

substring (graphemes), with any output string (phonemes). The rules are integrated into classes, a class being defined by the first character of the string to be transcribed. During the transcription process, the rules are examined in the order in which they are written and the first validated rule is then applied to the input text. Left and right contexts of the string to be transcribed may weight the rewriting process. A context is defined by:

- 1) a graphemic sub-string
- 2) a set defined in the declarations
- 3) the enumeration of 1) or 2) noted “,”
- 4) the concatenation of 1) or 2) noted “+”

The following illustration exemplifies the notion of conciseness mentioned in the preceding paragraph. These two rules:

(“#” + a) +s+ (“#”) = [s]

(the grapheme “s”, preceded by a word boundary + “a”, and followed by a word boundary is pronounced [s])

(“#” + bu) +s+ (“#”) = [s]

are merged into one:

(“#” + a,bu +s+ (“#”) = [s]

Then the enumeration of contexts is described by a lexicon, which is finally called on in the last rule:

(“#” + “lexica of those “s” pronounced at the end of the word”) +s+ (“#”) = [s]

It is apparent that the difference between rule and lexicon becomes fuzzy in this type of terminal rule because most of the mechanisms may be interpreted as the association (under contextual constraints) between graphic form and phonetic form, and therefore as lexica of associations, for example: the grapheme “o” is associated with [o] and the word “as” (“ace”) is associated with [as].

IV. CLASSIFICATION OF PROCEDURES AND SOLUTIONS

The process of phoneticization conveys linguistic information on the graphemic, lexical, morphophonemic, categorical, syntactic and semantic levels.

4.1 The graphemic level

The graphemic level corresponds to a regular rewriting of a graphemic string into a phonetic string. Thus it is a phonotactic level which ignores both the global form of the word and higher more complex levels. For example:

- the string “eau” (“water”) is transcribed as [o];
- if “s” is found between two vowels, it is rewritten [z]; otherwise it is rewritten [s].

4.2 The morphophonemic level

French lexical units are composed of at least one lexical base and a number of affixes. Thus a graphic word must be considered, as far as its phoneticization is concerned, as being made up of different types of morphemic units:

4.2.1 prefix + base

- “asocial” (“asocial”) is pronounced as “a” and “social”, “s” is transcribed [s] instead of [z];
- “polyacide” (“polyacide”) is pronounced “poly” + “acide”, “ya” is transcribed [ia] instead of [ja].

4.2.2 base + suffix

Here we meet the problem of the different pronunciations of the graphic forms “tie”, “tion”, where “t”, when not preceded by “s”, is transcribed either [s] or [t]. One may assume that “t” is pronounced [s] when “tie” and “tion” are suffixes. But when these suffixes are derivations, there subsists a problem of categorical ambiguity (see 4.4). For

example, “t” is pronounced [s] in “acrobatie” (“acrobatics”) and [t] in “sotie” (“farce”). When word lists taken from either Juilland’s dictionary [10] or the dictionary *Le 60 000* are analyzed, it is clear that this rule is not valid in all cases. To solve this kind of problem, we propose the systematic creation of specific lexica.

4.2.3 The construction base + base

This is a marginal construction in French, but one may find traces of it in, for example: “tournesol” (“sunflower”) which is pronounced as “tourne” (“turn”) and “sol” (“sun”, derived from “soleil”).

4.3 The lexical level

A set of words may share the same rewriting specificity. Very often this fact can be explained in that the set has undergone the same process under similar circumstances (loan words of like origin or similar use). These sets of irregular words will then be placed in the appropriate lexica. We will qualify these as exceptional words.

It is then possible to construct a complete list of lexica and to determine and verify their relation to linguistically identifiable classes.

We will present in section five the name and size of the more important lexica drawn from the ICP’s dictionary *Le 60 000*. The description of each is given in Belrhali and Libert [3]. Excluding some rare exceptions, the lexica contain loan or “learned” words. For example:

- “ch” is pronounced [k] in learned words (which are generally of Greek origin);
 - “er” is pronounced [e] in verbal inflection: “chanter” (“sing”) as well as at the end of polysyllabic nouns and adjectives: “boulangier” (“baker”), “léger” (“light”);
 - “er” is pronounced [ɛr] at the end of monosyllabic nouns and adjectives as in “cher” (“dear”) as well as at the end of loan words such as “bulldozer”, “bitter”, “bookmaker”.
- Exceptions are unusual and are brought about by historical accident or frequent use:
- “ch” is pronounced [ʃ] in “trachée” (“trachea”); current use has differentiated it from cognates where “ch” is pronounced [k]: “trachéotomie”;
 - “pt” is sometimes pronounced [pt] in “dompter” (“tame”), but the standard pronunciation given in the dictionary *Le Petit Robert 1* is [t]. The word’s pronunciation and typography has evolved diachronically: the initial form “dompter” (cf. Latin “domitare”) has given rise to the present form “dompter” by analogy to “compter” (“count”);
 - the dictionary *Le Petit Robert Oral-Ecrit* gives “dilemme” (“dilemma”) as [dilem], but also as [dilemn] by analogy to “indemne” (“unscathed”) [ɛdemn];
 - “carrousel” borrowed from Italian “carosello” in the 17th century has maintained the ‘s’ and is pronounced [karusel].

4.4 The categorical level

From the standpoint of phoneticization, several graphic markers (derivations and inflections) are ambiguous at the strictly phonotactic level. In particular, the endings “ent”, “tions” have different pronunciations following word category:

- “ent” is not pronounced when it constitutes the verbal inflection of the third person plural present: “président” (“preside”);
- “ent” is pronounced [ɑ̃] at the end of nouns, adjectives adverbs and a few verbal forms: “vent” (“wind”), “lent” (“slow”), “souvent” (“often”), “sent” (“feels”);
- “tions” is pronounced [tjɔ̃] when it is an inflectional marker of the third person plural of the imperfect tense:

“chantions” (“sing”), except in “balbutions” (“stammer”), “argutions” (“quibble”), “initions” (“initiate”) where it is pronounced [sjɔ̃];

• “tions” is pronounced [sjɔ̃] at the end of plural nouns, “rations” (“rations”).

If a linguistic description consists of reducing these ambiguities by categorizing the word under question, a practical solution would be to create lexica of “non-verbal” forms when the forms are not homographic (cf. section five). But when two linguistic units are homographic, the categorical ambiguity is multiplied by a lexical ambiguity that can not be reduced by the TOPH grammar:

- “se fier” (“trust”) [fje] verb vs. “fier” (“proud”) [fjɛʀ] adjective;
- “ent”: “président” (“preside”) [pʀezid] verb vs. “président” (“president”) [pʀezidã] noun;
- “tion”: “portions” (“carry”) [pɔʀtjɔ̃] vs. “portions” (“portions”) [pɔʀsjɔ̃] noun;
- inflectional ambiguity: “un os” (“a bone”) [œ̃n-ɔs] singular vs. “des os” (“bones”) [dez-ɔ] plural.

4.5 The semantic level

Finally, there exists a class of ambiguities which can not be resolved by the syntax due to categorical and lexical homomorphism. Only the semantic or *a fortiori* the pragmatic levels permit the reduction of this kind of ambiguity. Such is the case for the nouns:

- “les fils” (“the sons”) [fis] vs. “les fils” (“the threads”) [fil];
- “jet” (“gush”) [ʒɛ] vs. “jet” (“jet”) [ʒɛt].

V. THE FRENCH TOPH GRAMMAR

We have not included in the grammar the lexicon of words ending in “-ent” [ã] because it contains too many words (2,688) in comparison with other lexica. It will be replaced by the lexicon of verb forms in “-ent”. Having not yet been compiled from our dictionary *Le 60 000*, the precise number of words in this lexicon is not known, but it should be of much smaller size. The grammar now constitutes twenty lexica containing 1,650 words and a set of 1,270 rules.

Though it is difficult to decide whether a rule defines an exception or a regular rewriting, it is estimated that 600 of them describe about 200 exceptional cases and the other 670 characterize the basic rewriting process. Classified according to linguistic level, the following lists are examples of lexica classified according to the lexical, morphonemic, categorical and syntactic levels.

5.1 The lexical level

Lexica of words:

- beginning in “ai” pronounced [e] and not [ɛ], 25 words: “aider”, (“help”), “aigri” (“sour”)...
- ending in silent “-c”, 20 words: “banc” (“bench”), “blanc” (“white”)...
- ending in a pronounced “-d”, (loanwords), 42 words: “background”, “barmaid”...
- containing an “-e” pronounced [e], 52 words: “referendum”, “revolver”
- ending in “-et” pronounced [ɛt], 16 words: “cricket”, “gadjet”...
- ending in “-ey” pronounced [ɛ], 13 words: “jockey”, “bey” (“lord”)...
- containing a “-gu-” before “i” pronounced [g], 7 words: “anguille” (“eel”), “déguiser” (“disguise”)...
- ending in “-g” which is not pronounced, 24 words: “bast(a)ing” (“beam”), “bourg” (“town”)...

• beginning with an aspirated “h-”, 403 words:

“ha”, “hâbleur” (“boaster”)...

• ending in “-ing” pronounced [ɪŋ], 166 words:

“bowling”, “brainstorming”...

• ending in “-d” which is not pronounced, 22 words:

“cul” (“bottom”), “outil” (“tool”)...

• containing “-oo-” pronounced [u], 18 words:

“hollywood”, “lambswool”...

• ending in “-p” which is not pronounced, 15 words:

“beaucoup” (“many”), “camp” (“camp”)...

• beginning with “kw-” pronounced [kw], 53 words:

“quadragénaire”, “quadragésimal”...

• ending in “-s” pronounced [s], 510 words:

“bus”, “acinus”...

• ending in “-t” pronounced [t], 194 words:

“abject”, “abrupt”...

• ending in “-un-” pronounced [ɔ̃], 22 words:

“acupuncture”, “carborundum”...

• ending in “-x” which is not pronounced, 44 words:

“afflux”, “eaux” (“waters”)...

• beginning with “gn-” pronounced [ɲ], 8 words:

“gnon” (“punch”), “gniaf” (“cobbler”)...

• ending in “-er” pronounced [ɛʀ], 86 words:

“africamer”, “container”...

• ending in “-er” pronounced [œʀ], 33 words:

“bitter”, “bookmaker”...

5.2 The morphophonemic level

• Lexicon of prefixes ending in a vowel and which cause the following “s” to be transcribed as [s]:

“bi-, carbo-, entre-, homo-, ultra-”...

• Lexicon of prefixal bases (Greek stems) ending in a vowel and, as above, which cause a following “s” to be transcribed as [s]:

“bio-, éco-, hypo-, -, thio-, zoo-”...

• Lexicon of prefixes ending in “y” which cause the following “a” to be pronounced [ia]:

“poly-, oxy-, tachy-”...

• Lexicon of prefixes ending in “i” which cause a following “a” to be transcribed as [ia] and not [ija]:

“di-, bi-”...

5.3 The categorical level

• The lexicon of words ending in “-tions” pronounced [sjɔ̃] contains 1,780 nouns (i.e. “non-verbs”): “abdications”, “abductions”, “abréviations”... Given their number, we have therefore maintained the lexicon of verb endings with the same graphic form, 825 words:

verb forms ending in “-tions” (first person plural of the imperfect) pronounced [tjɔ̃], that is, all verbs in “-tions” with the exceptions of “argutions”, “balbutions”, “initions”;

verbs in “-ter”, “mentions” (“lie”), “portions” (“leave”)..., verbs in “-ter”, “abritions” (“shelter”), “absentions” (“absent”)...

• The lexicon of words ending in “-ent” pronounced [ã], that is all “non-verbs”, 2680 words:

“abaissément” (“lowering”), “abâtardissement” (“bastardizing”)...

5.4 The syntactic level

The lists are given simply as a supplement since a lexicon of these forms integrated into TOPH will not permit resolution of syntactic ambiguities:

• a lexicon of homographs ending in “-tions” pronounced [tjɔ̃] in verb forms and [sjɔ̃] for noun forms, “ablations”: (“ablations”) noun [-sjɔ̃] vs. (“we were ablating”) [-tjɔ̃] verb;

