# Formatting and Analysing a Learner Corpus with CHAT and CLAN

Kevin McManus & Nicole Tracy-Ventura

University of Southampton

# Introduction to CHILDES

- **Child Language Data Exchange System** / CHILDES

  (for a useful introduction to CHILDES see MacWhinney 2000 or http://childes.psy.cmu.edu/)

- **Talkbank**– The Database – primarily child language but also some language disorder data and bilingual data

- **CHAT** (Codes for the Human Analysis of Transcripts) – transcription procedures, a system for notation and coding

- **CLAN** (Computerised Language Analysis) – computer programs for searching and manipulating the data.

# Why do researchers use CHILDES?

- Well-supported; list-serves available

- International standing (over 1300 published studies)

- Flexible (not language-specific)

- Powerful (over 40 commands)

- Morphosyntactic tagger & direct search for morphosyntactic variables

- It's free!

# Example SLA Projects using CHAT/CLAN

- FLLOC (French Learner Language Oral Corpora) Project
  - 3.5 million words

- SPLLOC (Spanish Learner Language Oral Corpora) Project
  - 330,000 words

- LANG–SNAP (Language and Social Networks Abroad Project)
  - 680,000 words

# The FLLOC and SPLLOC projects
flloc.soton.ac.uk & splloc.soton.ac.uk

- Focus in on instructed learners of French and Spanish

- Dual aim:
    - constructing databases of oral learner data freely available to the research community
    - Substantive research agenda on French and Spanish SLA

- Collaboration between the Universities of Essex, Newcastle, Southampton and York

- Team Members: Florence Myles, Rosamond Mitchell, Laura Domínguez, Emma Marsden, Sarah Rule, Annabelle David, Maria Arche, Nicole Tracy-Ventura, Kevin McManus, Christophe dos Santos, Administrative and IT support

# LANG-SNAP project:
langsnap.soton.ac.uk

- French and Spanish L2

- Focus on speech and writing

- Longitudinal: 6 data collection cycles over 20 months

- Use of CHILDES Procedures: soundfiles, bulleted transcripts, tagged transcripts

- Advanced learners, native speaker controls

- Eventually accessible on the web via Talkbank and langsnap.soton.ac.uk/

- 680,000 words

# CHAT/CLAN Getting Started

- Download the CLAN program from the Childes website
  - childes.psy.cmu.edu/clan/
  - It is already on the computers in this lab.

- Download both the CHAT & CLAN manuals
  - CHAT: childes.psy.cmu.edu/manuals/chat.pdf
  - CLAN: childes.psy.cmu.edu/manuals/clan.pdf

# CHAT: Codes for the Human Analysis of Transcripts

- CHAT is the set of transcription conventions you need to follow.

Header

Start of transcript. Each t-unit goes on a separate line

# Bulleting: Linking the transcript and audio/video


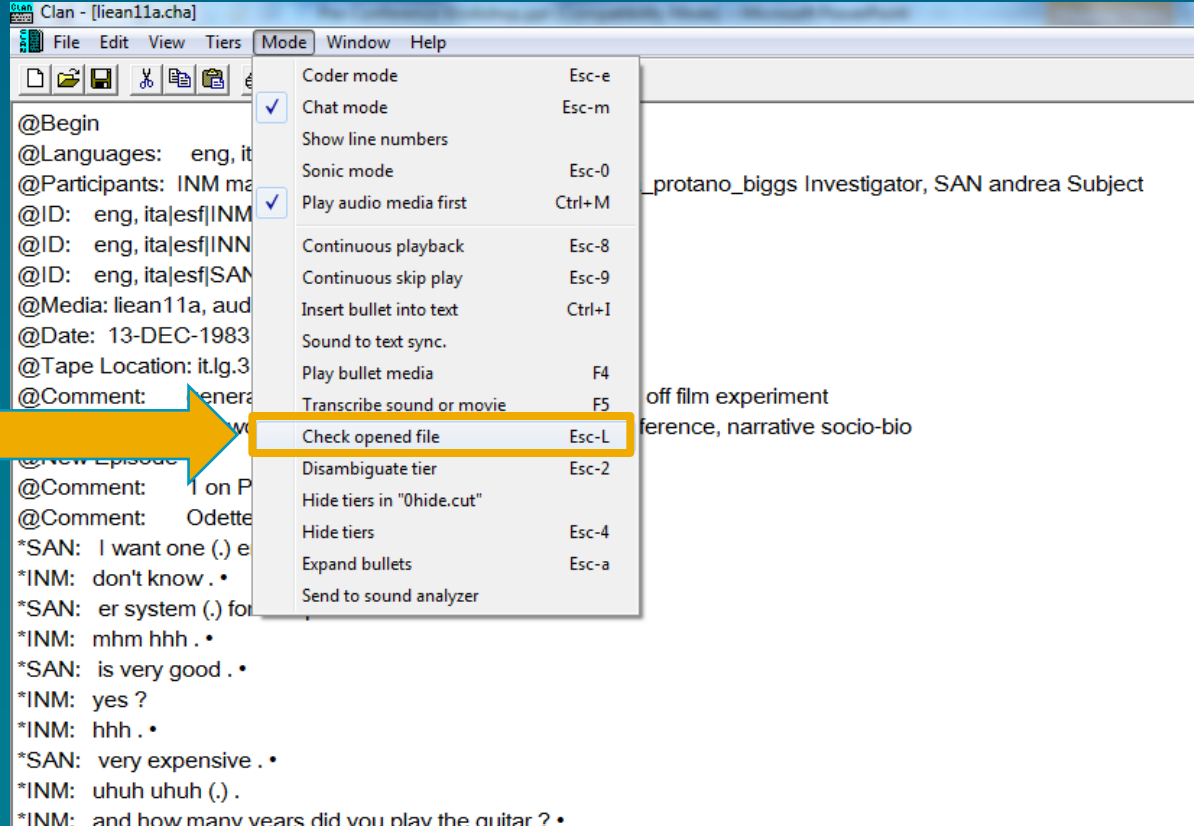
Bullets show that the audio and transcript are linked.

Now you can play any line in the transcript by pressing F4 when the cursor is over that line

# Transcribing in the CLAN Program

- 'Sonic Mode': transcribe and bullet at the same time.

- 'Transcriber Mode': transcribe with video and link already transcribed files with audio.

- 'Soundwalker': similar to the old transcriber foot pedal.

# Once you've finished transcribing: The 'Check' command

▪ Used to make sure your transcript doesn't include any technical errors.

# CLAN, the analysis program

- Over 40 built-in programs that work on CHAT files
  - Windows and Mac versions


- Some commands are run on the CHAT file.

- Some commands are run on the 'tagged' or 'MOR' file

# Obtaining a Frequency List

- Command: freq



To run commands you need to type them in the 'Commands Window' or select them from the 'Progs' menu

'Freq' is the command to get a frequency count. You also need to specify the file you want to use.

# Freq results

Type and token counts are provided too

```
2 up
3 used
13 very
1 visit
2 wait
1 waitress
2 want
1 warm
5 was
1 washer
1 watch
7 we
2 week
5 well
4 went
1 westminster
7 what
8 when
4 where
3 which
1 who
4 will
4 with
1 wonder
1 words
12 work
1 write
54 yeah
1 year
6 years
25 yes
69 you
6 your
2 youre
--------------------------------
 304  Total number of different item types used
1282  Total number of items (tokens)
0.237  Type/Token ratio

04dec12[E|TEXT]  148
```

Clan - [CLAN Output]

File  Edit  View  Tiers  Mode  Window  Help

# Saving results in a separate file: +f

- By adding +f to your command, the CLAN program will save the results in a separate file.

# Measure of Lexical Diversity: 'D'

- D is a measure that estimates lexical diversity while accounting for text length

- Command = vocd

```
50    100   0.7680   0.054   63.559

D: average = 59.200; std dev. = 2.244
D_optimum      <59.19; min least sq val = 0.000>

tokens  samples  ttr    st.dev    D
 35     100    0.8097   0.056   60.297
 36     100    0.7964   0.066   56.069
 37     100    0.7938   0.067   56.527
 38     100    0.7963   0.065   59.152
 39     100    0.7913   0.065   58.498
 40     100    0.7830   0.063   56.506
 41     100    0.7920   0.053   61.800
 42     100    0.7924   0.054   63.507
 43     100    0.7788   0.058   58.968
 44     100    0.7839   0.059   62.543
 45     100    0.7733   0.060   59.365
 46     100    0.7689   0.060   58.845
 47     100    0.7747   0.060   62.592
 48     100    0.7692   0.058   61.511
 49     100    0.7622   0.050   59.872
 50     100    0.7606   0.053   60.413

D: average = 59.779; std dev. = 2.179
D_optimum      <59.81; min least sq val = 0.000>

VOCD RESULTS SUMMARY
====================
  Types,Tokens,TTR:  <258,787,0.327827>
  D_optimum  values:  <59.58, 59.19, 59.81>
  D_optimum average:  59.53
```

You can also just get results for one speaker by specifying the tier
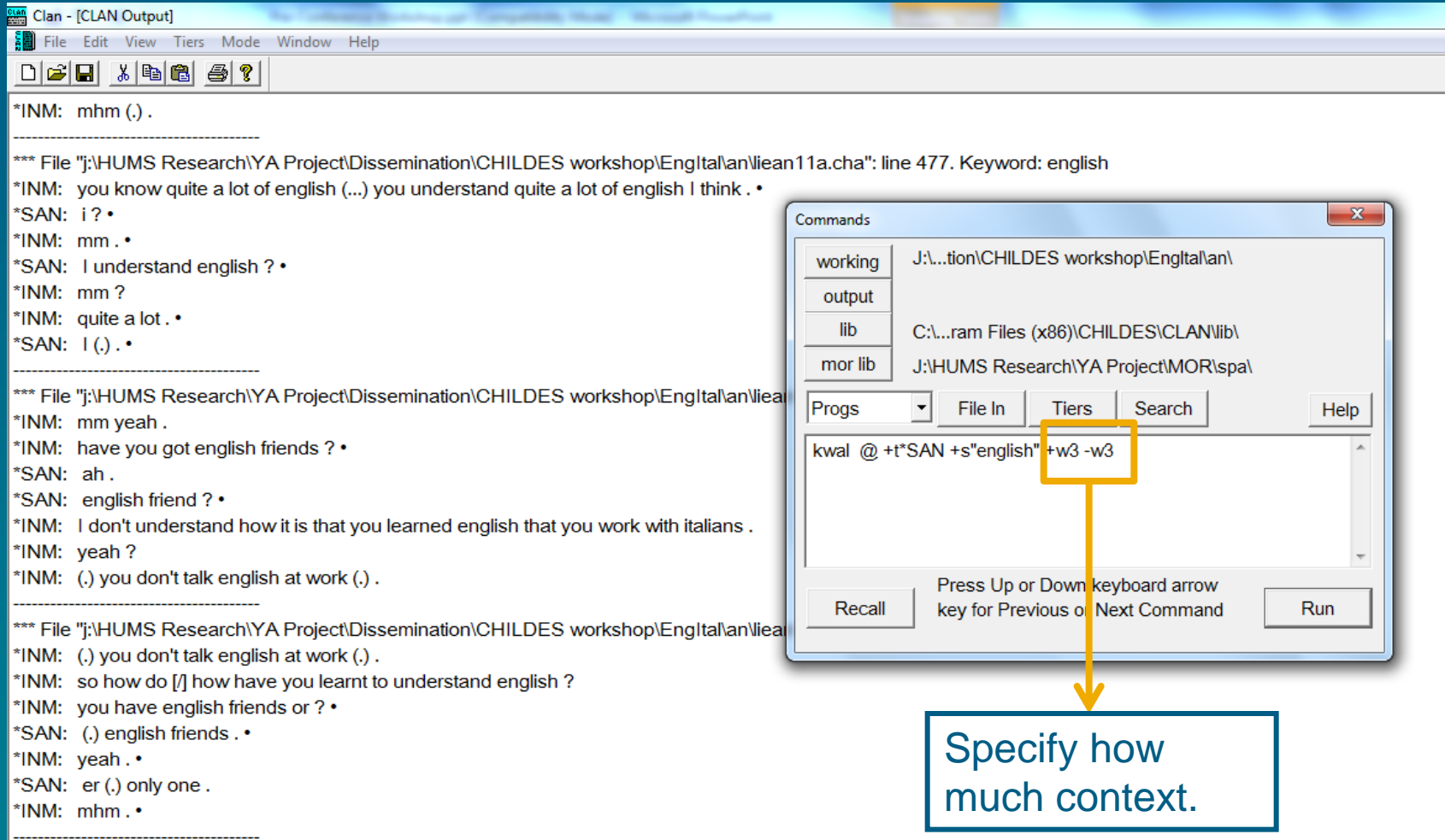
Commands

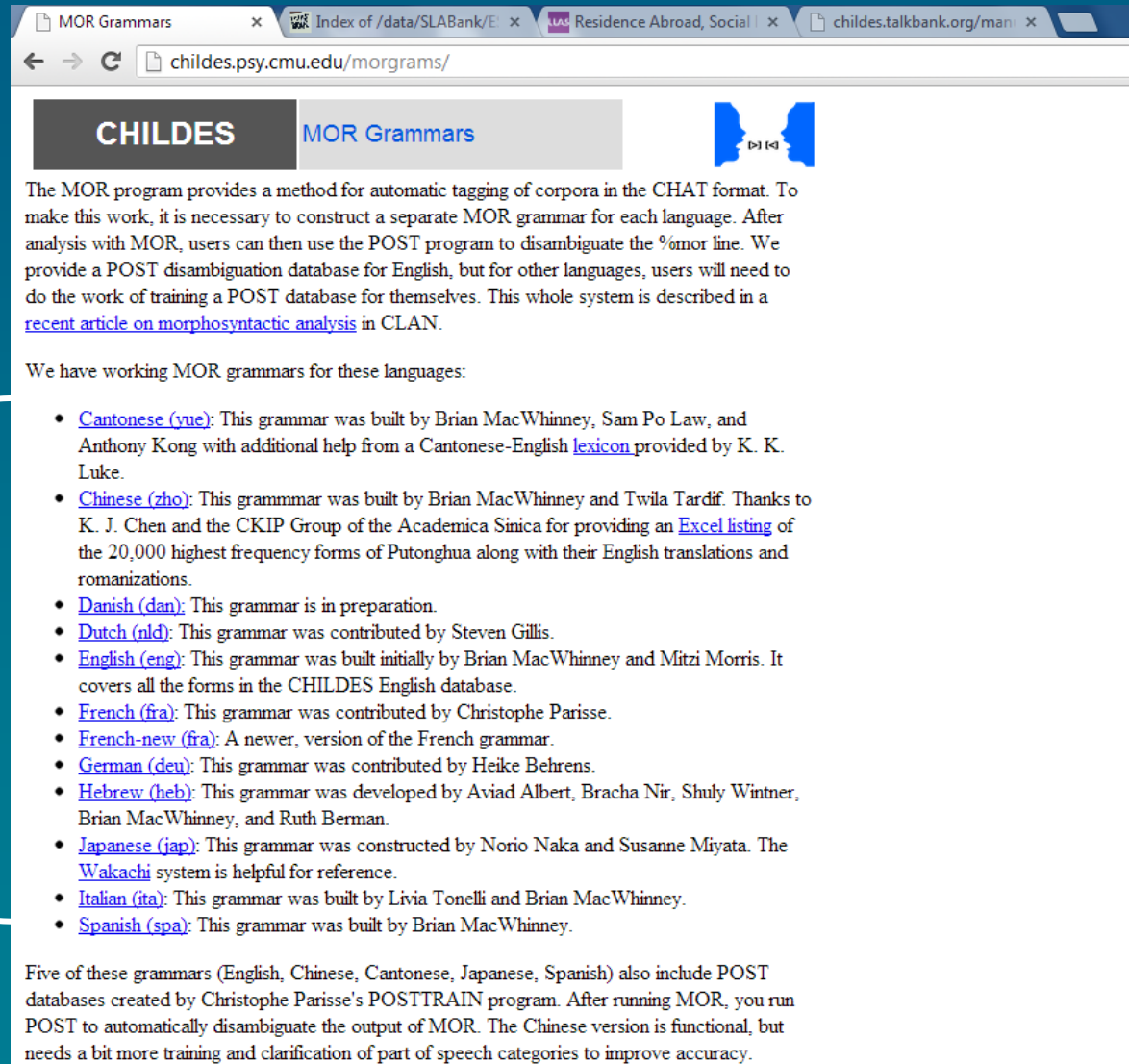| working | J:\...tion\CHILDES workshop\Engltal\an\ |
| output | |
| lib | C:\...ram Files (x86)\CHILDES\CLAN\lib\ |
| mor lib | J:\HUMS Research\YA Project\MOR\spa\ |

| Progs | File In | Tiers | Search | | Help |

```
vocd  @  +t*SAN
```

Recall    Press Up or Down keyboard arrow key for Previous or Next Command    Run

# Considerations for D

- Run freq to see what words get counted 'as words'. Spelling mistakes could impact the result

- You can have some 'words' excluded from the D calculation using some additional commands

- You may only want the result from a specific speaker so you need to specify in the command

# Search for words or sequence of words: COMBO

- Command= combo

- Use +s (string) and add the word in quotes: e.g., "english"

# Search for words in context: KWAL



Specify how much context.

# Part of Speech Tagger: MOR

Available in a number of languages

# English MOR

- You must first download the MOR onto your computer. Make sure you know where to locate it.

- MOR provides a morphosyntactic 'tag' for each word in the corpus, adding an additional layer of grammatical information

- The benefits of having a tagged text is that you can search by grammatical category (e.g., noun, verb, etc).

# Selecting the MOR Library
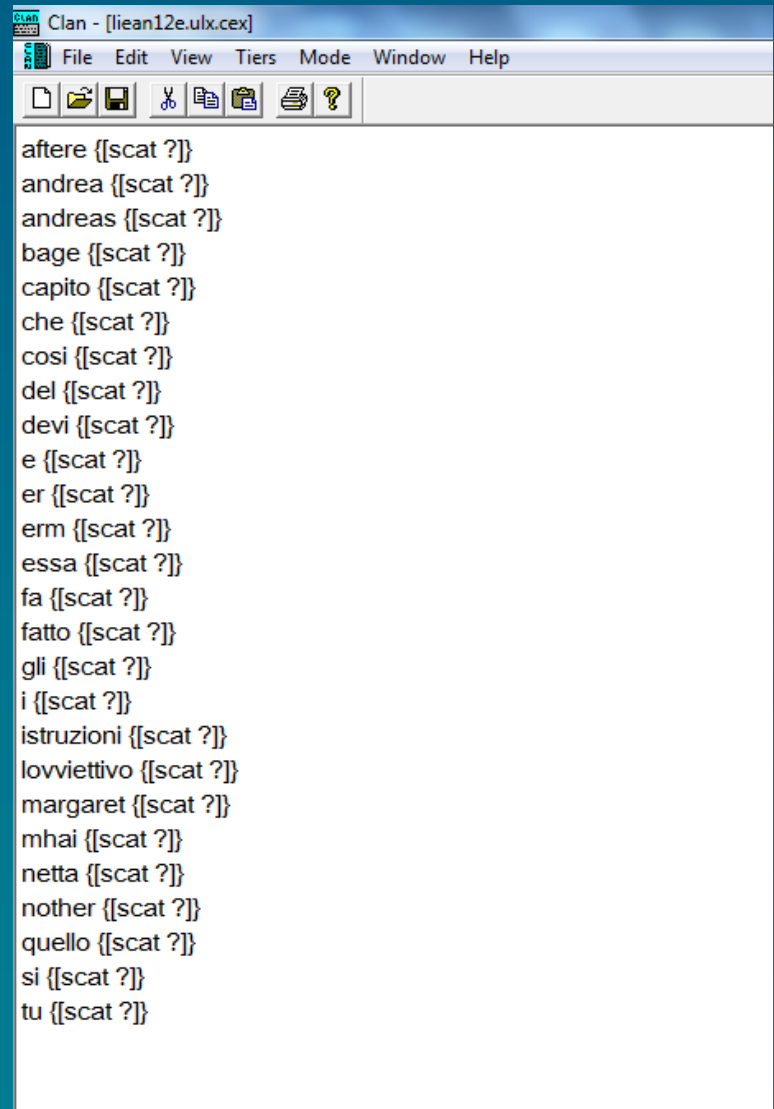
- Make sure you select the right location

# MOR–ing Files: Step 1

- Check for 'new' vocabulary

- The MOR dictionary or 'lexicon' that the program uses has been built word by word.

- Therefore, some real words in your corpus may not already be in the lexicon and you will need to add them manually.

# Checking for 'new' words

- Command=
  mor +xl

- A list of these words is automatically created in a new file.

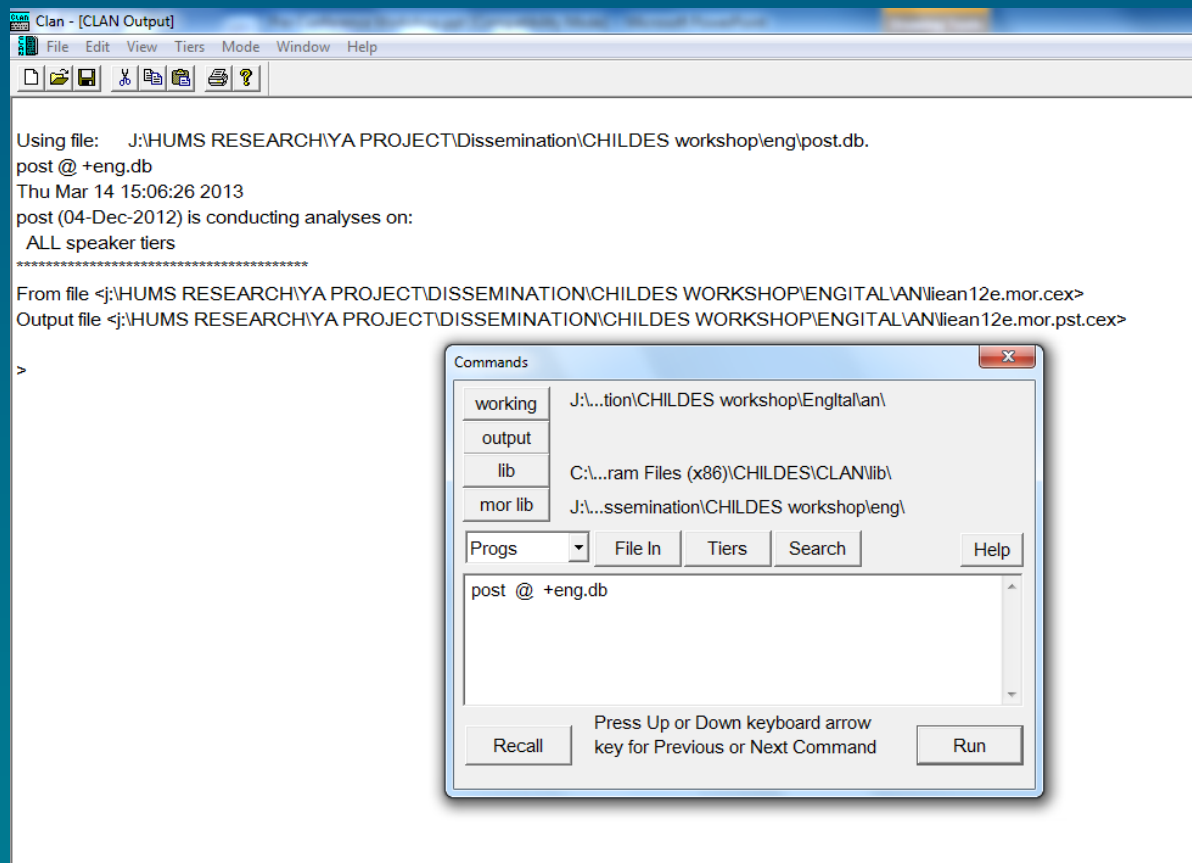- Some will be words to add to the lexicon. Others might be transcription errors

Clan - [liean12e.ulx.cex]

File  Edit  View  Tiers  Mode  Window  Help

aftere {[scat ?]}
andrea {[scat ?]}
andreas {[scat ?]}
bage {[scat ?]}
capito {[scat ?]}
che {[scat ?]}
cosi {[scat ?]}
del {[scat ?]}
devi {[scat ?]}
e {[scat ?]}
er {[scat ?]}
erm {[scat ?]}
essa {[scat ?]}
fa {[scat ?]}
fatto {[scat ?]}
gli {[scat ?]}
i {[scat ?]}
istruzioni {[scat ?]}
lovviettivo {[scat ?]}
margaret {[scat ?]}
mhai {[scat ?]}
netta {[scat ?]}
nother {[scat ?]}
quello {[scat ?]}
si {[scat ?]}
tu {[scat ?]}

# Step 2: MOR-ing

- Command: mor

- Specify the file (or files) and the tier.

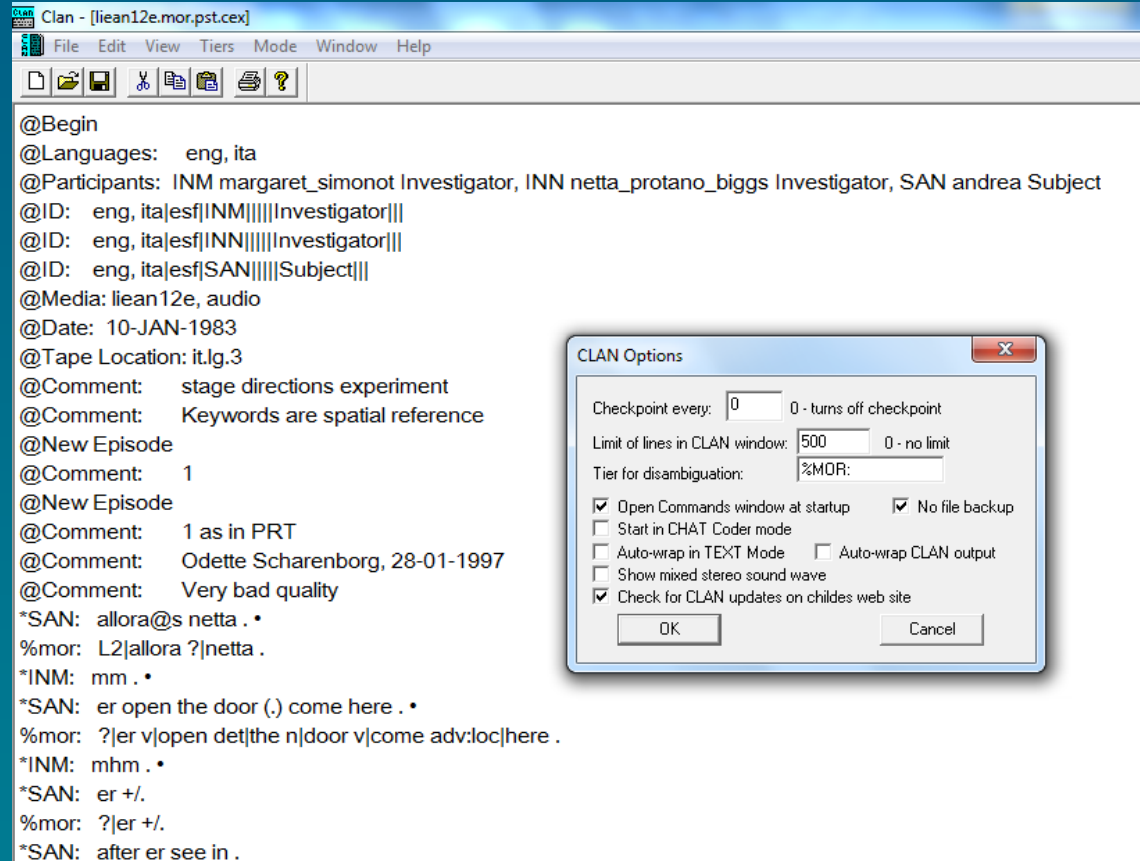- A new file is created with a new extension

# Step 3: Automatic Disambiguation

- Command:
  post +eng.db

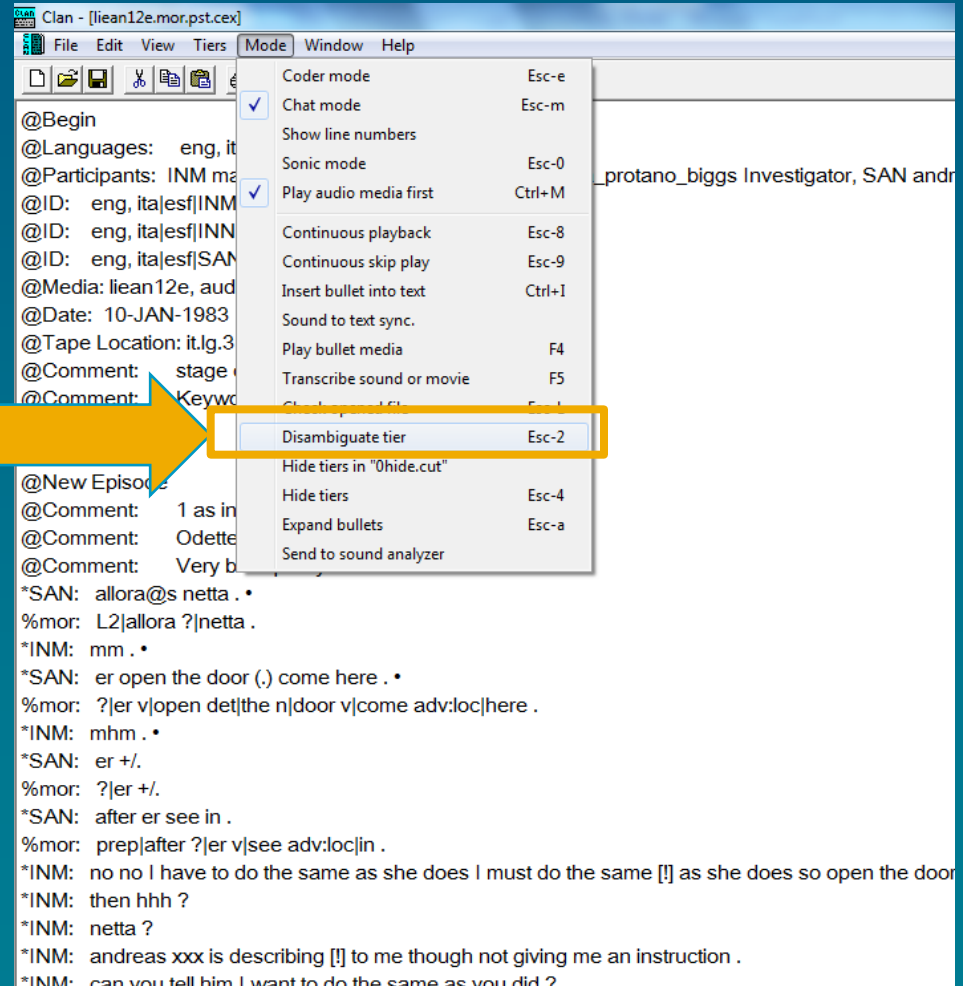- The program creates a new file with a new extension

# Step 4: Manual Disambiguation

- Open CLAN

- Under 'Edit', select 'CLAN Option'

- Make sure the 'tier for disambiguation' is %MOR
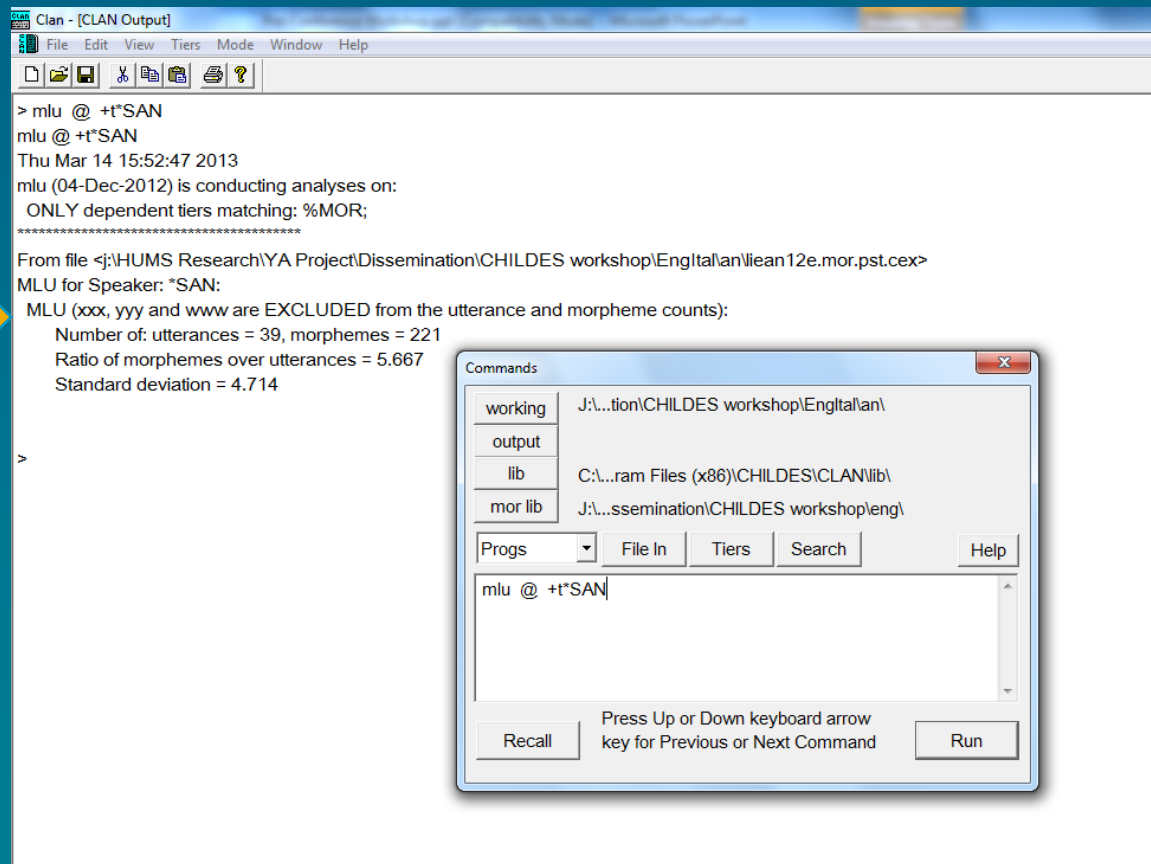
# Manual Disambiguation

- Under 'mode', select 'Disambiguation tier'

- If any tags need disambiguating, you'll be asked at the bottom of the page.

# Programs with a MOR file:
# Mean Length of Utterance (MLU)

- ## Command: mlu

Results are provided automatically or you can save in a separate file with +f command.

# Frequency list of tags

- **For example: We can use the freq command to get a list of verbs**

# EVAL: results of many commands all at once

- Currently works for English files only

  type, tokens, TTR, number of nouns, amount of retracing, repetition, etc



| | A | B |
|---|---|---|
| 1 | File | liean12e.mor.pst.cex |
| 2 | Speaker ID | eng, ita\|esf\|SAN\|\|\|\|\|Subject\|\|\| |
| 3 | Duration | 00:00:04 |
| 4 | Total Utts | 41 |
| 5 | MLU Utts | 39 |
| 6 | MLU Words | 5.385 |
| 7 | MLU Morphemes | 5.667 |
| 8 | types | 75 |
| 9 | tokens | 229 |
| 10 | TTR | 0.328 |
| 11 | Clause/Utt | 0.78 |
| 12 | Word Errors | 0 |
| 13 | Utt Errors | 0 |
| 14 | Nouns | 35 |
| 15 | Plurals | 8 |
| 16 | Verbs | 32 |
| 17 | 3S | 3 |
| 18 | 1S/3S | 0 |
| 19 | PAST | 1 |
| 20 | PERF | 0 |
| 21 | PROG | 1 |
| 22 | prep\| | 21 |
| 23 | adv\| | 13 |
| 24 | conj\| | 4 |
| 25 | det\| | 30 |
| 26 | pro\| | 3 |
| 27 | retracing[//] | 0 |
| 28 | repetition[/] | 6 |

# Exporting files to/from other programs

- Commands: CHAT2ELAN, CHAT2PRAAT

- Also, text files can be imported into CHAT:
  - Command: TEXTIN

# Let's do some practice!

- Download 'EngItal.zip' from Talkbank:
  - talkbank.org/data/SLABank/ESF/


- Move the files to an appropriate folder

- Open CLAN and click on 'working'. Locate the folder with these files.

# Run some commands!

- Try the ones we've demonstrated on the plain transcripts:
  - Freq
  - Vocd
  - Combo
  - Kwal

- You can also try MORing. First, download the English MOR:
  - childes.psy.cmu.edu/morgrams/

# To learn more:

- Join the CHILDES online community
  - childes.psy.cmu.edu/tools/email.html
  - Visit the FLLOC and SPLLOC websites

# Thank you!

K.McManus@soton.ac.uk

N.Tracy-Ventura@soton.ac.uk