

Outils Linguistiques

Jean-Yves Antoine
LI - Université Rabelais de Tours
Jean-Yves.Antoine@univ-tours.fr



Outils linguistique — Mastere L&R — © J.Y. Antoine — 1

Outils Linguistiques

TRANSCRIPTION DE CORPUS ORAUX



Outils linguistique — Mastere L&R — © J.Y. Antoine — 2

Corpus oraux et transcription

• Utilité des corpus transcrits

- Etudes linguistiques sur la langue parlée : évitent un retour au signal audio pour des études lexicales, syntaxiques, sémantiques ou pragmatiques.
- TALP : apprentissage de modèles de langage



• Deux types de transcription

- **Phonétique ou prosodique** : prononciation phonétique exacte et/ou le rythme adopté (longueur des pauses etc...)
- **Orthographique** : transcription moins fine rendant compte sous forme écrite de ce qui a été prononcé : restitution **littérale** ce qui a été prononcé

Outils linguistique — Mastere L&R — © J.Y. Antoine — 3

Corpus oraux et transcription

• Le dilemme de la transcription orthographique

- Transcription littérale sans interprétation : comment assurer ce passage objectif, sans variabilité, de l'oral à sa transcription écrite ?
- Séparation phonétique (prononciation) et linge (syntaxe de la langue parlée)

Exemples

douskipudonctan [Queneau, Zazie dans le métro]
transcription des toponymes dans la langue du colonisateur
fautes d'accord : * *des grands hommes* ou *des grands-t-hommes* ?
ambiguïté : * *il marche dans la rue* ou *ils marchent dans la rue* ?

Règle d'or de transcription orthographique

- Transcrire littéralement ce qu'on entend et non pas ce qu'on croit entendre*
- *ne jamais corriger ce qu'on entend pour le rendre plus « acceptable »*
 - *ne pas inventer de nouvelles formes écrites pour simuler la prononciation orale*

Outils linguistique — Mastere L&R — © J.Y. Antoine — 4

Corpus oraux et normalisation

• Réutilisabilité des ressources existantes

- Transcription : coût très important (20 à 40h d'écoute pour 1h de signal)
- Normaliser pour optimiser l'utilisation de ces ressources

• Normalisation linguistique : conventions de transcription

- Objectiver la transcription
- Éviter les biais méthodologiques : une réalité linguistique doit être partout transcrite de la même manière.

• Normalisation du codage

- Normalisation technologique : permet la réutilisation des mêmes outils informatique pour la gestion des corpus
- Besoins et contenus différents d'un corpus à un autre
- Codage sous format structuré : langage de balisage XML
- TEI (*Text Encoding Initiative*)

Outils linguistique — Mastere L&R — © J.Y. Antoine — 5

Conventions de transcription

• Objectifs

- Définir des règles systématiques de comportement face à des observations problématiques : limiter la variabilité entre transcribers et entre corpus

• Limites

- Limite de l'objectivité : toute convention répond à des a priori théoriques.
- Toujours joindre les conventions à un corpus afin d'expliquer les choix de transcription qui ont été adoptés
- Recours à l'audio reste toujours possible

• Exemples de conventions (français)

- GARS-DELIC [Blanche-Benveniste, 1991]
- Parole Publique [Antoine, 2002]
- Transcriber www.etca.fr/CTA/gip/Projets/Transcriber/ [Barras et al, 1998]

Outils linguistique — Mastere L&R — © J.Y. Antoine — 6

Transcription : ne pas céder à l'écrit

• Phrase et énoncé

- la notion de phrase n'a pas de sens à l'oral. On distingue simplement des énoncés qui correspondent à une prise de parole ininterrompue du locuteur
« et donc il y avait également sans doute et qui va avec tout ce que je viens de dire avant il y avait aussi un homme mythique complètement mythique et c'est la / première et là en plus c'est intéressant parce que c'est la première image que j'ai eue de la Tchécoslovaquie lorsque je suis arrivé à la frontière »
- La transcription ne comporte donc pas de points, virgules ou points virgules. Les énoncés ne débutent pas par une majuscule, ceci pour faciliter l'analyse syntaxique ou lexicale des corpus (majuscule réservée aux noms propres)

• Segmentation en mots

pour permettre l'analyse lexicale, il est préférable de n'utiliser le tiret que pour les unités polylexicales insécables, indépendamment des normes de l'écrit

peux tu et non pas ~~peux-tu~~
porte-feuille et non pas ~~porte-feuille~~

Outils linguistique — Master LAR — © J.Y. Antoine — 7

Transcription : éviter la caricature

• Ne pas céder à une vision caricaturale de l'oral

Exemple : contractions phonétiques

Eviter une transcription orthographique qui colle trop à la prononciation

il y a et non pas y'a
il part et non pas i'part
je vais et non pas j'vais

MAIS...

Rendre compte de toute élision complète d'une unité lexicale :

il pleut pas et non pas il ne pleut pas si élision du discordanciel ne

Outils linguistique — Master LAR — © J.Y. Antoine — 8

Conventions de transcription : exemples

Majuscules / minuscules

- Majuscules réservées aux noms propres, épellations, sigles
- Quelles règles pour les associations nom communs / noms propres ?

Exemples

| | | |
|------------------------------|------|-------------------------------------|
| mairie de Paris | ou ? | <u>Mairie de Paris</u> |
| amicale bouliste de Lamastre | ou ? | <u>Amicale Bouliste de Lamastre</u> |
| crédit lyonnais | ou ? | Crédit Lyonnais |
| parti socialiste | ou ? | Parti Socialiste |
| Toulouse football club | ou ? | <u>Toulouse Football Club</u> |
| médecins du monde | ou ? | <u>Médecins du Monde</u> |
| l'île de la tentation | ou ? | <u>L'île de la Tentation</u> |

Choix de convention

- Majuscules pour marques identifiées ou dénominations correspond à un acronyme
- Guillemets + minuscules si minuscules seules trompeuses

Exemple « l'île de la tentation »

Outils linguistique — Master LAR — © J.Y. Antoine — 9

Conventions de transcription : exemples

• Epellation

Rendre compte de la prononciation : majuscules séparées par des blancs

Exemple *My name's Cholomondely Featherstonough CH O L O M O N D E L Y ...*

• Acronymes

Si ambiguïté, distinction entre épellation ou non (voire mixte)

| | | | | |
|--------------------|------|----------------------|------|------------------------------|
| S N C F | ou ? | SNCF | | |
| SMIC | ou ? | Smic | ou ? | S M I C |
| O N U | ou ? | ONU | | |
| CDROM.... | ou ? | CD.Rom... | ou ? | C D R o m cédérom |

• Chiffres / Nombres / Dates

Pour les nombres (i.e.>10) privilégier l'écriture en chiffre

1200 se transcrit mille deux cents ou douze cents
73,5 se transcrit soixante treize virgule cinq ou septante trois virgule cinq

Certaines conventions limitent la règle aux nombres à prononciation ambiguë

Exemple (Transcriber) 60 mille 200

Outils linguistique — Master LAR — © J.Y. Antoine — 10

Conventions de transcription : exemples

• Troncatures

Marquer la troncature tout en essayant de rendre compte de l'intention du locuteur

Exemple (Transcriber)

On en recouse après-de(main) euh lundi aucune ambiguïté

Alors je vous() ambiguïté vous, voue, voulais, voudrais

Eviter toute interprétation non certaine dans la complétion

• Prononciations inattendues, lapsus

- Ne pas représenter orthographiquement la liaison
- Transcription phonétique (ou glose de prononciation) à part

Exemple (Transcriber) moi aussi [pron=moi-z-aussi]

• Ambiguïté

Transcriptions alternatives et non pas interprétation personnelle...

Exemple (DELIC) ils/il marchent/marche dans la rue

Outils linguistique — Master LAR — © J.Y. Antoine — 11

Transcription et dialogue : tours de parole

Du fait du chevauchement des interlocuteurs, il est difficile de donner une définition claire de la notion d'énoncé, appelée dans ce cas **tour de parole** (*speech turn*). Deux solutions envisageables :

- centrer la définition sur chaque locuteur (période limitée par une prise et une fin de parole) et coder les chevauchements temporels

Exemple (alignement de type GARS/DELIC)

U1 : on a maintenant un camping-car à la place de notre caravane

S1 : ah ouais

U2 : ouais c'est bien plus pratique

- centrer la définition sur le dialogue (tour de parole = période d'interlocution)

Exemple (conventions Transcriber, PAROLE PUBLIQUE)

T1 : <U = on a maintenant un camping-car>

T2 : <U = à la place de>

<S = ah ouais>

T3 : <U = notre caravane ouais c'est bien plus pratique>

Outils linguistique — Master LAR — © J.Y. Antoine — 12

