

Syllabation des représentations phonétiques de Brulex et de Lexique

Christophe Pallier*

17 février 1999 (mis à jour le 12 mai 2004)

Cette note décrit l'algorithme utilisé pour syllaber (ou "syllabifier") les représentations phonétiques des bases lexicales du français Brulex et Lexique.

Disons tout de suite qu'il existe plusieurs théories sur la syllabification des groupes de consonnes. Certaines se fondent sur des critères acoustico-phonétiques, d'autres sont d'inspiration plus phonologique. La table 1 (adaptée de Laeufer (1992) dans le chapitre 6 de ma thèse (Pallier, 1994)) présente six propositions.

Pratiquement toutes les théories s'accordent à syllabifier les groupes occlusive-occlusive (OO) entre les deux consonnes : /capture/ est syllabifié en /cap-ture/. Par contre, les théories sont en désaccord sur le statut des groupes occlusive-fricative (OF) : /capsule/ est-il /ca-psule/ ou /cap-sule/ ?

La syllabation que nous proposons pour Brulex et Lexique n'est donc qu'un choix parmi plusieurs possibles (essentiellement fondé sur nos intuitions personnelles !). Elle correspond à celle de Pulgram sur les exemples de la table 1. Notre idée n'a jamais été d'imposer, avec Lexique, une syllabation. Cette note décrit l'algorithme afin permettre à ceux qui le désirent de corriger selon leurs besoins la syllabation que nous avons proposé.

Si vous utilisez cet algorithme, ou le modifiez, ou détectez des problèmes, je vous serais reconnaissant de me le faire savoir.¹

*<http://www.pallier.org>. Voir aussi <http://www.lexique.org>. Tout commentaire est bien venu.

¹Mon email est disponible sur le site www.pallier.org

Table 1: Différentes propositions de syllabifications de mots français (d'après Laeufer, 1992)

		Grammont	Delattre		Pulgram Malmberg	Noske	Levin
			apt.	force			
OL	<i>caprice</i>	-pr	-pr	-pr	-pr	-pr	-pr
	<i>atlas</i>	-tl	-tl	-tl	t-l	t-l	t-l
ON	<i>technique</i>	-kn	-kn	-kn	k-n	k-n	k-n
OF	<i>adverbe</i>	-dv	-dv	-dv	d-v	d-v	d-v
OO	<i>structure</i>	-kt/k-t	k-t	k-t	k-t	k-t	k-t
FL	<i>casserole</i>	-sr	-sr	-sr	s-r	s-r	-sr
	<i>disloque</i>	-sl	-sl	s-l	s-l	s-l	s-l
	<i>influent</i>	-fl	-fl	-fl	-fl	f-l	-fl
FN	<i>transmis</i>	-sm	s-m	s-m	s-m	s-m	s-m
FF	<i>blasphème</i>	-sf/s-f	s-f	s-f	s-f	s-f	s-f
FO	<i>diphthongue</i>	f-t	f-t	f-t	f-t	f-t	f-t
NL	<i>minerais</i>	-nr	-nr	-nr	n-r	n-r	n-r
NN	<i>calomnie</i>	-mn/m-n	-mn	-mn	m-n	m-n	m-n
NF	<i>hameçon</i>	m-s	-ms	-ms	m-s	m-s	m-s
NO	<i>samedi</i>	m-d	m-d	-md	m-d	m-d	m-d
LL	<i>galerie</i>	-lr/l-r	-lr	-lr	l-r	l-r	l-r
	<i>berlue</i>	-rl/r-l	r-l	r-l	r-l	r-l	r-l
LN	<i>calmant</i>	l-m	l-m	-lm	l-m	l-m	l-m
LF	<i>répulsif</i>	l-s	l-s	-ls	l-s	l-s	l-s
LO	<i>culbute</i>	l-b	l-b	-lb	l-b	l-b	l-b

O = occlusives ; F = fricatives ; N = nasales ; L = liquides

Table 2: Règles de syllabation simplifiées

Chaîne	→	Syllabation	Exemple(s)
VV	→	V-V	poete [po-et]
VCV, VLV, VYV	→	V-CV	cadeau [ca-do], calot [ca-lo], aboyer [a-büa-ïé]
VYYV	→	VY-YV	bouilloire [buï-üaR]
VCYV, VLYV	→	V-CYV	tatouer [ta-tüé], polluer [po-lÿé]
V [td]RV V	→	V-CCV	autrement [o-tre-mâ]
VCLV (où C ∈ [pkbgfs/vzj])	→	V-CLV	caprice [ka-pRis]
VCCV, VLCV, VLLV, VYCV, VYLV	→	VC-CV	capture [kap-ture], costume [kos-tum], galerie [gal-Ri], vieillerie, [vi_ï-Ri], atlas [at-las], madeleine [mad-l_n]
VXXXV	→	VX-XXV	astrophysique [as-tRo...]
VXXXXV	→	VX-XXXV	obstruer [op-stRy-é], octroyer [ok-tRüa-ïé]

Notations:

V = voyelles (incluant schwa), *sans* les semi-voyelles

C = toutes les consonnes sauf liquides et semi-voyelles

L = liquides {l, R}

Y = semi-voyelles {ü, ï, ÿ}

X = C ou L ou Y

La table 2 résume les principales règles de syllabation que nous nous proposons d'appliquer. Plusieurs remarques méritent d'être soulignées :

- On est obligé de distinguer les semi-voyelles des voyelles: tuer [tüé] vs clouer [klu-é]
- On est obligé de distinguer les semi-voyelles des liquides: polluer [po-lüé] vs. galerie [gal-Ri]
- On a supprimé les schwa finaux des mots multisyllabiques: notre (notR[^]) → [notR] arbre (aRbR[^]) → [arbR]
- Certains choix de syllabation sont certainement discutables, tels que: stagner → [stag-né], castor → [kas-toR], astro → [as-tRo],

Le script de syllabation fournit en annexe explicite toutes les règles.

Nous n'avons pas systématiquement privilégié la règle de l'attaque maximal ou le principe de sonorité.

Nous avons décidé qu'une syllabe non finale contenait au plus une consonne en coda. Par exemple, 'exploit' devient 'ek-splait'. Ceci est conforme avec la règle de l'attaque maximale même si cela viole le principe de sonorité.

Il nous a aussi fallu prendre une décision vis-à-vis des schwa finaux. Nous avons décidé de les supprimer des représentations phonétiques. Ainsi 'arbre' est pour nous un monosyllabe.

Finalement, notez que les représentations phonétiques de Brulex et de Lexique ne sont pas strictement identiques (cf. par exemple 'fluide' qui contient une voyelle 'u' selon le premier, et un semivoyelle selon le second). C'est pourquoi les syllabations ne sont pas toujours en accord...

Réalisation pratique

L'algorithme est réalisé par un script `syllabation.awk` écrit dans le langage Awk. Le programme libre `gawk` (www.gnu.org/software/gawk/gawk.html) permet de le faire fonctionner.

L'intérêt de Awk est que les expressions régulières permettent d'exprimer les règles de façon très lisible.

Les représentations phonétiques de Brulex et Lexique utilisent des codages différents. Par défaut, le script suppose que le codage de Lexique (en fait de LAIPTTS) est utilisé ; en mettant 'brulex' dans la variable 'code', le cadage de Brulex est employé.

Ce script peut être utilisé interactivement, en tapant simplement :

```
gawk -f syllabation.awk
```

puis en entrant des représentations phonétiques (code Lexique).

Pour syllaber le fichier `brulex.txt`:

```
gawk -vphons=2 -vcode=brulex -f syllabation.awk brulex.txt
```

et pour syllaber `lexique260_graph.txt`:

```
gawk -vphons=2 -f syllabation.awk lexique260_graph.txt
```

Si l'utilitaire 'make' est disponible, il suffit de taper 'make' pour syllaber les deux fichiers.

Le commande 'make test' fournit les syllabations des mots listés dans le fichier `mots_test.txt`. Le résultat est dans le fichier `tests.txt`.

Les fichiers qui accompagnent `syllabation.awk` sont :

- `syllabation.pdf`: description (sommaire) des règles de syllabation

- brulex.txt et lexique260_graph.txt : fichiers contenant les formes orthographique et phonétiques de brulex et la table graphemes de lexique.
- Makefile: règles pour l'utilitaire 'make'.
- mots_test.txt : liste de mots dont on veut verifier les syllabations
- tests.txt : resultats de la syllabation (selon brulex, puis selon lexique)

Code source du script `syllabation.awk`

Voici le script awk qui réalise cette syllabification.

```

1  #!/usr/bin/gawk -f
2  #
3  # This script reads a tab-separated file and syllabifies the columns pointed to
4  #
5  #
6  # Author: Christophe Pallier (christophe.pallier@m4x.org)
7  #
8  # License: GNU (cf. http://www.gnu.org)
9  #
10 # Last update: 13 May 2004
11 # (original date: the first version of this script was written during
12 # my dissertation, in 1993)
13 #
14 # 2004/05/13: merge of the Brulex & Lexique versions
15 #             correction for the 'j'-'>'Z' (viellerie) in Lexique
16 #             add rule '[td]R' (comment on 'autrefois' by Sprenger-Charolles)
17 #
18 #
19 #
20 # Note: changed \377 into y-umlaut to run under DOS (bug gawk).
21
22 BEGIN {
23     FS="\t";
24     OFS="\t";
25
26     if (code=="brulex") {
27         V="[aiouyîâêôû^eEéAO_]"; # vowels
28         C="[ptkbgdgs/vzjmnN£]"; # consonants except liquids & semivowels
29         Cl="[pkbgfs/vzj]";
30         L="[lR]"; # liquids
31         Y="[iü\377]"; # semi-vowels \377 stands for y-umlaut
32         X="[ptkbgdgs/vzjmnN£xlRiü\377]"; # all consonants
33     } else { # code == LAIPTTS)
34         V="[iYeE2591a@oO$uy*]"; # Vowels
35         C="[pbmfvtdnNkgszxsZGh]"; # Consonants except liquids & semivowels
36         Cl="[pkbgfsSvzZ]";

```

```

37     L="[lR]"; # liquids
38     Y="[j8w]"; # semi-vowels
39     X="[pbmfvtdnNkgszSZGlRrhxGj8w]"; # all consonants, including semivowels
40 }
41 if (phons==0) phons=1;
42 }
43
44 {
45     a=$phons;
46     n=1
47 }
48
49 {
50     while (i= match (a, V V)) {
51         a=substr(a,1,i) "-" substr(a,i+1,length(a)); n++; }
52
53     while (i= match(a, V X V)) {
54         a=substr(a,1,i) "-" substr(a,i+1,length(a)); n++;
55
56     while (i=match(a, V Y Y V)) {
57         a=substr(a,1,i+1) "-" substr(a,i+2, length(a)); n++;
58
59     while (i=match(a, V C Y V)) {
60         a=substr(a,1,i) "-" substr(a,i+1, length(a)); n++;
61
62     while (i=match(a, V L Y V)) {
63         a=substr(a,1,i) "-" substr(a,i+1, length(a)); n++;
64
65     while (i=match(a, V "[td]R" V)) {
66         a=substr(a,1,i) "-" substr(a,i+1, length(a)); n++;
67
68     while (i=match(a, V "[td]R" Y V)) {
69         a=substr(a,1,i) "-" substr(a,i+1, length(a)); n++;
70
71     while (i=match(a, V C1 L V)) {
72         a=substr(a,1,i) "-" substr (a,i+1,length(a)); n++;
73
74     while (i=match(a, V X X V)) {
75         a=substr(a,1,i+1) "-" substr(a,i+2, length(a)); n++;
76
77     while (i= match(a, V X X X V)) {
78         a=substr(a,1,i+1) "-" substr(a,i+2,length(a)); n++;
79
80     while (i=match(a, V X X X X V)) {
81         a=substr(a,1,i+1) "-" substr(a,i+2,length(a)); n++;
82
83     while (i=match(a, V X X X X X V)) {
84         a=substr(a,1,i+1) "-" substr(a,i+2,length(a)); n++;
85

```

```

86 # suppress the final schwa (^) in some multisyllabic words
87 # notr^ -> notR
88 # ar-bR^ => aRbR
89 b=gensub(/-([\^-]+)\^$/, "\1", 1, a) ;
90 if (b!=a) { # there is a schwa to delete
91     a=b;
92     $phons=substr($phons,1,length($phons)-1);
93     n-;
94 }
95 # meme chose quand schwa='*'
96 b=gensub(/-([\^-]+)\*$$/, "\1", 1, a) ;
97 if (b!=a) { # there is a schwa to delete
98     a=b;
99     $phons=substr($phons,1,length($phons)-1);
100    n-;
101    }
102
103
104 # compute the CVY skeleton
105 sk= " ";
106 for (i=1;i<=length(a);i++) {
107     ph=substr(a,i,1);
108     if (ph~V) sk=sk"V";
109     else if ((ph~C)|| (ph~L)) sk=sk"C";
110     else if (ph~Y) sk=sk"Y";
111     else sk=sk ph;
112 }
113 }
114
115 { print $0,a,n,sk }

```