# Creating and exploiting multimodal annotated corpora

Philippe Blache, Roxane Bertrand & Gaëlle Ferré

*Laboratoire Parole et Langage*
CNRS & Université de Provence

LREC 2008

# Introduction

- Multimodality
  - Information comes from different sources
  - Modalities interaction
    - Each source is partial, incomplete
    - They have to be synchronized
- Multimodal annotation
  - Goals
    - Usually focus on gesture description
    - Mainly in the perspective of communication
  - Conventions and schemes
  - Tools (Praat, Anvil, Elan, etc.)
- Our project
  - Linguistic description
  - Study of interaction: annotation of all domains
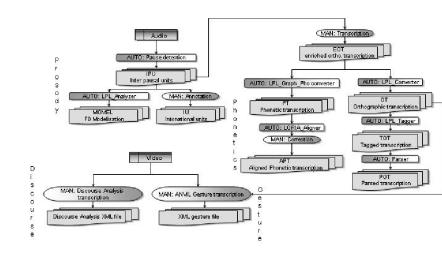  - Unrestricted data (natural situations)

# Outline

- The project
  - The CID corpus
  - The annotation process

- Results
  - Backchannels
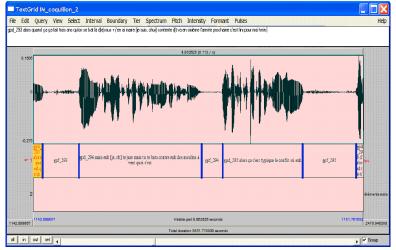  - Reinforcing gestures

- Perspectives

# The corpus

- *Corpus of Interactional Data*: 8 dialogs, 1 hour each ([Bertrand & al 07])

- Transcribed (orthographic, phonetic)

- Aligned

- Annotated
    - Prosody (intonation, units, contours, etc.)
    - Morphosyntax, syntax,
    - Discourse (markers, speech turns, etc.)
    - Gestures

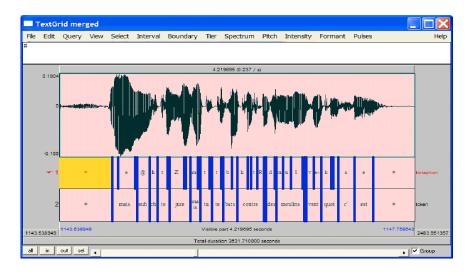# The annotation architecture

# Signal segmentation

- Interpausal units segmentation (IPUs)
- Syntactic units detection (pattern method)

## Transcription

- Precise transcription convention

- Transcription by 2 experts

- Enriched orthographic transcription (EOT), needed for different phenomena annotation and alignment (elisions, schwa, etc.)

- Generation of 2 transcription versions
  - Orthographic (for the NLP module)
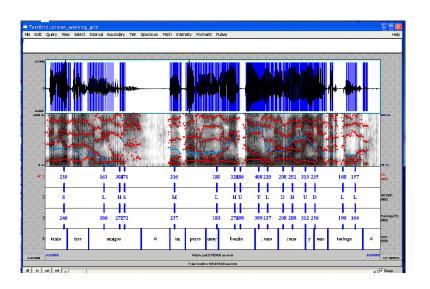  - Phonetic (for speech analysis)
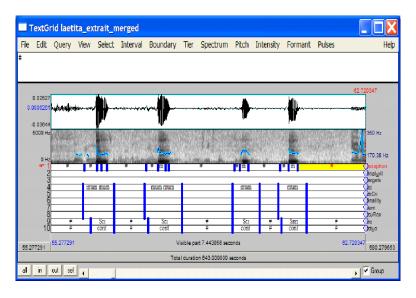
# Alignment

# Alignment

- Identifying the phoneme suite
    - Tokenisation
    - Grapheme-phoneme conversion

- Alignment tool
    - Input: list of phonemes + audio signal
    - Temporal localization of the phonemes in the signal

- Manual correction
    - Wrong boundaries
    - Overgeneration (false units)

- Tokens and phonemes are primary levels, used for anchoring other levels

# Intonation: INTSINT

# Discourse

# Gestures

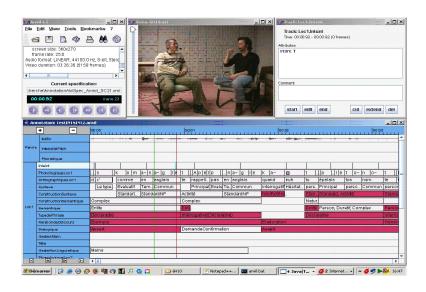# Summary of the tools

- Fully automatic
    - IPU segmentation
    - Phoneme alignment
    - Intonation
    - POS tagging

- Semi-automatic
    - Intonational units
    - Shallow parsing (still needs a segmentation tool)

- Manual
    - Transcription (we are experimented speech recognition as helping tool)
    - Other annotations

- Tools and resources available from the CRDO (http://crdo.fr/)

- Backchannels: minimal signal produced by the hearer. Vocal and gestural BCs (head movements, smiles and laughter, eyebrow movements, etc.), they have different functions

- *Example*:



| A | ah ouais nous on est rentré à (...) | dix heures dix heures et demi je crois du soir (...) |
| B | | nod |
| A | et elle a accouché à six heures je crois (...) | |
| B | | ah quand même ouais |
| | | head tilt |
| | | eyebrow raising |
| A | donc c'était ouais c'était quand même assez long quoi (...) | |
| B | | head tilt |

[A] oh yeah we were admitted at 10, 10.30 I think pm
[A] and she had the baby at 6 I think
[B] [oh yeah right?]
[A] so it was yeah it was quite long indeed

- **Question**: Do vocal and gestural BCs behave similarly? In what prosodic and morphological contexts do they appear?
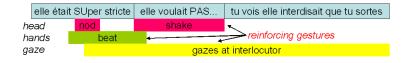
## Backchannels

- Vocal and gestural BCs show similar behavior but gestural BCs appear later than vocal ones

- Morphological and discursive context
  - After nouns, verbs and adverbs (words with semantic function)
  - Not after connectors (linking words between conversational units)

- Prosodic context
  - Gestural BCs: after accentual phrases (APs) and intonational phrases (IPs)
  - Vocal BCs: after IPs
  - Encouraged by specific contours (esp. rising), speakers gaze

- *Conclusion*: BCs occur at the end of some units, but not with possible turn change. They also play a role in the elaboration of discourse.

- Reinforcing gestures: eyebrow movements, gaze direction, head movements, highlighting discourse elements

- *Example*:



| elle était SUper stricte | elle voulait PAS... | tu vois elle interdisait que tu sortes |

head — nod / shake
hands — beat
gaze — gazes at interlocutor

*reinforcing gestures*

*she [the teacher] was super strict she didn't want... you see she forbade us to leave the room [during lessons]*

- **Questions**: What do gestures reinforce? Are they equivalent to known focalization phenomena?

# Reinforcing gestures: results

- No correlation with prosodic focalization, no gesture is associated with specific stress or contour

- Correlation with adverbs and connectors at the beginning of speech turns

- Correlation for metaphorics, no correlation for eyebrow movements

- *Conclusion*
  - Reinforcing gestures do not serve to express focus
  - Their role is more discursive than expressive

# Conclusion

- CID: large corpus, richly annotated
- Interest of multimodal annotated corpora
  - Study of natural language, in context
  - Study of interaction
- Problems
  - Standardisation: coding schemes
  - Synchronization of the different domains (+/- temporal)
  - Interfacing the different tools
- Perspectives
  - Information structure study
  - Description in terms of constructions (CxG)
  - Multimodal interaction for virtual reality