

Cette méthode ne permet aucune gestion de l’ambiguïté². Une suite de caractères est définie comme étant, ou non, une unité lexicale, indépendamment du contexte.

Il existe donc des regroupements incorrects, qui induisent des erreurs d’étiquetage. Ainsi, dans les exemples suivant, *rendez-vous* est systématiquement regroupé et étiqueté nom commun :

- Vous arrivez en retard à tous vos rendez-vous ? (E115)
- Faites comme lui, rendez-vous sourde à tous les cris, rejoignez la pierre pendant qu’il en est temps. (K244)

1 Définition des unités lexicales composées

À l’issue des travaux réalisés par divers membres de l’équipe, les unités composées ont été définies comme l’ensemble des formes appartenant aux sous-ensembles suivant :

1. Les formes de lemmes composés comportant une espace ou une apostrophe dans Morphalou 3 [1], après double validation manuelle selon 4 critères :
 - Si l’un des mots graphiques du composé n’existe pas seul (*jeun, afin, fur*) => une seule unité
 - Si le composé est ambigu et que les mots graphiques peuvent se suivre par ailleurs (*bien que*) => plusieurs unités
 - Si le moule syntaxique ne pose pas problème (*à l’aune de*) => plusieurs unités
 - Si le composé est une expression latine (*post scriptum*) => une seule unité.
2. Les formes composées comportant un trait d’union dans Morphalou 3³.
3. Les formes composées à l’aide de traits d’union commençant par un préfixe répertorié dans GLAWI [2], à l’exception de *très-* et *mon-*, qui ont été identifiés comme étant source de trop d’erreur de regroupement.

²La seconde méthode aurait pu nous offrir cette possibilité, mais aurait nécessité une masse de données importante.

³Morphalou 3 comporte 10 111 lemmes de ce type. Leur vérification manuelle n’a donc pas été réalisée, à la fois pour des questions de temps et de difficulté à maintenir un jugement uniforme sur un nombre aussi important de cas.

4. Les formes composées à l'aide de traits d'union commençant par un préfixe selon une liste additionnelle obtenue à partir des formes composées de Morphalou 3 et disponible en annexe (page 7).
5. Les nombres composés (*trois cent, quatre-vingt-dixième*).
6. Les formes débutant par un trait d'union appartenant à la liste de mots grammaticaux disponible page 8.
7. Les noms propres composés répertoriés dans nos lexiques⁴.

2 Avertissement

Comme le détaille Véronique Montémont [3], la base de données textuelles Frantext a été initiée dans les années 1960, dans le cadre d'un projet lexicographique. Diverses politiques de mécanographie puis de numérisation se sont succédées au cours de sa constitution. D'une façon générale, nous pouvons distinguer les côtes récentes - dont la numérisation a été réalisée depuis 2000 - des côtes plus anciennes.

récentes : dont le nom commence par la lettre B, E ou R ou se situe entre S300 et S999

anciennes : dont le nom commence par la lettre K, L, M, P, Q ou Z ou se situe entre S001 et S299

Dans les côtes les plus anciennes, une particularité majeure peut être relevée comme étant source d'erreurs : le mauvais encodage de tirets d'incise et de dialogue, couplé à l'absence malheureuse d'espace, provoquant une assimilation à des traits d'union.

- la maman interdira-à juste titre-toute autre tentative. (P455)

Certaines heuristiques ont été mises en œuvre pour limiter ces erreurs, mais elles n'ont pas pu être entièrement éliminées.

⁴Remarque 08/02/2018 : Le regroupement systématique des noms propres composés a été revu, 239 regroupements issus de noms de communes homographes de séquences DET + NC, tels que *Le Passage*, ont été abandonnés.

Références

- [1] ATILF. Morphalou, 2015. ORTOLANG (Open Resources and TOols for LANGuage) –www.ortolang.fr.
- [2] Nabil Hathout and Franck Sajous. Wiktionnaire’s Wikicode GLAWified : a Workable French Machine-Readable Dictionary. In Nicoletta Calzolari (Conference Chair), Khalid Choukri, Thierry Declerck, Marko Grobelnik, Bente Maegaard, Joseph Mariani, Asuncion Moreno, Jan Odijk, and Stelios Piperidis, editors, *Proceedings of the Tenth International Conference on Language Resources and Evaluation (LREC 2016)*, Portorož, Slovenia, may 2016. European Language Resources Association (ELRA).
- [3] Véronique Montémont. Bâtir une ressource lexicale : l’aventure Frantext. In Jean-Paul Bravard, Natalia Del Fatti, Ioan Negruțiu, and Cristina Vieira, editors, *Les ressources*, Les colloques de l’Institut universitaire de France [2], Saint-Étienne, 2011. Publications de l’Université de Saint-Étienne. Textes issus du colloque du 20e anniversaire de L’Institut universitaire de France, tenu à École nationale supérieure de Lyon du 30 au 31 mai 2011.
- [4] Assaf Urieli. *Robust French syntax analysis : reconciling statistical methods and linguistic knowledge in the Talismane toolkit*. PhD thesis, Université de Toulouse II le Mirail, 2013.

Annexes

A Lemmes de Morphalou 3 comportant une espace conservés

a b c	à partir de	à vau-l’eau
à brûle-pourpoint	a posteriori	à vrai dire
à califourchon	a priori	ab absurdo
a contrario	à priori	ab origine
à donf	à qui mieux mieux	ad hoc
a fortiori	à quia	ad hominem
à jeun	à tâtons	ad libitum
a latere	a tempera	ad limina
a minima	à tue-tête	ad litem

ad litteram	de plano	in pace
ad nauseam	de profundis	in petto
ad nutum	de visu	in praesenti
ad patres	delirium tremens	in praesentia
ad personam	deo gratias	in sillico
ad rem	deus ex machina	in situ
ad valorem	dies irae	in spiritu
ad vitam aeternam	don juan	in utero
afin de	don quichotte	in vitro
afin que	don quichottesque	in vivo
alter ego	don quichottisme	intra muros
au demeurant	don-juanisme	ipso facto
au fur et à mesure	ecce homo	joint venture
au visé	eh ben	kyrie eleison
auto-destructif	eh bien	lato sensu
b.a. ba	eh quoi	lingua franca
bel canto	en catimini	loc. cit.
bernard l'ermite	en delà de	locus citatus
bernard l'hermite	en outre	magister dixit
bernard lermite	en stand-by	manu militari
bernard lhermite	et caetera	mass media
bric-à-brac	ex abrupto	mea culpa
buen retiro	ex aequo	mezza voce
c'est à dire	ex cathedra	minus habens
cardio-péricardo-	ex nihilo	missi dominici
myopexie de gor	ex professo	modern style
carpe diem	ex voto	modus operandi
casus belli	exempli gratia	modus vivendi
compact disk	extra muros	motu proprio
compte rendu	garde champêtre	mutatis mutandis
cosa nostra	gratis pro deo	ne varietur
d'ores et déjà	grosso modo	nec plus ultra
dare dare	habeas corpus	nihil obstat
de auditu	hasta luego	no man's land
de bric et de broc	hic et nunc	noli me tangere
de commodo	honoris causa	nota bene
de cuius	id est	numerus clausus
de facto	in abstracto	oh la la
de guingois	in extenso	oh là là
de gustibus	in extremis	one woman show
de jure	in fine	one-man show

op. cit.	post scriptum	stabat mater
opéra bouffe	pourvu que	statu quo
operating system	pretium doloris	stricto sensu
opus citatum	prima donna	sui generis
osso buco	pro domo	tandis que
over arm stroke	pro forma	te deum
p et t.	punching ball	terminus a quo
p et t	quant à	terminus ad quem
p. et t	quat'z arts	terra incognita
parce que	quod erat demonstrem	terra rossa
patati et patata	dum	terza rima
patati patata	res militaris	tout de go
pax americana	rock and roll	ultima ratio
pax romana	rock and roller	verbi gratia
per capita	sainte nitouche	vice versa
persona grata	salvé regina	volens nolens
persona non grata	se frotti-frotter	vomito negro
peu ou prou	sedia gestatoria	vox populi
pique niquer	semper virens	vulgum pecus
plan plan	sine die	white spirit
pole- position	sine qua non	y compris
post mortem	souventes fois	

B Lemmes de Morphalou 3 comportant une apostrophe conservés

aujourd'hui	chefs-d'œuvre	entr'ouvrir
baha'i	ct'	fedda'i
baha'isme	d'arrache-pied	fedda'i
bernard l'hermite	d'emblée	hors-d'oeuvre
bernard-l'ermite	d'ores et déjà	hors-d'œuvre
bernard-l'hermite	dos-d'âne	j'm'en-fichiste
bin's	entr'aide	j'm'en-foutisme
c'est à dire	entr'aimer	je-m'en-fichisme
c'est-à-dire	entr'apercevoir	je-m'en-fichiste
c't'	entr'axe	je-m'en-foutisme
ch'timi	entr'égorgement	je-m'en-foutiste
chef-d'oeuvre	entr'égorger	jourd'hui
chef-d'œuvre	entr'ouvert	jourd'huis

jusqu'au-boutisme	no man's land	quô'c-ngũ'
jusqu'au-boutiste	pac'que	R'n'B
l'on	parc'que	rock'n'roll
m'amie	pin's	rock'n'roller
m'as-tu-vu	presqu'île	sot-l'y-laisse
m'as-tu-vue	presqu'île	tape-à-l'oeil
m'as-tuvuisme	presqu'îlette	tape-à-l'œil
main-d'oeuvre	prud'homal	tout-à-l'égout
main-d'œuvre	prud'homie	traveler's
mains-d'œuvre	prud'hommal	traveleur's
mam'selle	prud'homme	traveller's
mam'zelle	qu'en-dira-t-on	trompe-l'oeil
monte-en-l'air	quat'z arts	trompe-l'œil
nid-d'abeilles	quat'zarts	vau-l'eau
nids-d'abeilles	quelqu'un	

C Liste additionnelle de préfixes

abat-	anhydro-	broncholo-
abdomino-	animo-	cap-
aborto-	ano-	casse-
accroche-	antéro-	cellulo-
acido-	apico-	centre-
adipo-	argilo-	centro-
adiposo-	arrière-	cérébello-
adréno-	artério-	cérébro-
adynamico-	assyro-	cervici-
affectivo-	auriculo-	cervico-
aides-	avant-	cholédoco-
aigre-	bacill-	coccy-
aigres-	bacillo-	cortico-
albumino-	basi-	costo-
alcalino-	basio-	cubito-
aluminico-	bien-	égypto-
alvéolo-	bilio-	fisco-
ammoniaco-	bissau-	fluvio-
ammonio-	bronchio-	franc-
analytico-	bronchiolo-	fronto-
anarcho-	broncho-	glucido-

grand-
grande-
grandes-
grands-
grosso-
ilié-
ilio-
maxillo-

nord-
occipito-
ouest-
physico-
sud-
thyréo-
tibio-
touristico-

toxi-
trachélo-
trachéo-
urétéro-
vésico-
vésiculo-

D Liste de mots grammaticaux

-ce
-ci
-elle
-elles
-en
-il
-ils
-je
-la
-là
-le
-les

-leur
-lui
-m'
-m'en
-même
-mêmes
-moi
-nous
-on
-t-elle
-t-elles
-t-en

-t-il
-t-ils
-t-on
-t-y
-t'
-t'en
-toi
-tu
-vous
-y