

# Conventions de transcription

## LI – Université de Tours

Jean-Yves Antoine

LI – Université François Rabelais de Tours

Université François Rabelais Tours

[http://www.info.univ-tours.fr/~antoine/parole\\_publicue/](http://www.info.univ-tours.fr/~antoine/parole_publicue/)

## ANNEXE A — Conventions de transcription du corpus Accueil\_UBS

---

La transcription est strictement orthographique, avec mention minimale des événements acoustiques connexes (voir ci-après). D'une manière générale, les conventions de transcription s'inspirent des recommandations utilisées dans le projet SPEECHDAT (Gibbon *et al.*, 1997), ainsi que des conventions définies par la laboratoire DELIC pour le français.

### 1.1 Structuration de la transcription : tours de parole

Chaque dialogue est segmenté en tours de parole. La définition du tour de parole varie dans la littérature d'un auteur à l'autre. Dans le cadre de ce corpus, nous avons utilisé la définition opérative suivante : un nouveau de parole apparaît lorsqu'un nouveau locuteur se met à parler. Deux situations peuvent alors survenir :

**Tour de parole sans chevauchement** — Le tour de parole est délimité par (début) la prise de parole d'un locuteur et (fin) par la fin de sa production. Ce tour de parole ne concerne donc qu'un seul locuteur. Exemple de tour de parole sans chevauchement transcrit au format ASCII :

```
<03> institutrice  
i: quel film veux tu voir
```

**Tour de parole avec chevauchement** — Le tour de parole est délimité par le début et la fin du chevauchement. Ce tour de parole regroupe alors deux (voire plus) locuteurs. Leurs productions orales sont représentées simultanément dans ce tour de parole, en distinguant chaque locuteur. Exemple de tour de parole avec chevauchement transcrit au format ASCII :

```
<04> client + hôtesse  
c: d'accord  
h : on a simplement
```

Dans les dialogues, les périodes sans chevauchement succèdent bien entendu sans arrêt à des périodes avec chevauchement.

A titre d'exemple, supposons qu'un locuteur prononce un certains énoncé (par exemple « Tiens j'ai vu Paul hier ») tandis que le second locuteur se contente d'une marque d'étonnement (« ah ouais ») en milieu d'énoncé. Cette « tranche » de dialogue sera alors segmentée en 3 tours de parole :

- début d'énoncé sans chevauchement du locuteur 1,
- partie chevauchée avec prononciations des locuteurs 1 et 2,
- fin d'énoncé sans chevauchement du locuteur 2.

### 1.2 Conventions de transcription

La transcription est strictement orthographique, avec mention minimale des événements acoustiques connexes (voir ci-après). Elle suit les normes orthographiques standards du français. Notons cependant que tout mot sera séparé par un espace (blanc), le tiret entre deux mots n'étant conservé que si ceux-ci constituent un lemme insécable. Ainsi :

<i>puis-je</i>	sera transcrit	puis je	(2 mots)
<i>plate-forme</i>	sera transcrit	plate-forme	(1 mot)

La description des événements acoustiques ou prosodiques est limitée au minimum et est non exhaustive.

On se contente ainsi de marquer seulement les pauses longues, sans distinction de type. De même, la transcription ne comprendra aucune marque de ponctuation<sup>1</sup>.

---

<sup>1</sup> Les linguistes travaillant sur l'oral, tels les chercheurs du GARS/DELIC, dénie généralement toute pertinence de la notion de ponctuation dans le langage parlé.

### 1.2.1 Bruits

Ce corpus a été enregistré en conditions réelles avec un médiocre rapport signal sur bruit. Les bruits non humains n'ont pas été transcrits. Nous avons par contre opéré réalisé une annotation minimale de certains bruits de l'appareil phonatoire :

<i>rire</i>	annoté	[rire]
<i>bruits de bouche</i>	annoté	[bb]
<i>toux</i>	annoté	[tx]
<i>souffle</i>	annoté	[pf]

### 1.2.2 Majuscules / minuscules

De manière générale, les transcriptions ne comportent que des caractères minuscules. L'emploi de majuscules est néanmoins pertinent pour marquer les noms propres de la langue ainsi que les caractères épelés. D'une manière plus précise :

- les énoncés transcrits ne débutent pas par une majuscule (on retrouve ici l'absence de ponctuations),
- Les acronymes et les caractères épelés (ou sigles) sont transcrits en majuscule. Ils ne sont pas séparés par des points :

*S N C F* et non *S.N.C.F.*

- les noms propres commencent par une majuscule (par exemple : *Jospin*, *Grenoble*). L'application de cette règle est stricte afin d'éviter d'englober autant que possible des noms communs. Ainsi, on transcrit :

*monsieur Lionel Jospin* et non *Monsieur Lionel Jospin*  
*mairie de Grenoble* et non *Mairie de Grenoble*

A l'opposé, les noms propres correspondant à des sigles sont mentionnés à l'aide de majuscules. L'existence d'un acronyme correspondant à ce sigle est un bon indice de "capitalisation". Par exemple :

*Société Nationale des Chemins de Fer* (SNCF)  
*Transports de l'Agglomération Grenobloise* (TAG)

- les noms communs ayant fonction de nom propre (par exemple : titre de film) ne correspondant pas à un sigle sont transcrits entre guillemet et restent en minuscule. Lorsqu'on relève un nom propre dans ce type de nom commun, il prend bien entendu une majuscule. Par exemple :

*le bureau "info montagne"*  
*"l'amicale laïque de la ville de Massy"*

**Remarque** — Cette règle de transcription était optionnelle, la délimitation des situations sigle / nom commun ayant fonction de nom propre / nom commun étant relativement floue.

### 1.2.3 Nombres

A l'exception du nombre *un* qui peut être confondu avec l'article indéfini, les nombres ont été codés en chiffre lorsque leur prononciation suivait celle du français standard. Par exemple :

*128* et non *cent vingt huit*

Dans le cas contraire, les nombres ou séquences de nombres sont transcrites en caractères afin de refléter la prononciation exacte du locuteur. Par exemple :

*septante deux* et non *72*

### 1.2.4 Acronymes et sigles

La transcription des sigles, déjà évoquée, suit bien entendu la prononciation du locuteur :

- Intégralement s'il est prononcé mot à mot : *Société Nationale des Chemins de Fer*

- Sous forme de caractères épelés si son acronyme est prononcé lettre à lettre : S N C F
- Sous forme d'un nom propre particulier si son acronyme n'est pas épelé : Tag et non T A G

### 1.2.5 Prononciations incomplètes

Sont considérées ici les prononciations incomplètes de mots dues au caractère spontané de la parole : phénomènes de reprises ou répétitions, ou interruptions par l'autre locuteur. Elles seront marquées à l'aide des parenthèses placées en fin du fragment prononcé. Ce fragment sera transcrit sous forme orthographique en suivant les règles standard de prononciation. Lorsqu'il y a difficulté d'interprétation du fragment, la transcription complète du mot attendu est précisée entre les parenthèses. Par exemple :

*donne moi une po()* *une poire* ou encore  
*donne moi une po(pomme) une poire*

### 1.2.6 Délétions, contractions

Le français parlé présente de nombreuses occurrences de contractions ou de délétions de syllabes qui concernent en particulier les locutions fréquentes ou les petits mots outils. Ces délétions ne peuvent être considérées comme des prononciations incomplètes, puisqu'elles relèvent de la stratégie d'élocution et non du caractère spontané de la production.

Certaines transcription rivalisent de conventions particulières destinées à rendre compte le plus précisément possible de la prononciation réalisée (par exemple : *y' a ka* pour *il n'y a qu'à*). Au contraire, on s'est limité ici — à l'instar des recommandations du DELIC (ex-GARS) — à une transcription aussi proche que possible de l'écriture standard. Par exemple :

*je vais* pour *j'veis* (en phonétique : /jve/)  
*il y a* pour *y'a*

Dans le cas d'une délétion complète de mot (cas de la chute du discordantiel *ne*, par exemple), le mot ne sera pas transcrit.

### 1.2.7 Erreurs de prononciations, prononciations idiomatiques

Les formes correspondant à une erreur manifeste de prononciation (lapsus, par exemple), ou à une prononciation idiomatique, sont transcrites sous leur forme régulière, précédée d'un astérisque. La forme réellement prononcée est alors transcrite sous forme orthographique, en respectant les règles standard de prononciation du français, entre crochets après la forme corrigée. Exemple :

*je \*rêpète{récapépète} depuis le \*début{béduť}*

Si la forme inattendue ne peut se traduire fidèlement sous forme orthographique, on adopte la notation phonétique ajoutée en signes "/". On utilise pour cela la convention de notation SAMPA.

### 1.2.8 Événements acoustiques : pauses

Deux types de pause ont été distinguées :

- pauses remplies (hésitations du type *euuh*, *mmh* etc...) notées par le sigle e
- pauses silencieuses notées par le sigle #

## 2 ANNEXE B — Codage : formats de transcription en sortie

---

Trois formats de sortie ont été définis pour les fichiers de transcription

- codage XML,
- codage en format texte (ASCII),
- format PDF regroupant dans un seul fichier l'ensemble des transcriptions obtenues en format texte.

### 2.1 Codage XML

La transcription a été réalisée à l'aide du logiciel libre Transcriber. Le format XML de sortie suit donc la DTD définie par ce logiciel. Nous ne détaillerons pas ici cette DTD : le lecteur intéressé se référera à (Barras *et al.* 1998) ou consultera le site Internet consacré à Transcriber :

<http://trans.sourceforge.net/>

On notera simplement que ce format de sortie permet de décrire les chevauchements ainsi que l'alignement temporel des débuts et fin de tours de parole. La version de Transcriber utilisée (version Windows) présentait un bug de codage du « à » en Unicode. Dans le corpus distribué, ce codage erroné a été corrigé.

### 2.2 Codage ASCII

Ce codage est la traduction simplifiée en ASCII de la transcription XML. Dans ce format (figure 1):

- ne sont conservés que les informations concernant le dialogue par lui-même (pas d'entête à l'exception de l'étiquette du dialogue concerné),
- ne sont pas conservées les informations d'alignement temporel
- est par contre conservée la segmentation en tours de parole. Chaque tour de parole se voit accorder un numéro spécifique par incrément. Pour un tour de parole donné, on précise ensuite à la ligne l'identité du locuteur ainsi que l'énoncé prononcé. Ce format permet toujours une représentation des chevauchements : dans ce cas, deux énoncés sont donnés dans un tour de parole particulier, avec toujours en tête d'énoncé la mention de l'identité du locuteur correspondant.

```
<01> hotesse
    h: U B S bonjour
<02> client
    c: oui bonjour madame j'aurais voulu avoir des renseignements pour e l'
l'inscription en A E S administration économique et sociale
<03> hotesse
    h: oui [pi] oui conserver je vais vous passer la personne
```

**Figure 1** : Extrait du corpus Accueil\_UBS : transcription orthographique (format ASCII).

### 2.3 Formats, DOC, ODT PDF

Ces formats sont la compilation, sous la forme d'un fichier unique, des fichiers ASCII décrits ci-dessus.

```
Dialogue 060
```

```
<01> hotesse
    h: U B S bonjour
<02> client
    c: oui bonjour madame j'aurais voulu avoir des renseignements pour e l' l'inscription en A E S
administration économique et sociale
<03> hotesse
    h: oui [pi] oui conserver je vais vous passer la personne
```

**Figure 2** : Extrait du corpus Accueil\_UBS : transcription orthographique (format DOC/ODT/PDF).

Le numéro de dialogue précisé dans cette compilation reprend les trois premiers chiffres des fichiers de transcription et audio correspondant. Par exemple, sur la figure 2, la numérotation *Dialogue 060* correspond aux fichiers de transcription 060\_0000003d.xml.