



Language Manual

French

Pierre

Language Manual
French
Pierre
Edition 21
1 January 2005

This document was produced by Acapela Group. We welcome and consider all comments and suggestions. Please send them to:

Acapela Group
Box 1328
SE-171 26 Solna
Sweden

Phone +46 (0) 8 799 86 00
Fax + 46 (0) 8 799 86 01

Acapela Group
33, Boulevard Dolez
7000 Mons
Belgium

Tel: +32 (0)65 37 42 75
Fax: +32 (0)65 37 42 76

Acapela Group
3939, la Lauragaise
BP 58309
F-31683 Labège cedex
France

Tel: +33 (0)5 62 24 71 00
Fax: +33 (0)5 62 24 71 01

www.acapela-group.com

© Copyright Acapela Group 2005. All rights reserved.

List of contents

| | | |
|-------|---|----|
| 1 | General | 4 |
| 1.1 | Notational conventions | 4 |
| 2 | Letters in orthographic text | 5 |
| 2.1 | Characters with diacritics common in French orthography | 5 |
| 2.2 | Characters treated as letters in other languages | 6 |
| 3 | Non-alphanumeric characters | 7 |
| 3.1 | Punctuation characters | 7 |
| 3.1.1 | Comma, colon and semicolon | 7 |
| 3.1.2 | Quotation marks | 7 |
| 3.1.3 | Full stop | 7 |
| 3.1.4 | Question mark | 8 |
| 3.1.5 | Exclamation mark | 8 |
| 3.1.6 | Parentheses | 8 |
| 3.2 | Non-punctuation characters | 8 |
| 3.3 | Characters treated as special symbols | 9 |
| 3.4 | Characters whose pronunciation varies | 10 |
| 3.4.1 | Hyphen | 10 |
| 3.4.2 | Asterisk | 11 |
| 3.4.3 | Multiple occurrences of the same character | 11 |
| 3.5 | Control characters | 11 |
| 3.6 | Apostrophe | 11 |
| 3.7 | Characters ignored by the system | 11 |
| 4 | Number processing | 12 |
| 4.1 | Full number pronunciation | 12 |
| 4.2 | Leading zero | 12 |
| 4.3 | Decimal numbers | 13 |
| 4.4 | Monetary amounts | 13 |
| 4.5 | Ordinal numbers | 13 |
| 4.6 | Time of day | 14 |
| 4.7 | Arithmetic operators and other symbols | 14 |
| 4.8 | Mixed digits and letters | 14 |
| 5 | French Phonetic Text | 15 |
| 5.1 | Consonants | 16 |
| 5.2 | Vowels | 17 |
| 5.2.1 | Comments on phonetic symbols for vowels | 17 |
| 5.3 | Extra symbols for phonetic details | 18 |
| 5.3.1 | Stress | 18 |
| 5.3.2 | Punctuation marks | 18 |
| 5.3.3 | Hyphen | 18 |
| 6 | The RULSYS phonetic alphabet | 19 |
| 6.1 | RULSYS Consonants | 20 |
| 6.2 | RULSYS Vowels | 21 |
| 6.2.1 | Comments on phonetic symbols for vowels | 21 |
| 6.3 | Extra symbols for phonetic details | 22 |
| 6.3.1 | Stress | 22 |
| 6.3.2 | Punctuation marks | 22 |
| 6.3.3 | Hyphen | 22 |
| 7 | Liaisons | 23 |
| 7.1.1 | Liaisons in phonetic text | 23 |
| 8 | How to change pronunciation errors | 24 |
| 8.1 | Change the orthography | 24 |
| 8.1.1 | Spelling incorrectly | 24 |
| 8.1.2 | Use of hyphen | 24 |
| 8.1.3 | Expanding acronyms | 24 |
| 8.2 | Using phonetic text | 25 |
| 8.2.1 | Writing with phonetic text | 25 |
| 9 | Abbreviations | 26 |
| 9.1 | Entering abbreviations in the user lexicon | 26 |

1 General

This document discusses certain aspects of text-to-speech processing for the French text-to-speech system, in particular the different types of input characters and text that are allowed.

This version of the document corresponds to the High Density voice Pierre.

1.1 Notational conventions

The following notational conventions are used in this manual:

- For linguistic entities in general, **boldface** is used.
- Input text is written in a `non proportional font`.
- Output text is written in *italics*.
- Keyboard entities are written within angle brackets `< >`.
- Phonetic transcriptions are written within slashes or hash marks depending on the phonetic alphabet used.

The following abbreviations are used in this manual:

| | |
|----|---------------|
| LM | Letter mode |
| SM | Sentence mode |

See the User's Guide for a description of the two different reading modes. Note that Sentence mode is sometimes referred to as Normal mode.

2 Letters in orthographic text

Characters from A-Z and a-z may constitute a word. Characters that are used as letters in French and certain other European languages are also considered letters, see sections 2.1 and 2.2. The apostrophe character < ' > is not considered to be a letter, but may occur at the position of an elided vowel, in accordance with French orthographic conventions, e.g. *quelqu'un*, see section 3.6

Characters outside of these ranges, i.e. digits and non-alphanumeric characters such as punctuation characters and currency markers etc, are not considered as letters. If such a non-letter is included within a word, the word is ended where the non-letter appears and the following letters considered belonging to a new word.

2.1 Characters with diacritics common in French orthography

Characters written with the diacritics common in French orthography are mostly treated in a similar way. Insofar as they are used in French orthography they will be processed according to French pronunciation rules. Other letters containing the same diacritics will mostly be processed in a way that is reasonable from a French point of view. In LM the diacritic is read out as indicated in Table 1.

| Character | SM | LM |
|-----------|-----|-----------------------------|
| Ç | Ç | <i>(maj) C cédille</i> |
| ü | u | <i>u tréma</i> |
| é | é | <i>e accent aigu</i> |
| â | â | <i>a accent circonflexe</i> |
| ä | é/è | <i>a tréma</i> |
| à | à | <i>a accent grave</i> |
| ç | ç | <i>c cédille</i> |
| ê | ê | <i>e accent circonflexe</i> |
| ë | ë | <i>e tréma</i> |
| è | è | <i>e accent grave</i> |
| ï | ï | <i>i tréma</i> |
| î | î | <i>i accent circonflexe</i> |
| ì | i | <i>i accent grave</i> |
| Ä | É/È | <i>(maj) A tréma</i> |
| É | É | <i>(maj) E accent aigu</i> |
| ô | ô | <i>o accent circonflexe</i> |
| ö | œ | <i>o tréma</i> |
| ò | o | <i>o accent grave</i> |
| û | û | <i>u accent circonflexe</i> |
| ù | u | <i>u accent grave</i> |
| ÿ | y | <i>i grec tréma</i> |
| Ö | Œ | <i>(maj) o tréma</i> |
| Ü | u | <i>(maj) U tréma</i> |
| Œ | Œ | <i>(maj) E dans l'O</i> |
| œ | œ | <i>e dans l'o</i> |

Table 1 Non-French letters with diacritics

2.2 Characters treated as letters in other languages

The special letters of some other European languages are processed as indicated in Table 2.

| Character | SM | LM |
|-----------|------------|-------------------------|
| â | <i>a</i> | <i>a cercle</i> |
| Å | <i>A</i> | <i>Angstrœm</i> |
| Æ | <i>É/È</i> | <i>(maj) E dans l'A</i> |
| æ | <i>é/è</i> | <i>e dans l'a</i> |

Table 2 Non-French letters in the French system

3 Non-alphanumeric characters

The processing of non-alphanumeric characters varies, depending on the reading mode, context of the character, and its function within that context. These are the types of non-alphanumeric characters to be distinguished:

- Characters normally processed as punctuation, and having a direct effect on the intonation and pausing in SM.
- Non-punctuation characters that are always pronounced, with no effect on the intonation or pausing.
- Characters treated as special symbols
- Characters whose pronunciation varies according to context.
- Control characters

Below is a discussion of the characters grouped by type. For each character, the pronunciation is given in the three basic reading modes.

3.1 Punctuation characters

Table 3 lists punctuation characters permitted in the normal text input string and their readings. In SM they are all silent but they affect both rhythm and intonation as described in the sections below.

| Character | LM | SM |
|-----------|------------------------------|-----------|
| . | <i>point</i> | (silence) |
| ! | <i>point d'exclamation</i> | (silence) |
| ? | <i>point d'interrogation</i> | (silence) |
| , | <i>virgule</i> | (silence) |
| : | <i>deux points</i> | (silence) |
| ; | <i>point virgule</i> | (silence) |
| " | <i>guillemet</i> | (silence) |
| (| <i>parentèse gauche</i> | (silence) |
|) | <i>parentèse droite</i> | (silence) |
| « | <i>guillemet gauche</i> | (silence) |
| » | <i>guillemet droit</i> | (silence) |

Table 3 Punctuation characters

3.1.1 Comma, colon and semicolon

Comma < , >, colon < : > and semicolon < ; > cause a short pause to occur where the respective character is in the sentence. Comma is used as decimal mark in numbers, see section 4.3.

3.1.2 Quotation marks

Quotation mark < “ > and opening and closing quotation marks < > and < > may be used around a single word or a group of words in a sentence. In SM they cause a small pause to be produced where they appear in the text (unless they happen to be right at the beginning of a sentence, or followed immediately by another punctuation mark). In LM their names are read out (differently in the two modes, see Table 3).

3.1.3 Full stop

A full stop < . > is a sentence terminal punctuation mark which causes a falling end-of-sentence intonation pattern and is accompanied by a somewhat longer pause.

3.1.4 Question mark

A question mark <?> causes two different intonation patterns:

- For Yes/No type questions (closed questions, “questions fermées”), there is a rising intonation contour.
- For questions beginning with **qui, que, ou, quand, pourquoi, comment**, so-called WH-questions (open questions, “questions ouvertes”), there is a falling intonation pattern.

Both are accompanied by a pause.

3.1.5 Exclamation mark

The exclamation mark <!> has the same effect as the full stop, causing a falling intonation pattern followed by a pause. It is often helpful to use a word stress mark to signal the key word of the phrase, thus adding more expression to the sentence, see section 6.3.1.

3.1.6 Parentheses

Parentheses <()> appearing around a single word or a group of words cause a brief pause before and after the bracketed text.

3.2 Non-punctuation characters

The characters listed below in Table 4 are processed as non-letter, non-punctuation characters which are pronounced at all times in all reading modes.

| Character | LM/SM |
|-----------|-------------------------------|
| @ | <i>a commercial</i> |
| / | <i>barre oblique</i> |
| { | <i>accolade gauche</i> |
| } | <i>accolade droite</i> |
| \ | <i>barre oblique inversée</i> |
| | <i>barre verticale</i> |
| [| <i>crochet gauche</i> |
|] | <i>crochet droite</i> |
| \$ | <i>dollar</i> |
| & | <i>et</i> |
| ^ | <i>circonflexe</i> |
| ` | <i>apostrophe inversé</i> |
| < | <i>inférieur à</i> |
| > | <i>supérieur à</i> |
| % | <i>pourcent</i> |
| — | <i>souligné</i> |
| ~ | <i>tilde</i> |
| = | <i>égal (see 3.4.3)</i> |
| + | <i>plus (see 3.4.3)</i> |
| # | <i>numéro</i> |

Table 4 Non-punctuation characters

3.3 Characters treated as special symbols

The characters in Table 5 are not considered letters. Many of these are treated as special symbols and as such are processed by the system. Apart from the special cases noted below, each of the characters in Table 5 will be read as a separate word, whether or not they occur together with letters or digits.

| Character | SM | LM |
|--------------|---|--------------------------------------|
| ¢ | <i>centime</i> | <i>centime</i> |
| £ | <i>livre Sterling</i> | <i>livre Sterling</i> |
| ¥ | <i>yen</i> | <i>yen</i> |
| Pt | <i>pesetas</i> | <i>pesetas</i> |
| ^a | <i>petit a</i> | <i>petit a</i> |
| ° | <i>petit o</i> (see special cases below) | <i>petit o</i> |
| ¿ | (silence) | <i>point d'interrogation inversé</i> |
| ½ | <i>un demi</i> | <i>un demi</i> |
| ¼ | <i>un quart</i> | <i>un quart</i> |
| ¡ | (silence) | <i>point d'exclamation inversé</i> |
| « | (silence) / <i>ouvrez les guillemets</i> (see 3.1) | <i>guillemet gauche</i> |
| » | (silence) / <i>fermez les guillemets</i> (see 3.1) | <i>guillemet droit</i> |
| α | <i>alpha</i> | <i>alpha</i> |
| β | <i>bêta</i> | <i>bêta</i> |
| π | <i>pi</i> | <i>pi</i> |
| Σ | <i>maj sigma</i> | <i>maj sigma</i> |
| σ | <i>sigma</i> | <i>sigma</i> |
| Ω | <i>maj oméga</i> | <i>maj oméga</i> |
| δ | <i>delta</i> | <i>delta</i> |
| ∞ | <i>infini</i> | <i>infini</i> |
| ≡ | <i>identique à</i> | <i>identique à</i> |
| ± | <i>plus ou moins</i> | <i>plus ou moins</i> |
| ≥ | <i>supérieur ou égal</i> | <i>supérieur ou égal</i> |
| ≤ | <i>inférieur ou égal</i> | <i>inférieur ou égal</i> |
| ≈ | <i>peu différent d'eux</i> | <i>peu différent d'eux</i> |
| ° | <i>degré</i> (see special cases below) | <i>degré</i> |
| • | <i>gros point</i> | <i>gros point</i> |
| · | <i>point supérieur</i> | <i>point supérieur</i> |
| √ | <i>racine carrée</i> | <i>racine carrée</i> |
| ² | <i>au carrée</i> (see special cases below) | <i>au carrée</i> |

Table 5 Non-letter characters recognised by the French system

The symbols m^2 , km^2 , mm^2 , cm^2 are read as *mètre(s) carré(s)*, *kilomètre(s) carré(s)*, etc.

km , mm and cm are also read as abbreviations, see 9.1, but m on its own is read as the name of the letter; to have it read differently it is necessary to enter it in the user lexicon. $\langle 2 \rangle$ will be read as *au carrée* in other cases.

The character $\langle \circ \rangle$, which in some fonts is written with a line underneath is read as *petit o* except in the collocations N° and N^o s which are read as *numéro(s)*. The character $\langle \circ \rangle$ is read as *degré(s)*, but this is treated the same way as $\langle \circ \rangle$ in the collocations N° and N^o s which are also read as *numéro(s)*.

3.4 Characters whose pronunciation varies

The pronunciation of the characters listed below varies according to their context.

| Character | LM | SM |
|-----------|----------------------|-----------------------|
| - | <i>trait d'union</i> | (see 3.4.1 and 3.4.3) |
| * | <i>astérisque</i> | (see 3.4.2 and 3.4.3) |

Table 6 Characters with varying pronunciation

See section 3.4.3 regarding the reading of multiple occurrences of the characters < * + - = >.

3.4.1 Hyphen

The reading of hyphen <-> follows the following principles:

- If surrounded by digits, it is pronounced *tiret*.
- When used to link words together, it is not pronounced.
- It is discarded at the end of a line, and causes the two parts of the hyphenated word to be joined into a single word.
- If a word normally contains a hyphen, and the hyphen terminates the line, the hyphen will be discarded and the two parts of the word will be merged together. To avoid the words being merged, keep the hyphenated word together on the same line.
- If surrounded by spaces, it is pronounced as a short pause.

Expression

555-4758

Donnez-le-moi

A la fin d'une ligne,
on est obligé par-<CR>
fois de couper un mot en deux.

Comment allez-<CR>
vous, Madame?

Comment allez-vous,<CR>
Madame?

A propos - comment le saviez-vous?

Reading

555 *tiret* 4758

donnez-le-moi

*à la fin d'une ligne, on est obligé parfois
de couper un mot en deux*

Comment allezvous, Madame?

Comment allez vous, Madame?

*à propos (pause)
comment le saviez-vous?*

3.4.2 Asterisk

Asterisk < * > is pronounced *fois* in SM if a digit immediately precedes and follows it; it is pronounced *astérisque* in all other cases.

| Expression | Reading |
|------------|------------------------|
| 2*3 | <i>deux fois trois</i> |
| *bc | <i>astérisque B C</i> |

3.4.3 Multiple occurrences of the same character

In SM, if more than three of the same character occur in sequence without a space separating the characters, only the first three occurrences will be pronounced. This is only valid for the following characters: < * + - = >.

| Expression | Reading |
|------------|---|
| ***** | <i>astérisque astérisque astérisque</i> |
| +++++ | <i>plus plus plus</i> |
| ----- | <i>tiret tiret tiret</i> |
| ===== | <i>égal égal égal</i> |

3.5 Control characters

In LM (only) most control characters are read out. In most cases they are read as *contrôle* + the appropriate letter, e.g. ^T is read as *contrôle T*. However, some of the control characters are read in accordance with the function that they are most commonly given in text applications, as follows:

| Character | LM |
|-------------------|-----------------------------|
| ^H <BACKSPACE> | <i>caractère précédant</i> |
| ^I <TAB> | <i>tabulation</i> |
| ^J <LINE FEED> | <i>nouvelle ligne</i> |
| ^K <VERTICAL TAB> | <i>tabulation verticale</i> |
| ^M <RETURN> | <i>retour</i> |
| ^? <DELETE> | <i>effacement</i> |

Table 7 Control characters with special names read in letter mode

3.6 Apostrophe

An apostrophe, < ' >, is allowed in the text to mark elision of one vowel in front of another, for example, *j'aime*.

3.7 Characters ignored by the system

All characters that are not described in section 2 and 3 and that are not phonetic symbols or digits, are ignored by the system. Normally, these characters are omitted but some of them may cause the sentence they appear in to be silent.

4 Number processing

Strings of digits that are sent to the text-to-speech converter are processed in several different ways, depending on the reading mode, format of the digit string, and the immediately surrounding punctuation or non-numeric characters. To familiarise the user with the various types of formatted and non-formatted strings of digits that are recognized by the system, we provide below a brief description of the basic number processing along with examples.

Number processing is subdivided into the following categories:

- Full number pronunciation
- Leading zero
- Decimal numbers
- Monetary amounts
- Ordinal numbers
- Time of day
- Arithmetic operators and other symbols
- Mixed digits and letters

Note that the system does not handle date formats.

4.1 Full number pronunciation

When LM is enabled, digit strings are read as single digits, and all punctuation marks are read.

In SM, full number pronunciation is given for the whole number part of the digit string. Optional periods < . > may be used in the whole number portion to separate three-digit groups into hundreds, thousands, millions, and billions. Comma < , > may be used only to indicate a decimal amount, see section 4.3.

The highest number read is 9999999999 (twelve digits). Numbers higher than this are read as separate digits, with pauses between three-digit groups.

| Number | Reading |
|---------------|--|
| 2425 | <i>deux mille quatre cent vingt-cinq</i> |
| 2.425 | <i>deux mille quatre cent vingt-cinq</i> |
| 1.000.000.000 | <i>un milliard</i> |
| 1234567890123 | <i>un deux trois (pause) quatre cinq six (pause) sept huit neuf (pause) zéro un (pause) deux trois</i> |

4.2 Leading zero

Digit strings that begin with 0 (zero) are read digit by digit, with pauses between groups of digits if there are four or more digits.

| Number | Reading |
|--------|--|
| 09253 | <i>zéro neuf deux (pause) cinq trois</i> |
| 0210 | <i>zéro deux (pause) un zéro</i> |

4.3 Decimal numbers

In decimal amounts, the digits to the left of the decimal mark < , > observe the rules for full number pronunciation and leading zero, see section 4.1 and 4.2 respectively.

Digits occurring to the right of the decimal mark are read as full numbers or as single digits depending on the number of digits. If there are 2 or 3 digits following the decimal mark, and the first digit is not a zero, the decimal portion is read as a full number. In all other cases the decimal portion is read as single digits, with pauses occurring between groups of digits if there are four or more digits.

| Number | Reading |
|---------|--|
| 16,234 | <i>seize virgule deux cent trente-quatre</i> |
| 3,1415 | <i>trois virgule un quatre (pause) un cinq</i> |
| 1251,04 | <i>mille deux cent cinquante-et-un virgule zéro quatre</i> |
| 2,50 | <i>deux virgule cinquante</i> |
| 2.50 | <i>deux (pause) cinquante</i> |
| .65 | <i>soixante-cinq</i> |

4.4 Monetary amounts

Dollar amounts using the < \$ > sign are the only monetary amounts currently supported in this version of the French language. French, Belgian and Swiss franc amounts are not supported.

In a dollar amount, optional commas may be used in the dollar portion of the string to delimit thousands, hundred thousands, etc.

| Expression | Reading |
|------------|--|
| \$1988.45 | <i>mille neuf cent quatre-vingt-huit dollar et quarante-cinq cents</i> |
| \$1,988.45 | <i>mille neuf cent quatre-vingt-huit dollar et quarante-cinq cents</i> |
| \$200 | <i>deux cent dollars</i> |
| \$.45 | <i>quarante-cinq cents</i> |

4.5 Ordinal numbers

When the ordinalizer *ième* or *ieme* (or *ème* or *eme* or even *e*) is appended to a digit string, the ordinal version of the full number is spoken in SM only. **1er** and **1ere** or **1ere** are spoken as *premier* and *première* respectively.

| Expression | Reading |
|------------|-------------------------------|
| 21ième | <i>vingt et unième</i> |
| 21ieme | <i>vingt et unième</i> |
| 21ème | <i>vingt et unième</i> |
| 21eme | <i>vingt et unième</i> |
| 21e | <i>vingt et unième</i> |
| 42ième | <i>quarante-deuxième</i> |
| 123ieme | <i>cent vingt-troisième</i> |
| 95eme | <i>quatre-vingt-quinzième</i> |

4.6 Time of day

Time of day is read in SM if the following format is observed:

HhM

H represents the hours

M represents the minutes

The hours and minutes must be separated by the letter h. The minutes field is optional.

| Time | Reading |
|-------------|----------------------------------|
| 10h30 | <i>dix heures trente</i> |
| 1h15 | <i>une heure quinze</i> |
| 08h45 | <i>huit heures quarante-cinq</i> |
| 14h | <i>quatorze heures</i> |

4.7 Arithmetic operators and other symbols

Digit strings with arithmetic operators and miscellaneous symbols are processed according to the examples below.

| Expression | Reading |
|-------------------|---|
| 25% | <i>vingt-cinq pourcent</i> |
| 3,4% | <i>trois virgule quarante-cinq pourcent</i> |
| ,05% | <i>virgule zéro cinq pourcent</i> |
| -12 | <i>minus douze</i> |
| +24 | <i>plus vingt-quatre</i> |
| 2*3 | <i>deux fois trois</i> |

4.8 Mixed digits and letters

In SM, no alphabetic characters are allowed in a digit string. The occurrence of one or more letters in a digit string will cause the digit string to be terminated and output according to the number processing described above.

| Expression | Reading |
|-------------------|--|
| 208FR | <i>deux cent huit F R</i> |
| 77B84Z3 | <i>soixante-dix-sept B quatre-vingt-quatre Z trois</i> |
| 0092B87-B | <i>zéro zéro (pause) neuf deux B quatre-vingt-sept tiret B</i> |

5 French Phonetic Text

In the current version of the text-to-speech system, SAMPA (Speech Assessment Methods Phonetic Alphabet) is used when making lexicons or using phonetic strings within texts. In earlier versions, RULSYS was used. For the voices based on RULSYS, a conversion is made automatically from SAMPA to RULSYS inside the system.

We recommend new users to use only SAMPA since this is the notation that will be used in future development. Users who are already familiar with the RULSYS alphabet still have the possibility to use it when making user lexicons for all RULSYS-based voices (among them the French voice Pierre). There will be a description of RULSYS in the next section.

For the sake of clarity, SAMPA transcriptions are written within slashes (/ /) and RULSYS transcriptions within hash marks (# #). Note that neither the slashes nor the hash marks are part of the actual transcription.

The French system uses a phonetic alphabet similar to the French subset of SAMPA. The phonetic alphabet is described below.

If the pronunciation is incorrect the user may write phonetic transcriptions in the text. Then, a PRN-tag is needed to switch to phonetic mode, see User's Guide. It is also possible to make user lexicons (see User's Guide), or change the orthography of a word (see section 8) in order to achieve the preferred pronunciation.

5.1 Consonants

The table below lists the phonetic symbols in SAMPA used for the French consonants along with example words (the letters corresponding to the consonant sound are in boldface) and their transcriptions.

| Consonant symbol | Example | Transcription |
|------------------|----------------|---------------|
| b | baobab | /b a O b a b/ |
| d | dinde | /d e~ d/ |
| f | fife | /f i f/ |
| g | gag | /g a g/ |
| j | vieille | /v j e j/ |
| k | coq | /k O k/ |
| l | lille | /l i l/ |
| m | même | /m E m/ |
| n | nonne | /n O n/ |
| p | pipe | /p i p/ |
| R | rare | /R a R/ |
| s | sauce | /s o s/ |
| t | toute | /t u t/ |
| v | vive | /b i v/ |
| w | oui | /w i/ |
| z | oser | /o z e/ |
| S | cherche | /S E R S/ |
| Z | juge | /Z y Z/ |
| H | suis | /S H i/ |
| N | dancing | /d a~ s i N/ |
| n j | ligne | /l i n j/ |

Table 8 French consonant symbols

5.2 Vowels

The table below lists the phonetic symbols in SAMPA used for the French vowels along with example words and their transcriptions.

| Vowel symbol | Example | Transcription |
|--------------|---------------|---------------|
| a | bat | /b a/ |
| _A1 a | bas | /b _A1 a/ |
| a~ | banc | /b a~/ |
| e | bée | /b e/ |
| E | baie | /b E/ |
| e~ | bain | /b e~/ |
| @ | brebis | /b R @ b i/ |
| i | bis | /b i s/ |
| O | bol | /b O l/ |
| o | beau | /b o/ |
| o~ | bon | /b o~/ |
| u | bout | /b u/ |
| y | bu | /b y/ |
| 9 | beure | /b 9 r/ |
| 2 | boeufs | /b 2/ |
| 9~ | un | /b 9~/ |

Table 9 French vowel symbols

Note that /_A1 a / is not a proper SAMPA symbol but it may still be used when making transcriptions. There is a corresponding symbol in RULSYS, see next section.

5.2.1 Comments on phonetic symbols for vowels

The vowel sound referred to as “e-muet” or “e-caduc” is represented by the phonetic symbol /@/. It is found in words like **le**, **prenons**, and **petit**. The symbol /@/ should only be used when transcribing orthographic **e**:

petit /p @ t i/
monsieur /m @ s j 2/

The phonetic symbol /9/ should be used to represent the full vowel represented by the letters **eu** in normal orthography:

peu /p 9/

The French rules take care of @-deletion/retention within a word, at the end of a word, and across word boundaries:

samedi /s a m d i/
barque /b a R k/
une amie /y n _spc a m i/

@ is retained if more than two consonants would occur together due to @ deletion:

table de travail /t a b l @ _spc d @ _spc t R a v a j/

If @-deletion is undesirable in a particular word, transcribing the word phonetically usually avoids the deletion. For example, **devenir** is normally pronounced *devnir* by the French system. To force the

second e to be pronounced, the word can be phonetically transcribed as /d @ v @ n i R/ and all three syllables will be pronounced.

5.3 Extra symbols for phonetic details

In the current version of the French synthesis certain phonetic details can be specified in phonetic text. This can be exploited in case the user wishes to achieve an unusual pronunciation, or if the transcription automatically generated by the system is inaccurate.

5.3.1 Stress

In words with more than one syllable, one (and normally only one) of the syllables is more prominent than the others. This is referred to as word stress, or lexical stress. Words of one syllable also have word stress when spoken in isolation, although many may lose the stress in certain contexts.

In French, the word-level stress normally occurs on the last full vowel of a word in isolation, e.g. **petite**. Not all words have stress. The so-called *function* words (prepositions, pronouns, articles, etc.) are usually unstressed. When a word is combined with other words in a phrase, the stress that occurs when the word is in isolation may be lost in favour of stressing the last word of a phrase, e.g. **la petite fille**.

A word can also be stressed in a sentence to emphasise its importance, e.g. **La petite fille, je veux dire, pas la grande**. The text-to-speech system automatically determines the stress pattern of orthographically entered text. To alter the stress pattern produced by the system one can indicate phrase-level stress. In orthographic text, this is done by placing <_X>, where X represents a single digit between 0 and 9, within a PRN-tag (see User's Guide) immediately before the word whose prominence is to be altered. The emphasis mark can also be used in transcriptions in a user lexicon.

| | |
|---------|--|
| _2 | normal stress for most words |
| _0 | makes a word non-stressed |
| _1 | gives stress to a normally unstressed word |
| _3 - _9 | gives levels of emphatic stress |

Compare how the meaning is changed when the emphatic stress is varied in the sentence below:

Examples C'est la fille de nos voisins.
C'est la /_4/ fille de nos voisins.
C'est la fille de /_6/ nos voisins.

Note that the necessary PRN-tag is not included in the examples above.

When using RULSYS notation, it is also possible to alter the word level stress, see section 6.3.1, this is not possible when using SAMPA.

5.3.2 Punctuation marks

The punctuation marks <. ! ? , > used in phonetic text have the same effect on intonation as when appearing in orthographic text. In SAMPA the punctuation marks are denoted /_/, /_! _/, /_? _/, and /_com _/ respectively.

5.3.3 Hyphen

In phonetic text, hyphen (in SAMPA underscore + hyphen, <_->) can be used to separate parts of a compound word. If the hyphen separating two parts of a word comes at the end of a line, the word is not read until the second part on the next line is typed. For a description of the use of the hyphen character in normal orthographic text, see section 3.3.1.

6 The RULSYS phonetic alphabet

Note that we recommend new users to use only SAMPA since this is the notation that will be used in future development. Note also that it is only possible to use RULSYS when making user lexicons, not in the input text string.

The following differentiates RULSYS from SAMPA in the French system:

- no spaces are used within words in transcriptions
- it is possible to denote word stress

Note that the hash marks (# #) are used to indicate RULSYS transcriptions and to differentiate them from SAMPA transcriptions; the hash marks are not part of the actual transcriptions.

If the pronunciation is incorrect the user may write phonetic transcriptions in the text. Then, a PRN-tag is needed to switch to phonetic mode, see User's Guide. It is also possible to make user lexicons (see User's Guide), or change the orthography of a word (see section 8) in order to achieve the preferred pronunciation.

6.1 RULSYS Consonants

The table below lists the phonetic symbols in RULSYS used for the French consonants along with example words and their transcriptions.

| Consonant symbol | Example | Transcription |
|------------------|----------------|---------------|
| B | baobab | #BAOBAB# |
| D | dinde | #DE9D# |
| F | fife | #FIF# |
| G | gag | #GAG# |
| J | vieille | #VJEJ# |
| K | coq | #KOK# |
| L | lille | #LIL# |
| M | même | #ME1M# |
| N | nonne | #NON# |
| P | pipe | #PIP# |
| R | rare | #RAR# |
| S | sauce | #SO1S# |
| T | toute | #TUT# |
| V | vive | #VIV# |
| W | oui | #WI# |
| Z | oser | #O1ZE# |
| SH | cherche | #SHE1RSH# |
| ZH | juge | #ZHYZH# |
| J1 | suis | #SJ1I# |
| NG | dancing | #DA9SING# |
| NJ | ligne | #LINJ# |

Table 10 RULSYS consonants

6.2 RULSYS Vowels

The table below lists the phonetic symbols in RULSYS used for the French vowels along with example words and their transcriptions.

| Vowel symbol | Example | Transcription |
|--------------|---------------|---------------|
| A | bat | #BA# |
| A1 | bas | #BA1# |
| A9 | banc | #BA9# |
| E | bée | #BE# |
| E1 | baie | #BE1# |
| E9 | bain | #BE9# |
| E0 | brebis | #BRE0BI# |
| I | bis | #BIS# |
| O | bol | #BOL# |
| O1 | beau | #BO1# |
| O9 | bon | #BO9# |
| U | bout | #BU# |
| Y | bu | #BY# |
| \ | beure | #B\R |
| \1 | boeufs | #B\1# |
| \9 | un | #\9# |

Table 11 RULSYS vowels

6.2.1 Comments on phonetic symbols for vowels

The vowel sound referred to as “e-muet” or “e-caduc” is represented by the phonetic symbol #E0#. It is found in words like **le**, **pre**nons, and **pe**tit. The symbol #E0# should only be used when transcribing orthographic e:

| | |
|-----------------|-----------|
| petit | #PE0TI# |
| monsieur | #ME0SJ\1# |

The phonetic symbol #\1# should be used to represent the full vowel represented by the letters **eu** in normal orthography:

| | |
|------------|-------|
| peu | #P\1# |
|------------|-------|

The French rules take care of E0-deletion/retention within a word, at the end of a word, and across word boundaries:

| | |
|-----------------|----------|
| samedi | #SAMDI# |
| barque | #BARK# |
| une amie | #YN AMI# |

#E0# is retained if more than two consonants would occur together due to E0-deletion:

| | |
|-------------------------|---------------------|
| table de travail | #TABLE0 DE0 TRAVAJ# |
|-------------------------|---------------------|

If E0-deletion is undesirable in a particular word, transcribing the word phonetically usually avoids the deletion. For example, **devenir** is normally pronounced *devnir* by the French system. To force the second e to be pronounced, the word can be phonetically transcribed as #DE0VE0NIR#, and all three syllables will be pronounced.

6.3 Extra symbols for phonetic details

In the current version of the French synthesis certain phonetic details can be specified in phonetic text. This can be exploited in case the user wishes to achieve an unusual pronunciation, or if the transcription automatically generated by the system is inaccurate.

6.3.1 Stress

For a description of stress, see section 5.3.1. In the current version of the French system, it is possible in RULSYS, but not in SAMPA, to denote word stress. This is done by placing an apostrophe < ' > before the vowel that is to receive the stress. However, this is not compulsory, as the system will stress this vowel automatically.

| | | |
|---------|-------------------|----------------|
| Example | ami | #AM' I# |
| | impossible | #E9POS' IBLE0# |

Remember that only vowels are stressed, i.e., a stress mark must be followed by a vowel written in phonetic characters. Be sure, for example, that you do not leave a real apostrophe in phonetic text.

It is also possible to denote phrase level stress, see section 5.3.1.

6.3.2 Punctuation marks

When using RULSYS, punctuation marks are also permitted in phonetic text, and have the same effect as in normal text, affecting both the rhythm and intonation of the sentence. These punctuation characters are permitted in phonetic text:

, . ? ! -

The character < ' > has a completely different function when writing in phonetic text than in ordinary text. It is a reserved character that may be used to mark stress in a word, see section 6.3.1. It cannot be used to quote text or mark the place of an elided vowel in phonetic text.

6.3.3 Hyphen

Hyphen, < - >, in phonetic text can be used to separate parts of a compound word.

| | | |
|---------|------------------------|-------------------|
| Example | rez-de-chaussee | #RE-DE0-SHO1S' E# |
|---------|------------------------|-------------------|

If the hyphen separating two parts of a word comes at the end of a line, the word is not spoken until the second part on the next line is also read in. A word written in phonetic text that contains one or more hyphens is spoken as a complete word in SM. For a description of the use of the hyphen character in normal orthographic text, see section 3.4.1.

7 Liaisons

In spoken French, a “liaison” links the final consonant of a word not normally pronounced to the following word if this begins with a vowel. For example, in the phrase **le petit chat** the **t** of **petit** is not pronounced since **chat** begins with a consonant, while in **son petit ami** a liaison between **petit** and **ami** is made and the **t** is pronounced.

In the French text-to-speech system, liaisons are generated automatically. Sometimes the system makes an undesirable liaison between two words. If this happens, the liaison can be “blocked” by inserting a special phonetic symbol between the two words where the liaison occurs. To block a liaison, insert the symbol /_l/ within an PRN-tag (see User’s Guide) between the two words.

It is also possible to create a liaison when the system does not automatically generate one. To force a liaison between two words, use an ordinary hyphen < - >.

| Input | Reading |
|-------------------------|--------------------------|
| Quand /_l/ est-il venu? | <i>Quan est-il venu?</i> |
| tout-à l'heure | <i>toute à l'heure</i> |

Note that the PRN-tag is replaced by slashes < / / > in the example above.

7.1.1 Liaisons in phonetic text

When writing in phonetic text, a liaison can be blocked by simply omitting the linking character from the phonetic text string.

Example **Quand est-il venu?** can be phonetically transcribed in RULSYS as #KA9 ETIL VEONY?#

A liaison can be created by including the linking consonant in the phonetic text string.

Example **Pied a terre** can be transcribed as #PJET A TE1R#

8 How to change pronunciation errors

Words that are not pronounced correctly by the text-to-speech converter can be entered in the user lexicon (see User's guide). There are two ways to do this: either, the user enters a phonetic transcription of the word (see section 5), or, the user rewrites the word orthographically. Phonetic transcriptions can also be entered directly in the text, using a PRN-tag (see User's guide).

8.1 Change the orthography

8.1.1 Spelling incorrectly

It is possible to intentionally misspell a word by trying to spell a word in a more phonetic manner, i.e., choosing non-ambiguous letter combinations to represent difficult sounds. For example, the letters **ch** in **archaque** might better be represented by the letter **k**, which is closer to the actual sound in the word.

Examples **computer** can be misspelled **compiouteur**
 Andrew can be misspelled **Androu**

Adding an **e** to the end of a word usually causes a word-final consonant to be pronounced. If the final consonant is **m** or **n**, the addition of a final **e** will change the preceding vowel from a nasal vowel to an oral vowel.

Example The noun **est** denoting the direction (east) can be misspelled **este** to avoid confusion with the verb **etre** (the word **est** will sometimes be correctly pronounced as the direction word, for instance in the collocation **a l'est de**).

PROM can be misspelled **promme** to rhyme with **pomme**.

8.1.2 Use of hyphen

A hyphen character can be used within a word to separate two letters that might otherwise be incorrectly pronounced together.

Example **deshonneur** can be written **des-honneur**

8.1.3 Expanding acronyms

Very few acronyms are handled by the current French system (see section 9). Therefore, it may be very useful to expand them in the user dictionary. Since acronyms should be expanded to more than one word it may be difficult to enter a proper transcription. It is much easier to enter the words in question orthographically. The examples below show some acronyms and their expanded readings.

Examples **UE** Union europeenne

8.2 Using phonetic text

When you are unable to correct a pronunciation error by misspelling the word, phonetic text should be used to produce the desired pronunciation. When phonetic text is used, the system bypasses the normal pronunciation rules, and pronounces each phonetic symbol “literally”, according to the examples listed in Table 8 and Table 9.

8.2.1 Writing with phonetic text

A helpful way to transcribe in phonetic text is to work with a dictionary. Often, dictionaries give a pronunciation guide for each word. They also provide a pronunciation key to show how to pronounce the special symbols used in the pronunciation guide. Similarly, Table 8 and Table 9 in this document give the pronunciation key for the special phonetic symbols used in French for the text-to-speech converter.

To transcribe a word phonetically, “sound out” the word slowly. Working with Table 8 and Table 9, find the symbols that most closely correspond to each of the sounds in the word you want to transcribe. To mark stress, decide which vowel is the most prominent, and place a primary stress mark < 1 > after the vowel.

Example

Suppose we want to transcribe the word **schooner**. We consult a dictionary and find an entry that looks like this:

schooner [skuno:r] *n.m.* Petit navire a deux mats.

Among the symbols listed in the dictionary’s pronunciation key are:

s (salle), k (cou), u (tout), n (nous), o: (soeur), r (rue)

Using Tables 8 and 9, we find the following corresponding symbols:

s (sauce), k (coq), u (bout), n (nonne), 9 (beurre), R (rare)

Now we can transcribe **schooner** using the French phonetic symbols for the text-to-speech converter.

/s k u n 9 R/

9 Abbreviations

In the current version of the French text-to-speech system, the abbreviations in Table 10 are recognised in all contexts in SM only. These abbreviations are case-insensitive, and do not require a period in order to be processed as an abbreviation. In SM, if a period accompanies the abbreviation, the sentence is terminated at the abbreviation and output.

The user lexicon may be used to redefine any of the abbreviations below, or to create your own.

| Abbreviation | LM | SM |
|--------------|----------------|----------------------|
| dr | <i>DR</i> | <i>docteur</i> |
| mlle | <i>M L L E</i> | <i>mademoiselle</i> |
| mme | <i>M M E</i> | <i>madame</i> |
| st | <i>ST</i> | <i>saint</i> |
| ste | <i>ST E</i> | <i>sainte</i> |
| mm | <i>MM</i> | <i>millimètre(s)</i> |
| cm | <i>CM</i> | <i>centimètre(s)</i> |
| km | <i>KM</i> | <i>kilomètre(s)</i> |
| ml | <i>ML</i> | <i>millilitre(s)</i> |
| cl | <i>CL</i> | <i>centilitre(s)</i> |

Table 12 Abbreviations in the French system

As mentioned in section 4.6, the abbreviation **h** represents **heures** when placed immediately after a digit:

| Input | Reading (SM) |
|---------|----------------------------|
| 20h30 . | <i>Vingt heures trente</i> |
| 20 h . | <i>Vingt H</i> |

9.1 Entering abbreviations in the user lexicon

The single letters **m**, **l**, **s**, etc. are read as the names of the letters unless they are given a special transcription in the user lexicon. In LM they will always be read out as the names of the letters, even if they have been given a different transcription in the user lexicon. Thus, if one wishes the letter **M** on its own to be read as **Monsieur** or as **metre/metres** it is possible to enter suitable transcriptions in the user lexicon. Likewise one can give any other single letter, or combination of letters, a user-defined pronunciation by entering them in the user lexicon.

In order to avoid the transcription given to a single consonant letter in the user lexicon being applied also in front of an apostrophe, it is necessary to use a phonetic transcription and to start this with an equals sign **< = >**. When an apostrophe follows a word which has been given a transcription that begins with **< = >** followed by a consonant, this consonant will be read and the rest of the transcription discarded. Transcriptions are given below in both RULSYS and SAMPA notation.

| User lexicon entry | Input | Reading (SM) |
|-------------------------------|----------------------|-----------------------------------|
| M #ME0SJ\1# /m @ s j 2/ | Je m'appelle Pierre. | <i>Je monsieur appelle Pierre</i> |
| M #ME0SJ\1# /_= m @ s j 2/ | Je m'appelle Pierre. | <i>Je m'appelle Pierre .</i> |
| S #SE9# /s e~/ | Il s'appelle Pierre. | <i>Il saint appelle Pierre .</i> |
| S #SE9# /_= s e~/ | Il s'appelle Pierre. | <i>Il s'appelle Pierre</i> |

Since the letter **M** is often used as an abbreviation for both **Monsieur** and **metre(s)**, there are special rules for these readings. The system contains rules for dealing with the following entries in a user lexicon as readings for the single letter **m**:

| User lexicon entry | Reading |
|--|--|
| M #===ME1TR+E0# /_= _= _= m E t R _+ @/ | <i>mètre(s)</i> only when a number precedes <i>Monsieur</i> otherwise <i>m</i> before apostrophe |
| M #==ME1TR+E0# /_= _= m E t R _+ @/ | <i>mètre(s)</i> only when a number precedes <i>M</i> otherwise but <i>m</i> before apostrophe |
| M #=ME1TR+E0# /_= m E t R _+ @/ | <i>mètre(s)</i> in all cases <i>m</i> before apostrophe |
| M #ME1TR+E0# /_= m E t R _+ @/ | m on its own will be read as <i>mètre</i> in all cases, even when followed by apostrophe |
| M #=ME0SJ\1# /_= m @ s j 2/ | <i>Monsieur</i> in all cases except before apostrophe |

The two-letter word **MM** will by default be read as *millimètre(s)* by the system. To have this read as *Messieurs* it is necessary to enter it in the user lexicon. As for the single-letter abbreviation **M** there are special formats allowing different pronunciations of this abbreviation.

| User lexicon entry | Reading |
|---|--|
| MM #===MILIME1TR+E0# /_= _= _= m i l i m E t R _+ @/ | <i>millimètre(s)</i> only when a number precedes <i>Messieurs</i> otherwise but <i>m</i> before apostrophe |
| MM #==MILIME1TR+E0# /_= _= m i l i m E t R _+ @/ | <i>millimètre(s)</i> only when a number precedes <i>MM</i> otherwise <i>m</i> before apostrophe |
| MM #=MILIME1TR+E0# /_= m i l i m E t R _+ @/ | <i>millimètre(s)</i> in all cases but <i>m</i> before apostrophe |

For **L** the following user lexicon entries can be used:

| User lexicon entry | Reading |
|---------------------------------------|---|
| L #==LITR+E0# /_= _= l i t R _+ 2/ | <i>litre(s)</i> in all cases other than before an apostrophe |
| L #LITR+E0# /l i t R _+ 2/ | <i>litre(s)</i> in all cases including before an apostrophe |

The text-to-speech system does not normally distinguish between small (lower-case) letters and capital (upper-case) letters. However, a user lexicon can be made case-sensitive (in some systems, the user lexicon is always case-sensitive; see the User's Guide for further information). This can be exploited in order to distinguish between for instance **M** as *Monsieur* and **m** as *mètre(s)*. A useful procedure would be to have a file containing the following items (see next page):

```

M      #===ME1TR+E0#
      /_ = _ = _ = m E t R _+ @/

MM     #===MILIME1TR+E0#
      /_ = _ = _ = m i l i m E t R _+

Mm     #===MILIME1TR+E0#
      /_ = _ = _ = m i l i m E t R _+

m      #===ME1TR+E0#
      /_ = _ = _ = m E t R _+ @/

```

Then the desired items can be copied from this file into the user lexicon. Note: use either the RULSYS transcription (indicated by # #) or the SAMPA transcription (indicated by / /).