# TalkBank and CLARIN

Brian MacWhinney
CMU - Psychology
send to clarin@clarin.eu by
Oct 21 — 15+10 minutes

# Basic Questions

- How did language emerge in the species?

- How does it change?

- How is it learned?

- How is it processed?

- What are the results of damage and variation?

# Areas

| | | | | |
|---|---|---|---|---|
| **Children:** | CHILDES | PhonBank | Narrative | Bilingual |
| **Clinical:** | AphasiaBank | FluencyBank | Dementia | TBIBank |
| **Adult:** | CABank | TutorBank | GestureBank | ClassBank |
| **Multilingualism:** | BilingBank | SLABank | Online Tutors | DOVE |

# Funded Projects

| | CHILDES | TalkBank | AphasiaBank | PhonBank | FluencyBank | LangBank | HomeBank |
|---|---|---|---|---|---|---|---|
| Age of Project | 28 | 12 | 8 | 6 | 0.2 | 1 | 1 |
| Words (millions) | 59 | 47 | 1.8 | 0.8 | 0.5 | 2 | audio |
| Linked Media (TB) | 2.8 | 1.1 | 0.4 | 0.7 | 0.3 | - | 3.5 |
| # Languages | 41 | 22 | 6 | 18 | 4 | 2 | 2 |
| # Publications | 7000 | 320 | 256 | 480 | - | 4 | 5 |
| # Users | 2950 | 930 | 390 | 182 | 25 | - | 22 |
| # Web Hits (millions) | 4.1 | 1.3 | 0.3 | 0.1 | - | - | |

# 41 languages (including Cantonese)

*MIC:  睇 吓5 Sophie 畫 咗 啲 咩 圖畫 .

%mor: v|tai2=look_at asp|haa5=tentative n:prop|Sophie v|waak6=draw asp|zo2=perfective

 cl|di1=some wh|me1=what n|+n|tou4+n|waa2=drawing .

*CHI:  呢 啲 我哋 彈 琴 架 .

%mor: sfp|ne1=how_about cl|di1=some pro|ngo5-PL=I v|daan6=bounce n|kam4=piano

 sfp|gaa3 .

*SIS:   你 鍾意 Alicia 定係 呀 呀 Lulu 定係 Sophie 定係 Timmy ?

%mor: pro|nei5=you v|+v|zung1+n|ji3=like n:prop|Alicia conn|ding6hai6=or sfp|aa3

 sfp|aa3 n:prop|Lulu conn|ding6hai6=or n:prop|Sophie conn|ding6hai6=or n:prop|Timmy

 ?

# TalkBank Principles

- Community Driven

- Open access to data, media, derived corpora, and programs

- Standard format — CHAT, CHAT-XML, CHAT-CA

- CLAN programs running on CHAT format

- Transcripts linked to media

- Interoperable with other resources: R, Elan, Praat, SALT, Annis, CONLL, SpeechKitchen/Kaldi for ASR, LENA

- CHAT/PHON incorporates Praat

# CLARIN Principles in TalkBank

- CLARIN-B center

  - InCommon login through Shibbolet

  - OAI-PMH server for OLAC, VLO

  - DOI through HandleServer, EZ-Cite

- CLARIN-K center

  - focus on analysis of spoken language

  - tutorials for CLAN, video tutorials

  - help desk, 5 Google Groups discussion boards

- SamtaleBank as a CLARIN illustration

# Let's take a look on the web

❖ childes.talkbank.org

❖ talkbank.org

❖ homebank.talkbank.org

❖ sla.talkbank.org

❖ talkbank.org/access/SamtaleBank

❖ childes.talkbank.org/browser — Alicia at 3;3

❖ downloadable materials

# Major Methods

1. Corpus Analysis

2. Profiling

3. Microanalysis

# 1. Corpus Analysis

- FREQ - Frequency analysis

  - wild cards

  - word files (morality words, LIWC, medical)

- KWAL - Key word and line

  - matches highlighted

- COMBO - Regular expression matching

- Hits can be triple-clicked to go back to transcript and play

# LENA2CHAT

- ❖ 24 hour/day recordings in the home

- ❖ Much like Deb Roy's database and the "water" example, but open

- ❖ Huge ITS files reduced automatically to manageable CHAT files

- ❖ Check out http://homebank.talkbank.org

# Looking under the Hood

# MOR, POST, GRASP

- ❖ 41 languages, but only 11 have MOR/POST

- ❖ Cantonese, Danish, Dutch, English, French, Italian, Hebrew, Japanese, German, Mandarin, Spanish

- ❖ GRASP for English, German, Hebrew, Spanish, Mandarin

# MOR

- More declarative than FST

- Part-of-speech tuned to spoken language

- Easy to use once there is a grammar

- Hard to build the grammar (A-rules, C-rules)

- 98% accuracy for English

- POSTMORTEM rules for German declension

# Bilingual MOR

- *CHL: +" [- spa] <yo no la> [/] yo no la desmentí porque. [+ break]

- *CHL: what's my word against hers &ladadada .

- *CHL: +" [- spa] todos estamos con un calor and@s working@s .

- All words are tagged implicity; can be made explicit.

- Coding system makes code-switching junctures evident.

- Run English MOR, excluding [- spa], then Spanish MOR including [- spa]

# Dependency Graphs

Web service runs by triple-clicking on %gra line



Gallia est omnis divisa en partes tres.

# 2. Profiling – EVAL/KIDEVAL

- ❖ This all depends on MOR and GRASP

- ❖ Comparison database with s.d. scores

- ❖ IPSyn, DSS

- ❖ MLU, MLT

- ❖ TTR, vocD, MATTR

- ❖ Brown's 14 morphemes

- ❖ TIMEDUR

# EVAL

MLU, TTR
Verbs/Utt
% errors
% N, V, Aux, Adv, Conj,
  Pro
% PAST, PASTP, PL
Retracing, repetition



Select eval options

PLEASE SELECT AT LEAST ONE SPEAKER
Speaker: ● *PAR   ○ *INV   ○ *CLI

Database types:   [Deselect Database]   [Update Database]

○ Anomic   ○ Broca        ○ TransSensory
○ Global   ○ Wernicke     ○ TransMotor
○ Control  ○ Conduction   ○ NotAphasicByWAB
[Fluent]   [Nonfluent]    [All Aphasia]

Age range: [          ]   ○ Male only  ○ Female only

Gem choices:   [Deselect all gems]   [Select all gems]

○ Speech   ○ Cinderella   ○ Important_Event
○ Cat      ○ Umbrella     ○ Stroke
○ Flood    ○ Sandwich     ○ Window

[Cancel]                  [OK]

# Sample Output

# Error Analysis

- ❖ [*p] phonological p:w, p:n, p:m

- ❖ [* s] semantic  s:r, s:ur. s:uk, s:per

- ❖ [* n] neologism n:k, n:uk, n:k:s, n:uk:s

- ❖ [* d] dysfluency

- ❖ [* m] morphology m:a:0es  etc.

- ❖ [* f] formal lexical

- ❖ [+ gram] [+ jar] [+ es] [+ per] [+ cir]

# 3. Microanalysis (CA and Gesture)

# CHAT2ELAN

# CHAT2PRAAT - sociophonetics

- Highlight utterance bullet
- Send to sound analyzer
- Extracts audio from video
- In Praat, draw a picture

# CHAT2PHON

# CHAT in ANNIS

# CA Coding

```
↑    shift to high pitch; F1 up-arrow
↓    shift to low pitch; F1 down-arrow
↗    rising to high; F1 1
↗    rising to mid; F1 2
→    level; F1 3
↘    falling to mid; F1 4
↘    falling to low; F1 5
∞    unmarked ending; F1 6
≈    ≈continuation; F1 +
·    inhalation; F1 .
≈    latching≈; F1 =
≡    ≡uptake; F1 u
⌐    top begin overlap; F1 [
¬    top end overlap; F1 ]
L    bottom begin overlap; F1 {
⌐    bottom end overlap; F1 }
Δ    Δfaster Δ; F1 right-arrow
∇    ∇slower∇; F1 left-arrow
⁎    ⁎creaky⁎; F1 *
⁇    ⁇unsure⁇; F1 /
°    °softer°; F1 0
⊙    ⊙louder⊙; F1 )
=    ˍlow pitch ˍ; F1 d
     ˉhigh pitch ˉ; F1 h
☺    ☺smile voice☺; F1 l
⊗    ⊗breathy voice⊗ marker; F1 b
ʃ    ʃwhisperʃ; F1 w
ÿ    ÿyawnÿ; F1 y
♯    ♯singing♯; F1 s
§    §precise§; F1 p
∾    constriction∾; F1 n
◡    ◡pitch reset; F1 r
Ⱨ    laugh in a word; F1 c
„    Tag or sentence final particle; F2 t
‡    ‡ Vocative or summons; F2 v
```

# Gestural Detail

- Interaction / Sequence / Segment

- Each participant coded through sequence

  - Deedee 1a-1b-1c

  - Nina 1a-1b-1c

- Bullet links each segment back to transcript

- Coding: gaze direction, action, classification, meaning, language

- Rapport coding through gaze, smile, language

# Discourse Analysis

- CHAINS, KEYMAP, DIST

- CHIP

- PD (Propositional density)

- CI (Complexity index)

- SCRIPT + Speech Kitchen ASR

# Time Series – corpora to R

Alberto and
Jorge — I no go.

# Collaborative Commentary



```
@Begin
@Languages: en
@Participants: MOT Mother, CHI David Target_Child
@ID: en|rollins|MOT|||||Mother||
@ID: en|rollins|CHI|1;8.|||||Target_Child||
@Activities: book
*MOT: ahhah: look we can read books Tim .        Commentary (5)
%spa: $DHA:YY $DHA:RP

*MOT: it's a look and see <book> [>] .            Commentary (7)
%spa: $DHA:ST

*MOT: <ahhah> [>] we open it up and there are a set of eyes and there is a bird looking at David .
%spa: $DJF:ST $DHA:ST                                           Commentary (2)

*MOT: <the bear has a baby> [>] bottle .          Commentary (1)
%spa: $DHA:ST

*MOT: yes # David has baby <bottles> [>] .        Commentary (3)
%spa: $DRP:ST

*MOT: <oh> [>] .         Commentary (0)
%spa: $DHA:MK

*MOT: <there's a mirror> [>] .      Commentary (4)
%spa: $DJF:ST

*MOT: can David see <David> [>] .      Commentary (7)
%spa: $DHA:RQ

*CHI: 0 .
```

31

# Messages for CLARIN

❖ Importance of open access

❖ Importance of uniform transcription format linked to analysis programs

❖ Importance of focus on specific research communities for:

  ❖ corpus development

  ❖ tool development

  ❖ FUNDING

# Conclusions

- We need to expand TalkBank

- CLARIN can make wider use of TalkBank methods

- We can promotate TalkBank-CLARIN integration