleveland, 1981). Elles confirment le fait que l'on a plus de réductions anaphoriques en catégorie 1 (vulgarisation) qu'en catégorie 2 (scientifique pédagogique) que 3 (scientifique spécialisé).

Nous constatons également une différence entre les textes français et italiens. Dans le premier cas on remarque un meilleur regroupement : il y a moins de dispersion pour les catégories 1 et 2 et la catégorie 3 se trouve plutôt sous la courbe de la catégorie 1. Dans le cas italien, il y a beaucoup plus de dispersion et aucun regroupement n'est possible. On peut en déduire que la langue italienne est terminologiquement moins stable que la langue française : les termes italiens sont moins figés et leur comportement est plus chaotique, l'alternance entre forme complète et forme réduite a tendance à ne pas suivre un schéma régulier. En français, en revanche, on trouve que les termes sont plus réguliers.

5.3.3. Troisième hypothèse

Notre troisième hypothèse peut être reformulée de la manière suivante : il existe moins de documents avec au moins une réduction cataphorique que de documents avec au moins une réduction anaphorique mais sans aucune réduction cataphorique.

Cette version de l'hypothèse se vérifie aisément à partir de nos données (toutes catégories de texte confondues), cf. tableau 5.

Domaine	Langue	RA	RCA
Vétérinaire	italien	37,47 %	7,00 %
Cancer	français	11,9 %	1,49 %
Parcs	italien	48,48 %	4,71 %
Parcs	français	40,95 %	4,76 %

Tableau 5. *Présence de réductions anaphoriques et cataphoriques dans les documents. La colonne RA (resp. RCA) représente les pourcentages de documents avec au moins une réduction anaphorique (resp. cataphorique).*

6. La réduction lexicale

6.1. Hypothèses

Nous formulons les hypothèses suivantes.

(1) La réduction lexicale dépend du domaine, et est plus fréquente dans les domaines techniques à évolution rapide.

(2) Elle résulte de l'inutilité progressive d'un composant et peut aboutir à une lexicalisation à part entière de la forme réduite. Dans ce cas, on peut observer un « cycle de vie » de la forme réduite : naissance, mise en concurrence avec la forme pleine, remplacement de la forme pleine.

Pour le point (2), nous nous sommes inspirés des articles de Dury & Drouin sur la nécrologie des termes (Dury et Drouin, 2010) et de Møller (1998) sur la terminochronie.

Revue	Domaine	Langue	Années	Volume
<i>Klēsis</i>	philosophie	français	2006–2010	7,25 Mo
<i>Dialegesthai</i>	philosophie	italien	1999–2010	19,43 Mo
<i>Emballages Magazine</i>	emballages	français	2002–2009	1,16 Mo
<i>Italia Imballaggio</i>	emballages	italien	2002–2007	14,05 Mo
<i>Espaces naturels</i>	parcs	français	2003–2009	5,49 Mo
<i>Parchi</i>	parcs	italien	1990–2009	18 Mo

Tableau 6. Liste comparative des revues utilisées pour l'étude de la réduction lexicale

Dans cette section, nous appellerons « termes 3-complexes », les termes dont les composants sont formés de plus de deux mots.

6.2. Validation expérimentale

Contrairement à la réduction anaphorique qui est interne au texte, la réduction lexicale est un mécanisme discursif. Il est donc important d'utiliser des corpus thématiquement stables pour l'étudier. D'ailleurs, la confirmation ou infirmation de la deuxième hypothèse nécessite une datation précise des données textuelles. Ainsi, plutôt que de récupérer des documents arbitraires du Web contenant la forme pleine et des éventuelles formes réduites, nous avons choisi de travailler sur des corpus plus « conventionnels » : nous avons étudié le contenu textuel des six revues spécialisées, dans les domaines suivants : les techniques d'emballage, les parcs nationaux et la philosophie (qui nous servira de contre-exemple puisqu'elle est l'antithèse même d'un domaine « technique à évolution rapide »).

Le lecteur trouvera dans le tableau 6 la liste comparative des revues⁵ considérées. Voici les étapes du protocole utilisé :

- 1) extraction des termes 3-complexes à l'aide du logiciel Acabit (cf. § 5.2.2) ;
- 2) obtention des formes réduites lexicales potentielles par transformation des formes pleines ($A B C \rightarrow A C$, $A B C D \rightarrow A B D$, $A B C D \rightarrow A C D$, etc.) ;
- 3) analyse de la distribution des occurrences des formes pleines du (1) et des formes réduites lexicales du (2) dans le corpus, en tenant compte de leur datation et de leurs positions dans les documents.

5. *Klēsis* <http://www.revue-klesis.org/>, ISSN 1954-3050; *Dialegesthai* <http://mondodomani.org/dialegesthai/>, ISSN 1128-5478; *Espaces naturels* http://www.espaces-naturels.fr/a_la_une/la_revue_espaces_naturels, ISSN 1637-9896; *Parchi* <http://www.parks.it/federparchi/rivista/>, sans ISSN; *Emballages Magazine* <http://www.industrie.com/emballage/>, ISSN 0013-6573; *Italia Imballaggio* http://www.italiaimballaggio.it/italiaimballaggio/05_00/index.html, ISSN 2037-2183.

6.2.1. *Extraction des termes 3-complexes*

Le nombre de termes 3-complexes varie énormément d'un corpus à l'autre. Nous n'avons pas tenu compte des termes complexes de la forme « nom adverbe adjectif(s) ». Ceux-ci ne sont pas pertinents pour notre étude puisque l'adverbe en modifie le sens. Ainsi, la variation qui consiste à omettre l'adverbe ne constitue pas une forme réduite lexicale du terme, puisque elle ne renvoie pas au même concept.

Le nombre de termes 3-complexes ainsi obtenus, triés par corpus, est donné dans le tableau 7. Nous constatons dans le tableau 7 que le nombre de formes pleines et

Revue	<i>t</i>	<i>r</i>	Occurrences <i>t</i>	Occurrences <i>r</i>
<i>Klēsīs</i>	2	3	3	72
<i>Dialegesthai</i>	6	6	12	293
<i>Emballages Magazine</i>	33	15	93	111
<i>Italia Imballaggio</i>	127	96	218	2 015
<i>Espaces naturels</i>	51	45	209	1 576
<i>Parchi</i>	88	52	1 829	24 953

Tableau 7. *Termes 3-complexes obtenus selon le corpus. *t*, *r* dénotent, respectivement, le nombre de termes complexes distincts obtenus (étape 1) et celui de formes réduites (étape 2) attestées dans le corpus. D'autre part « Occurrences *t* » (resp. « Occurrences *r* ») dénote le nombre d'occurrences de formes pleines (resp. réduites) dans le corpus.*

réduites n'est pas directement corrélé avec le nombre d'occurrences de celles-ci dans le corpus : ainsi, dans *Parchi*, on a 24 953 occurrences des 52 formes réduites attestées, alors que dans *Italia Imballaggio* on a presque deux fois plus de formes réduites attestées (96), avec douze fois moins d'occurrences (2 015).

Sur la figure 6 on présente la densité des formes pleines (courbe bleue en pointillé) et réduites (courbe continue) en fonction du temps. On voit que les formes réduites sont bien plus nombreuses que les formes pleines, phénomène dû à leur économicité. D'ailleurs dans les cas de *Espaces naturels* et *Parchi* on assiste à une croissance quasi constante de la densité de formes réduites et à une décroissance du nombre de formes pleines, au fil du temps.

Ces exemples montrent que la réduction lexicale est corrélée à la technicité du domaine, ce qui prouve la première hypothèse.

6.2.2. *Deuxième hypothèse, cycle de vie de terme*

Selon Dury et Drouin (2010) et Møller (1998), les termes complexes ont un cycle de vie lié à celui de l'objet qu'ils dénomment. Lorsque l'objet est encore peu connu, la forme pleine est indispensable pour le dénommer. Mais dans les cas où, progressivement, l'objet se répand, les réductions lexicales deviennent possibles, puisque le risque d'ambiguïté est moindre. Enfin, dans les cas où certains constituants de la forme pleine n'ont aucune charge sémantique réelle, la réduction lexicale finit par remplacer complètement la forme pleine.

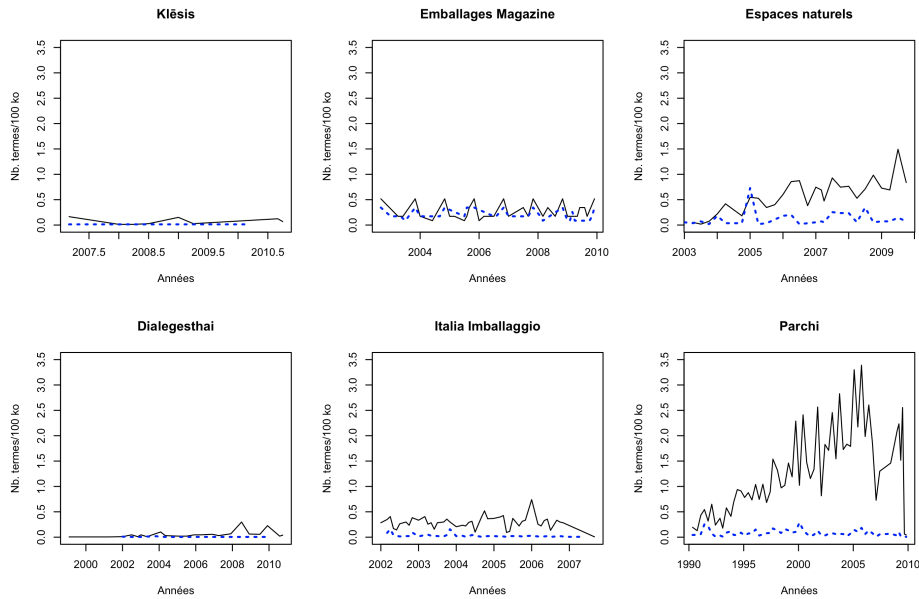


Figure 6. Densité des formes pleines (courbe en pointillé) et réduites (courbe continue) en fonction du temps

Pour déceler des traces du phénomène de « cycle de vie de terme », il a d’abord fallu « unifier » les mesures de datation et de position dans le texte. Pour cela, nous avons suivi l’approche suivante :

1) nous définissons un corpus \mathcal{C} comme étant une suite de N documents $(D_i)_{1 \leq i \leq N}$ datés. La datation est une fonction strictement croissante $T : \mathcal{C} \rightarrow \mathbb{Q}^N$ où $T(D_i) = \text{année}(D_i) + \frac{\text{mois}(D_i) - 1}{12}$;

2) en considérant un document D_i comme une suite de caractères $(c_{i,j})$ nous définissons la *fonction de datation généralisée* T^* comme la fonction linéaire par morceaux qui étend T de au niveau des caractères :

$$\begin{array}{ccc}
 (D_*) & \xrightarrow{T} & \mathbb{Q}^N \\
 \downarrow \pi & \nearrow T^* & \\
 (c_{*,*}) & &
 \end{array}$$

où π est la projection qui envoie le document D_i sur son premier caractère $c_{i,1}$.

Autrement dit, si le document D_5 date du 1^{er} mai 2005 et D_6 du 1^{er} juin de la même année, alors $T(D_5) = 2005 + \frac{5-1}{12} = 2005,333$, $T(D_6) = 2005,417$;

si une occurrence de terme se trouve à la position 37 238 de D_i (dont la taille est, par exemple, de 57 642 caractères) alors sa valeur de datation généralisée sera $T^*(t_i) = T(D_5) + \frac{37\,238}{57\,642} \cdot (T(D_6) - T(D_5)) = 2\,005,3873$. Ainsi toute occurrence a une valeur de datation unique, compatible avec la datation des fichiers et avec l'ordre linéaire du texte⁶ ;

3) la fonction de datation généralisée nous permet de définir une distance temporelle d_T entre les occurrences de termes dans le corpus tout entier : $d_T(t_i, t_j) = T^*(t_j) - T^*(t_i)$;

4) ainsi, pour un terme donné, les occurrences de formes pleines p_* et de formes réduites r_* deviennent des ensembles ordonnés de points de la droite temporelle.

La modélisation étant faite, posons-nous le problème de représentation du phénomène de « cycle de vie de terme ». On ne peut prendre simplement les premières occurrences de formes pleine et réduite, puisque celles-ci peuvent très bien être des *outliers* (données aberrantes). Prendre la moyenne arithmétique des valeurs temporelles de toutes les formes pleines (resp. réduites) ne serait pas une solution non plus, puisqu'on s'intéresse aux débuts de l'apparition d'une forme et non pas à son historique complet. Il convient de traiter séparément les cas où le nombre d'occurrences prend de l'ampleur et ceux où il reste limité. Dans le premier cas, on limitera le nombre d'occurrences à considérer, dans le deuxième cas on prendra la totalité des occurrences.

Comme nous voulons éviter l'aberration causée par des éventuels *outliers*, tout en gardant un maximum d'information pour le cas où le nombre d'occurrences est très limité, nous avons choisi de considérer la moyenne géométrique des N premières occurrences. En effet, la moyenne géométrique est mieux indiquée pour diminuer l'impact des *outliers* que la moyenne arithmétique. Dans les résultats présentés ci-dessous, on a pris $N = 100$.

Nous allons donc calculer, pour chaque terme de chaque corpus, la quantité $\xi = \bar{r} - \bar{t}$, où \bar{r} et \bar{t} sont les moyennes géométriques des 100 premières occurrences des formes réduites (resp. de la forme pleine) du terme. La figure 7 représente la densité de ξ . On constate que la médiane de cette densité est positive dans tous les corpus, ce qui valide notre deuxième hypothèse.

7. Conclusions et perspectives

En traitant des corpus dans divers domaines et dans les deux langues (français et italien) nous avons validé expérimentalement nos hypothèses : que la réduction anaphorique dépend du type de discours, mais ne dépend pas du domaine ou de la langue, que les formes réduites anaphoriques ont tendance à suivre les formes pleines,

6. Suite à la remarque d'un relecteur de l'article, notons que cette fonction convient à des corpus figés. En effet, dans le cas d'un corpus dynamique, l'ajout de documents supplémentaires changerait les valeurs obtenues pour les anciens documents.

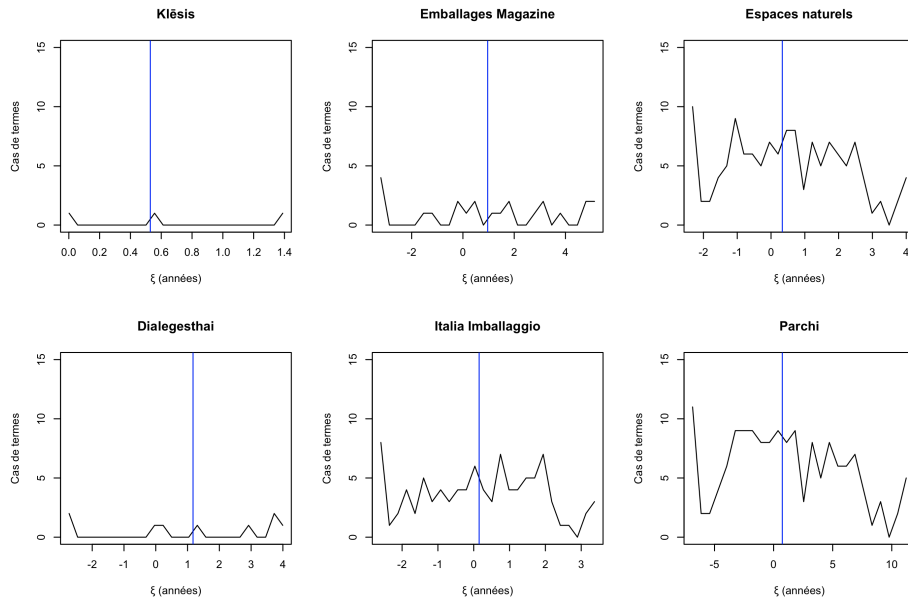


Figure 7. Densité des valeurs de ξ . La ligne bleue indique l'emplacement de la médiane

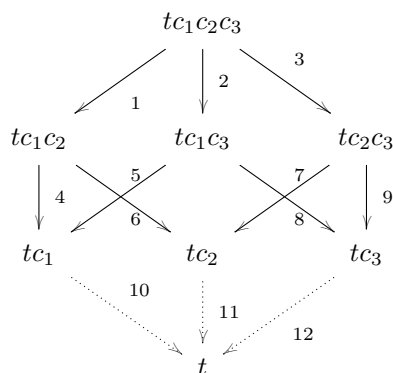
plutôt que de les précéder, que la réduction lexicale dépend du domaine, et est plus fréquente dans les domaines techniques à évolution rapide, et que les formes pleines suivent un cycle de vie et sont souvent remplacées par des formes réduites lexicales.

Les logiciels de type *Acabit* (Daille, 1996) ou *TerMine* (Frantzi *et al.*, 1998) extraient efficacement des termes complexes à partir de données textuelles en langue de spécialité, mais ne nous fournissent aucune indication sur les éventuelles relations sémantiques entre eux. Et pourtant, dans les langues de spécialité les relations d'hyponymie et hyperonymie entre les termes sont d'une importance capitale.

Dans ce travail nous avons essayé de fournir quelques indications sur la probabilité d'avoir, dans certains cas que nous nous proposons de décrire, des relations (quasi) synonymiques.

Un terme complexe peut être décrit en tant que $tc_1 \dots c_n$ où t est la tête et c_i les composants. Prenons, par exemple, $n = 3$. En réduisant le terme $tc_1c_2c_3$ nous

obtenons le treillis suivant (où les flèches sont des réductions purement formelles, sans aucune hypothèse sur les relations sémantiques entre les termes) :



Selon les définitions données dans la première partie de cet article, si réduction anaphorique (ou cataphorique) il y a, elle correspond forcément aux flèches 1, 4 et 10, c'est-à-dire le bord gauche du treillis. Les autres flèches du treillis peuvent, en revanche, être des réductions lexicales.

Les hypothèses que nous avons démontrées expérimentalement nous ont conduits à formuler les règles suivantes :

1) si une réduction est du type $tc_1 \dots c_n \rightarrow tc_1 \dots c_{n-1}$, si la forme pleine $tc_1 \dots c_n$ apparaît avant la forme réduite $tc_1 \dots c_{n-1}$, et si le texte est de catégorie 1 (vulgarisation), alors il y a des fortes chances que la réduction soit anaphorique.

Dans ce cas, s'agissant d'une anaphore, les mécanismes habituels de résolution d'anaphore peuvent être utilisés pour contribuer à la confirmation ou infirmation de l'hypothèse de réduction anaphorique ;

2) pour le même type de réduction, si la forme pleine apparaît après la forme réduite et/ou si le texte est de catégorie 3 (discours scientifique spécialisé ou texte législatif), alors il y a peu de chances que ce soit une réduction anaphorique ;

3) pour tout autre type de réduction dans le treillis ci-dessus, si le texte appartient à un domaine technique à évolution rapide, la possibilité d'une réduction lexicale (due à l'inutilité progressive d'un composant) existe. Pour la confirmer ou l'infirmier, il faudra utiliser des méthodes sémantiques. On pourra, par exemple, former des vecteurs de concepts environnant les deux termes, dont les coefficients seraient pondérés par la fréquence d'occurrences au niveau du corpus. En comparant les vecteurs (par exemple, en calculant leur cosinus), on aurait une indication plus forte sur une éventuelle synonymie.

Partant de là, nombreuses sont les pistes qui peuvent être suivies :

1) en se fondant sur notre corpus annoté, se poser la question de la pertinence de l'application des algorithmes traditionnels de résolution d'anaphore pour confirmer ou infirmer la réduction anaphorique ;

2) idem, pour les mesures de similarité sémantique et la réduction lexicale. Notons que si la résolution d'anaphore et la mesure de similarité s'avèrent pertinentes dans ce domaine, celles-ci étant des techniques lourdes, notre approche peut fournir des critères pour évaluer la pertinence de leur utilisation dans un texte ;

3) étendre notre étude à d'autres langues comme l'anglais (dont les mécanismes de réduction lexicale sont bien différents), l'allemand (qui fourmille de mots composés), les langues idéographiques (où le caractère est porteur de sens et où on assiste à un triangle terme, caractère et concept), etc. ;

4) étudier l'interaction des réductions lexicales et anaphoriques : en effet, nous avons constaté (sans en tenir compte dans nos calculs) le phénomène suivant : une forme pleine (avec arbre anaphorique) était suivie d'une réduction lexicale possédant son propre arbre anaphorique mais dont les feuilles représentaient les mêmes formes que celles des feuilles de l'arbre de la forme pleine. Autrement dit, on peut considérer qu'une forme réduite anaphorique se réfère à une forme pleine ou à une réduction lexicale de celle-ci. Cela nous incite à fusionner les deux arbres (ceux des formes pleine et réduite lexicale) en un seul graphe qui puisse modéliser les deux types de réductions ;

5) étudier la corrélation des différentes quantités définies dans cet article avec la C-valeur de Frantzi *et al.* (1998) ;

6) détecter les éventuels marqueurs anaphoriques dans le voisinage des formes réduites anaphoriques et se servir d'algorithmes de résolution d'anaphore pour obtenir un modèle plus riche de l'arbre anaphorique, dans lequel les arêtes représenteront, dans la mesure du possible, le lien entre anaphore et antécédent ;

7) si l'« indice de réductibilité » est la probabilité qu'un terme complexe soit la réduction (lexicale ou anaphorique) d'un autre terme complexe, développer un algorithme qui calcule cet indice en se fondant sur différents facteurs (type de texte, domaine, position dans le texte, etc.). Cette information serait d'une grande utilité pour l'analyse sémantique, l'indexation ou la traduction automatique.

8. Bibliographie

- Abreu-Garcia J.-M., Une enquête terminologique en espagnol dans le domaine des nouvelles technologies de l'information et de la communication, PhD thesis, Institut Télécom-Télécom Bretagne, 1992.
- Adelstein A., « Condiciones de reductibilidad léxica de los sintagmas terminológicos », *Estudios de Lingüística del Español*, 2002.
- Alber-DeWolf R., Étude sur la création néonymique, analyse comparée des procédés morphologiques et morphosyntaxiques de formation des termes du domaine de la spectroscopie en anglais, en allemand, en français et en russe, PhD thesis, Université Laval, Québec, 1984.
- Alexandru C., Gaudin F., « Les contextes : à la source du terme », *Mots, termes et contextes. Préactes des VII^{es} journées scientifiques AUF-LTT*, 2005.
- Bourigault D., Slodzian M., « Pour une terminologie textuelle », *Terminologies nouvelles*, vol. 19, p. 29-32, 1999.

- Cardero García A. M., « En torno a la frecuencia de algunas estructuras sintácticas en terminología », *Terminología e Industrias da Lingua. Actas do VII Simpósio Ibero Americano de Terminología*, ILTEC, 2003.
- Cleveland W., « LOWESS : A program for smoothing scatterplots by robust locally weighted regression », *The American Statistician*, vol. 35, p. 54, 1981.
- Collet T., « La réduction des unités terminologiques complexes de type syntagmatique », *Méta : journal des traducteurs*, vol. 42, n° 1, p. 193-206, 1997.
- Collet T., La réduction des unités terminologiques complexes de type syntagmatique, PhD thesis, Université de Montréal, 2000.
- Collet T., « Esquisse d'une nouvelle microstructure de dictionnaire spécialisé reflétant la variation en discours du terme syntagmatique », *Méta : journal des traducteurs*, vol. 49, n° 2, p. 247-263, 2004.
- Conceição M. C., « Terminologie textuelle : reformulations et accès aux concepts », *De la mesure dans le termes*, Presses Universitaires de Lyon, p. 269-305, 2005.
- Condamines A., *Linguistique de corpus et terminologie*, CNRS, Équipe de Recherche en Syntaxe et Sémantique, Toulouse, 2005.
- Daille B., « ACABIT : une maquette d'aide à la construction automatique de banques terminologiques monolingues ou bilingues », in A. Clas, P. Thoiron, H. Béjoint (eds), *Lexicomatique et Dictionnairiques*, FMA, Beyrouth, p. 123-136, 1996.
- Daille B., « Conceptual structuring through term variation », *Proceedings ACL 2003 Workshop on Multiword Expressions : Analysis, Acquisition and Treatment*, p. 9-16, 2003.
- Daille B., Gaussier É., Langé J.-M., « An Evaluation of Statistical Scores for Word Association », in J. Ginzburg, Z. Khasidashvili, C. Vogel, J.-J. Levy, E. Vallduvi (eds), *The Tbilisi Symposium on Logic, Language and Computation : Selected Papers*, CSLI Publications, p. 177-188, 1998.
- de Bessé B., « Le domaine », in H. Béjoint, P. Thoiron (eds), *Le sens en terminologie*, Presses universitaires de Lyon, 2000.
- de Vecchi D., Estachy L., « Pragmateterminologie : les verbes et les actions dans les métiers », *Actes de la conférence TOTh 2008*, 2008.
- Delmonte R., « Struttire sintattiche dall'analisi computazionale di corpora di italiano », *Intorno all'italiano contemporaneo. Tra linguistica e didattica*, Franco Angeli, 2004.
- Desmet, « Terminologie, culture et société. Éléments pour une théorie variationniste de la terminologie et des langues de spécialité », *Cahiers du Rifal*, 2001.
- Desmet I., « Variabilité et variation en terminologie et langues spécialisées : discours, textes et contextes », *Septième journées scientifiques du Réseau « Lexicologie, Terminologie, Traduction » de l'Agence universitaire de la Francophonie : Mots, termes et contextes*, Bruxelles, p. 235-247, 2006.
- Dury P., Drouin P., « L'obsolescence des termes en langues de spécialité : une étude semi-automatique de la « nécrologie » en corpus informatisés, appliquée au domaine de l'écologie », *Online proceedings of the XVII European LSP Symposium 2009*, p. 1-11, 2010.
- Frantzi K. T., Ananiadou S., Tsujii J., « The C-value/NC-value method of Automatic Recognition for Multi-Word Terms », *Springer LNCS*, vol. 1513, p. 585-604, 1998.

- Freixa A. J., La variació terminològica : anàlisi de la variació denominativa en textos de diferent grau d'especialització de l'àrea de medi ambient, PhD thesis, Universitat Pompeu Fabra de Barcelone, 2002.
- Guilbert L., *La créativité lexicale*, Larousse, 1975.
- Jacquemin C., « Syntagmatic and Paradigmatic Representations of Term Variation », *ACL '99 Proceedings of the 37th annual meeting of the Association for Computational Linguistics on Computational Linguistics*, p. 341-348, 1999.
- Jacquemin C., *Spotting and Discovering Terms Through Natural Language Processing*, The MIT Press, 2001.
- Jacques M.-P., « L'emploi de termes réduits comme révélateur de la centralité dans le domaine », *Septième journées scientifiques du Réseau « Lexicologie, Terminologie, Traduction » de l'Agence universitaire de la Francophonie : Mots, termes et contextes*, Bruxelles, p. 299-308, 1996.
- Jacques M.-P., « La réduction du syntagme terminologique au fil du discours », *Cahiers de grammaire*, vol. 25, p. 93-114, 2000.
- Jacques M.-P., Approche en discours de la réduction des termes complexes dans les textes spécialisés, PhD thesis, Université de Toulouse, 2003.
- Kocourek R., *La langue française de la technique et de la science*, Oskar Brandstetter, 1991.
- Lavagnino E., Terminologie et variations discursives : la réduction des termes complexes à l'épreuve de la linguistique de corpus, PhD thesis, Université de Gênes, 2011.
- Lerat P., « La combinatoire des termes. Exemple : *nectar de fruits* », *Journal of Language and Communication Studies*, vol. 42, p. 211-232, 2009.
- Møller B., « À la recherche d'une terminochronie », *Méta : journal des traducteurs*, vol. 43, n° 3, p. 426-438, 1998.
- Nenadić G., Ananiadou S., McNaught J., « Enhancing automatic term recognition through recognition of variation », *COLING '04 : Proceedings of the 20th international conference on Computational Linguistics*, 2004.
- Ngonga Ngomo A.-C., « SIGNUM : A Graph Algorithm for Terminology Extraction », *Springer LNCS*, vol. 4919, p. 85-95, 2008.
- Pazienza M. T., Pennacchiotti M., Zanzotto F. M., « Terminology extraction : an analysis of linguistic and statistical approaches », *Studies in Fuzziness and Soft Computing*, 2005.
- Portelance C., « Syntagmes et paradigmes », *Méta : journal des traducteurs*, vol. 34, n° 3, p. 398-404, 1989.
- Portelance C., « Fondements linguistiques de la terminologie », *Méta : journal des traducteurs*, vol. 36, n° 1, p. 64-70, 1991.
- Portelance C., « De la nomination : catégorisation et syntagmatique », *Applied Semiotics/Sémiotique appliquée*, vol. 1, p. 55-63, 1996.
- Prandi, *Le regole e le scelte. Introduzione alla grammatica italiana*, UTET, 2006.
- Sabatini F., « "Rigidità-esplicitezza" vs. "elasticità-implicitezza" : possibili parametri massimi per una tipologia dei testi », in G. Skytte, F. Sabatini (eds), *Linguistica testuale comparativa*, Museum Tusulanum Press, p. 142-172, 1999.
- SanJuan E., Ibekwe-SanJuan F., « Terminologie et classification automatique des textes », *Actes 6^{es} Journées internationales d'Analyse statistique des données textuelles (JADT'2002)*, p. 13-15, 2002.

Savary A., Jacquemin C., « Reducing Information Variation in Text », *Springer Lecture Notes in Artificial Intelligence*, vol. 2705, p. 145-181, 2003.

Soglia S., « Origine, sviluppo e tendenze della terminologia moderna », in M. Magris *et al.* (eds), *Manuale di terminologia*, Hoepli, 2002.

Tran T. M., Trancart M., Servent D., « Littéracie, SMS et troubles spécifiques du langage écrit », in J. Durand, B. Habert, B. Laks (eds), *Congrès mondial de linguistique française*, Institut de Linguistique Française, 2008.