

## **CLAN pour les Néophytes**

## Découvrir, comprendre et utiliser un nouvel outil de transcription et d'analyse des interactions

Stéphanie Caët Université Sorbonne Nouvelle stephanie.caet@gmail.com





Outils et Recherches pour les Corpus d'Acquisition

stephanie.caet@gmail.com





# Processus de collecte de données



Introduction





# Processus de collecte de données

Enregistrement 1h voire 1/2 journée ! de la séance Numérisation et Quelques heures Compression du film Transcription de la vidéo et 30h+10h correction de la transcription Toute la vie Analyse de la transcription et de la vidéo d'un chercheur ! **Processus long et coûteux** Données partiellement exploitées Introduction

3/79



## Dans cette présentation ...

- 1. Le Projet CHILDES
- 2. Que peut-on faire avec CLAN ?
  - 1. Lire et analyser des transcriptions de la base de données CHILDES
  - 2. Transcrire
  - 3. Aligner la transcription avec la vidéo ou le son
  - 4. Faire des analyses automatiques sur les transcriptions



## Dans cette présentation ...

## 1. Le Projet CHILDES

- 2. Que peut-on faire avec CLAN ?
  - 1. Lire et analyser des transcriptions de la base de données CHILDES
  - 2. Transcrire
  - 3. Aligner la transcription avec la vidéo ou le son
  - 4. Faire des analyses automatiques sur les transcriptions



## 1. Le Projet CHILDES CHIId Language Data Exchange System

### **Brian MacWhinney**



### \_\_\_\_\_

**Catherine Snow** 



1er corpus : Brown Corpus (Adam, Sarah & Eve) Au début, pas d'audio ...

Le Projet CHILDES



## 1. Le Projet CHILDES CHIld Language Data Exchange System

- Child Language Data ...
  - Interactions naturelles ou semi-guidées
  - 130 corpora partagés
  - dans 26 langues différentes
- ... Exchange System
  - Base de données CHILDES avec fichiers vidéo/son et transcriptions
  - Logiciel de lecture et d'analyse : CLAN
  - Conventions de transcriptions communes : format CHAT



### **1. Le Projet CHILDES** Site Internet : http://childes.psy.cmu.edu



Le Projet CHILDES



## Dans cette présentation ...

### 1. Le Projet CHILDES

## 2. Le logiciel CLAN

- 1. Lire et analyser des transcriptions de la base de données CHILDES
- 2. Transcrire
- 3. Aligner la transcription avec la vidéo ou le son
- 4. Faire des analyses automatiques sur les transcriptions





Le logiciel CLAN



#### For Windows:

CLANWin is for Windows XP/2000/NT/Vista. Windows 95, 98, or ME are no longer supported. Windows installation involves clicking on the installation file and following the directions given by InstallShield. If you have an older version of CLAN on your machine, InstallShield will overwrite it with the newer version.

Also, you will need to install QuickTime and Unicode fonts.



#### Links Manuals The TalkBank Database **CHAT Transcription IASCL** information **CLAN Programs** Other Child Language sites **Database Manuals** Research based on CHILDES BTS sign transcription system **Related Software Training Videos Phonology and Fonts** Teaching with CHILDES Phon & PhonBank Topics in language acquisition. Unicode and IPA for Mac Teaching Tips and Resources. Unicode and IPA for Windows Child Language Bibliographies **Special Procedures** Morphology and Lexicon Procedures and tools for CA analysis Part of Speech Analysis by MOR MRC lexical dictionary Working with digitized video Working with digitized audio Syntactic analysis by GRASP The Computerized Comprehension ChildFREQ Site and Paper Task

Télécharger la grammaire de la langue étudiée

Le logiciel CLAN



#### CHILDES MOR Grammars



The MOR program provides a method for automatic tagging of corpora in the CHAT format. To make this work, it is necessary to construct a separate MOR grammar for each language. After analysis with MOR, users can then use the POST program to disambiguate the %mor line. We provide a POST disambiguation database for English, but for other languages, users will need to do the work of training a POST database for themselves. This whole system is described in a recent article on morphosyntactic analysis in CLAN.

We have working MOR grammars for these languages:

- <u>Cantonese (yue)</u>: This grammar was built by Brian MacWhinney relying on a Cantonese-English lexicon provided by K. K. Luke.
- <u>Chinese (zho)</u>: This grammar was built by Brian MacWhinney and Twila Tardif. Thanks to K. J. Chen and the CKIP Group of the Academica Sinica for providing an <u>Excel listing</u> of the 20,000 highest frequency forms of Putonghua along with their English translations and romanizations.
- · Danish (dan): This grammar was written by Brian MacWhinney
- Dutch (nld): This grammar was contributed by Steven Gillis.
- English (eng): This grammar was built initially by Brian MacWhinney and Mitzi Morris. It covers all the forms in the CHILDES English database.
- French (fra): This grammar was contributed by Christophe Parisse.
  - . German (deu): This grammar was contributed by Heike Behrens.
  - <u>Hebrew (heb)</u>: This grammar was contributed by Sigal Uziel-Karl and Bracha Nir-Sagiv of Tel-Aviv University in Israel.
  - Japanese (jap): This grammar was constructed by Norio Naka and Susanne Miyata. The <u>Wakachi</u> system is helpful for reference.
  - Italian (ita): This grammar was built by Livia Tonelli and Brian MacWhinney.
  - Spanish (spa): This grammar was built by Brian MacWhinney and Monica Sanz









## 2. Le logiciel CLAN Une interface d'éditeur de texte

### Quand on ouvre CLAN ...

	newfi	ile.cha				- 5
ace de traiter	ment de texte,	dans laquell	e on peu	t écrire		0
ment.						-
f	face de traite ment.	face de traitement de texte, ement.	face de traitement de texte, dans laquell ment.	face de traitement de texte, dans laquelle on peu ement.	face de traitement de texte, dans laquelle on peut écrire ement.	face de traitement de texte, dans laquelle on peut écrire ement.

04nov10[EICHAT] \* 2







## 2. Les menus

	New	ЖN	Undo		жz	Ж1	Co	der mode {Eso	:-e}		Comma	nds	¥р
	Open	жο	Redo		ЖY	¥2	✓ Ch	at mode {Esc-	m}		Toggla	Movie Text	961
	Close	ЖW				ж3	✓ She	ow line numbe	ers		Toggie I	Novie-Text	<del>м</del> /
	Save	9995	Cut		ЖX	<b>#</b> 4	So	nic mode {Esc-	-0}		walker	Controller	
	Save As	<del>6</del> 5	Сору		жс	Ж5	Pla	v audio media	first	ЖМ	Special	Characters	
	Save Last Clip As		Paste		жv	ж6		.,			WEB Dat	ta	
	sure cust enp / is		Select All		ЖA	¥7	Co	ntinuous play	back {Esc-	8}		ale a	
	Page Setup		Et al.			¥8	Co	ntinuous skip	play {Esc-	9}	✓ newfile.	cha	
	Print	ЖP	Find		#F	¥9	Ins	ert bullet into	text	жı			
	Print Selection		Enter Selection		ЖE	жc	So	und to text sy	nc.				
	Ouit	<b>#</b> 0	Find Same		жG	Undata	Pla	y bullet media	a	F4			
	Quit	0002	Replace		ЖR	Update	Tra	anscribe sound	d or movie	F5			
			Replace and Find	Next	жн	ID neaders	She	ow movie thur	nbnails				
			Go To Line		жL		Ch	eck opened fil	le {Esc-L}				
Help			To Upper Case				Dis	sambiguate tie	er {Esc-2}				
Therp			To Lower Case				Hic	de tiers in "Ohi	ide.cut"				
Col	nmands and Shorte	uts	To Lower Case				Hic	de tiers {Esc-4	}				
Sav	e Them to File		CLAN Options		₩;		Ex	pand bullets {	Esc-a}				
			Select sound ana	lyzer			Ser	nd to sound a	nalyzer				
			Select F5 option										
			Define line numb	bers									
			Set WEB Data URL	L									
			Thumbnails Setti	nas									

1

### Menus généraux

### Menus spécifiques à CLAN



## 2. Le logiciel CLAN La fenêtre de commandes



Pour ouvrir la fenêtre de commandes : Windows>Commands ou ctrl/ d

Le logiciel CLAN



## Dans cette présentation ...

### 1. Le Projet CHILDES

## 2. Le logiciel CLAN

- 1. Lire et analyser des transcriptions de la base de données CHILDES
- 2. Transcrire
- 3. Aligner la transcription avec la vidéo ou le son
- 4. Faire des analyses automatiques sur les transcriptions



#### CHILDES Child Language Data Exchange System



CHILDES is the child language component of the <u>TalkBank</u> system. TalkBank is a system for sharing and studying conversational interactions.

#### System

Ground rules

**Guidelines for Contributors** 

**Overviews and Introductions** 

Membership list

How to subscribe to Mailing Lists

Links

The TalkBank Database

**IASCL** information

Other Child Language sites

#### **Programs and Database**

Downloadable Database

Browsable Database

The CLAN Program

Derived Corpora and Counts

Manuals

**CHAT Transcription** 

**CLAN Programs** 

**Database Manuals** 



La base de données

15/79



### CHILDES The Database Manuals



The guides to the CHILDES database are divided by subject matter into these files.

- Introduction to the Database
- American English
- British English
- Bilingual Acquisition
- Language Disorders
- Narratives
  - Germanic Languages
  - Romance Languages
  - <u>Slavic Languages</u>
  - East Asian Languages
  - <u>Celtic Languages</u>
  - Other Languages

$\Theta \Theta \odot$	Opening 08romance.doc
You have cho	osen to open
08romanc	e.doc
which is a:	doc File
from: http	://childes.psy.cmu.edu
What should	I Firefox do with this file?
O Open wi	ith Choose)
• Save File	2
Do this	automatically for files like this from now on.
	Cancel OK

#### Télécharger les manuels de base de données



#### Quand on ouvre le fichier 08Romance.doc ... Romance Corpora



This is a guide to CHILDES data on the acquisition of Romance languages. For a general introduction to the CHILDES database, please consult <u>intro.pdf</u>. The links in the table below are clickable, as are the thumbnails to the left.

+ Age Range N Corpus Comments French -1:9.18-2:5.27 1 Longitudinal study of language developmentin French with recordings made over a 28-month Champaud on page 4 period French - Geneva 1 Case study of the child Marie page 9 Test samples from a wide preschool population French -3:6 -5:6 489 Hammelrath page 11 Dizygotic twins French -1-3 2 Hunkeler page 12 French - Kern 0:7-2:0 4 Phonological development page 13 2:1.19-3:3.12 Longitudinal study with weekly recording French -1 Leveille page 14 sessions French - Lyon on 1-3 6 Longitudinal study with linked media for page 16 studying phonological development French -5:6-11:6 Cross sectional study focusing on speech acts 36 s in Kindergarten, third- and fourth-grade chil-Montreal page 17 dren playing "veterinarian" and "assistant" 0:7-3:5 Longitudinal video study of three children in French - Paris on 3 page 21 Paris French – Pauline Case study 1,2-2,6

Télécharger les manuels de base de données







English - USA

#### 31. Providence

Katherine Demuth Department of Linguistics Macquarie University Sydney, NSW 2109 Australia katherine.demuth@ling.mq.edu.au Références de CHILDES (MacWhinney, 2000)

Références corpus

+

Katherine Demuth and her research assistants at Brown University complication Providence corpus from 2002-2005. The corpus contains longitudinal audio/video recordings of 6 monolingual English-speaking children's language development from 1-3 years during spontaneous interactions with their parents (usually their mothers) at home. The aim of the study was to provide a corpus of phonetically transcribed data, with linked acoustic files, for the purpose of studying early phonological and morphological development. Collection and transcription of the Providence Corpus (and the similar Lyon Corpus for French) was supported by NIMH grant #1ROIMH60922.

#### Participants

The participants included 3 boys (Alex, Ethan, William) and 3 girls (Lily, Naima, Violet). Each child was recorded for 1 hour every 2 weeks beginning at the onset of first words. Two of the children have denser corpora, with weekly recordings from 1;3-2;10 (Naima) and 2;0-3;0 (Lily). The three girls (Lily, Naima, and Violet) were also recorded monthly from 3-4 years. The total corpus consists of 364 hours of speech. Audio is available for all children and video is available for all children except Ethan, who was diagnosed with Aspergers Syndrome at the age of 5.

#### Transcription

Both adult and child utterances were orthographically transcribed using CLAN conventions, with the audio/video files linked. Trained transcribers then carried out a broad phonemic (SAMPA > Unicode) transcription of the child utterances. A second trained coder then retranscribed 10% of each recording. Reliability scores ranged from 80%-98% (discounting voicing errors) on this second segment.

Please cite the following reference when using the corpus:

Références

du corpus

Demuth, K., Culbertson, J. & Alter, J. 2006. Word-minimality, epenthesis, and coda licensing in the acquisition of English. Language & Speech, 49, 137-174.

#### Télécharger les manuels de base de données

18/79





19/79

Télécharger les fichiers de la base de données



### CHILDES The CHILDES Database



The CHILDES database contains transcript and media data collected from conversations between young children and their playmates and caretakers. Conversations with older children and adults at available from <u>TalkBank</u>. All of the data is transcribed in CHAT and CA/CHAT formats. The use of all CHILDES and TalkBank data is governed by the <u>Gnu</u> <u>Public License (GPL)</u>. Please remember to read the <u>database manuals</u> and to follow the guidelines for <u>data-sharing</u>.

### **Downloading TalkBank Data**

- Downloadable Transcripts
- Downloadable Audio and Video

!! Mettre transcriptions et vidéos dans le même dossier



#### Using local transcripts and local media

- You need to download the transcripts from "zipped transcripts" below directory of <u>zipped XML data files</u>.
- 2. If the corpus is linked to audio or video media, you need to download
- 3. You need to download and install the CLAN program.
- 4. To open a transcript, you double-click on it. If there is associated med

	Zipped Transcripts
	English - USA
	English - UK
	<u>Celtic</u>
	East Asian
	Germanic
$\sim$	Romance
	Slavic
	Other Languages
	Bilinguals
	Clinical
Télécharger les trar	nscriptions de la base de données





### Index of /data/Romance







$\Theta \Theta \Theta$	Paris	C		
		Q		
Network     Name       Imain disk     Imain disk       Imain disk <t< th=""><th>&gt;</th><th>Date Modified May 3, 2010, 10:36 AM Yesterday, 10:05 AM Today, 2:38 PM Yesterday, 10:05 AM</th><th></th><th></th></t<>	>	Date Modified May 3, 2010, 10:36 AM Yesterday, 10:05 AM Today, 2:38 PM Yesterday, 10:05 AM		
Applications	000	📁 made	leine	$\bigcirc$
P Documents			Q	
Movies Music Pictures 4 items	Network main disk NidjhaWD Audio CD Desktop Cp	Name           MADELEINE-02-1_00_05.cha           MADELEINE-03-1_01_10.cha           MADELEINE-04-1_02_14.cha           MADELEINE-05-1_02_28.cha           MADELEINE-06-1_03_18.cha           MADELEINE-06-1_04_18.cha           MADELEINE-08-1_06_04.cha           MADELEINE-09-1_07_15.cha           MADELEINE-10-1_09_03.cha           MADELEINE-10-1_09_03.cha           MADELEINE-11-1_10_07.cha	<ul> <li>Date Modified</li> <li>Jun 5, 2010, 6:53 PM</li> <li>Jun 5, 2010, 6:53 PM</li> <li>May 30, 2010, 3:22 PM</li> <li>May 30, 2010, 3:22 PM</li> <li>Jun 5, 2010, 6:53 PM</li> </ul>	Size         Kind           100 KB         C           104 KB         C           64 KB         C           52 KB         C           104 KB         C           92 KB         C           100 KB         C           120 KB         C           116 KB         C           144 KB         C           100 KB         C
	Applications Documents Movies Music Pictures	MADELEINE-13-2_01_02.cha MADELEINE-14-2_02_06.cha MADELEINE-15-2_03_05.cha MADELEINE-16-2_04_15.cha MADELEINE-17-2_05_12.cha MADELEINE-18-2_06_10.cha MADELEINE-19-2_07_07.cha	Jun 5, 2010, 6:53 PM Jun 5, 2010, 6:53 PM	108 KB C 100 KB C 96 KB C 108 KB C 116 KB C 100 KB C 120 KB C

Télécharger les transcriptions de la base de données







Madeleine + Anaé bientôt disponibles jusqu'à 4 ans sur CHILDES Déjà disponibles sur <u>http://colaje.risc.cnrs.fr</u> :





#### Télécharger les transcriptions de la base de données





### CHILDES The CHILDES Database



The CHILDES database contains transcript and media data collected from conversations between young children and their playmates and caretakers. Conversations with older children and adults at available from <u>TalkBank</u>. All of the data is transcribed in CHAT and CA/CHAT formats. The use of all CHILDES and TalkBank data is governed by the <u>Gnu</u> <u>Public License (GPL)</u>. Please remember to read the <u>database manuals</u> and to follow the guidelines for <u>data-sharing</u>.

### **Downloading TalkBank Data**

- Downloadable Transcripts
- Downloadable Audio and Video

!! Mettre transcriptions et vidéos dans le même dossier



### Index of /media

Name	Last modified	Size	Descriptio Index of /med	ia/Romance		
Parent Directory		-				
Biling/	02-Jul-2009 12:52	-	Name	Last modified	Size	Description
Clinical/	07-Feb-2006 13:31	-				
EastAsian/	18-Sep-2010 13:10	-	Farent Directory		-	
Eng-UK/	30-Jan-2010 12:35	-	French/	05-Mar-2009 15:45	-	
Eng-USA/	19-Aug-2010 12:23	20	Italian/	19-Nov-2006 15:19	-	
Germanic/	02-Jun-2008 18:30	-	Portuguese/	20-May-2005 15:17	-	
Narrative/	20-May-2005 14:26	-	Spanish/	15-Sep-2008 17:33	-	
Other/	11-Nov-2010 16:09	_				
PDF/	08-Dec-2007 16:12	-				
Password/	30-Jan-2010 13:13	-	Index of /medi	a/Romance	/Fre	nch
PhonBank-wav/	16-Sep-2010 11:34	-	index of / mou			
PhonBunk/	16-Sep-2010 12:21	-				
Romance/	10-Oct-2010 15:01	_	Name	Last modified	Size	Description
Slavic/	11-Oct-2005 17:06	-	Parent Directory		-	
Unlinked/	10-Oct-2010 15:00	-	Lyon/	20-Sep-2007 19:05	$\simeq$	
			Paris/	05-Mar-2009 14:15	2	

#### Télécharger les média de la base de données

### Index of /media/Romance/French/Paris

Name	Last modified	Size	Description
Parent Directory		-	
Loonard/	12-Mar-2009 17:45	- 6	
Madeleine/	12-Mar-2009 17:45	- i	
Theophile/	12-Mar-200	ov o	f /media/Ro

!! Mettre transcriptions et vidéos dans le même dossier

Mar-200 Index of /media/Romance/French/Paris/Madeleine

	Name	Last modified	Size	Description
2	Parent Directory		-	
Ш	01.mov	05-Apr-2006 10:08	3 2.1G	
Ц	02.mov	21-Dec-2008 23:59	9 915M	
Щ	<u>03.mov</u>	22-Dec-2008 00:00	927M	
Ħ	04.mo	Window 08 00:00	5 804M	
H	05.mo Open Link in New	Tab 08 08:2	848M	
H	06.mo	08 00:11	915M	
H	07.mc Save Link As	08 00:07	7 742M	
H	08.mo Send Link	08 00:09	882M	
	09.mo	08 00:10	890M	
U	10.mo Properties	07 11:53	939M	

Télécharger les média de la base de données



000	MADELEINE	0
	<b>\$₹</b>	
Network	Name	A Date Modified
	MADELEINE-11-1_10_07.cha	Mar 26, 2010, 10:52
main disk	MADELEINE-11-1_10_07.mov	Feb 24, 2007, 8:29 PN
	圖 MADELEINE-12-1_11_13.cha	Dec 16, 2009, 3:47 PN
Desktop	MADELEINE 12 1_11_13.mov	Apr 8, 2007, 12:00 AM
1 co	MADELEINE-13-2_01_02.cha	Nov 17, 2010, 3:35 PM
	MADELEINE-13-2_01_02.mov	Jun 28, 2007, 12:45 P
Applications	副 MADELEINE-14-2_02_06.cha	Nov 6, 2009, 12:54 PM
P Documents	MADELEINE-14-2_02_06.mov	Jun 28, 2007, 10:37 A
bocuments	MADELEINE-15-2_03_05.cha	Oct 25, 2009, 9:11 PN
Movies	MADELEINE-15-2_03_05.mov	Sep 21, 2007, 12:04 F
la Music	MADELEINE-16-2_04_15.cha	Oct 28, 2010, 10:42 A
( Music	MADELEINE-16-2_04_15.mov	Sep 21, 2007, 2:29 PN 🔺
Pictures	IIII MADELEINE-17-2 05 12.cha	Dec 15. 2009. 9:49 Al
	1 of 66 selected, 40,76 GB available	11

!! Mettre transcriptions et vidéos dans le même dossier





Quand on ouvre un fichier CLAN

> En-têtes (headers)

/Users/cp/Desktop/Paris/madeleine/MADELEINE-13-2\_01\_02-@Begin @Languages: fra @Participants: CHI Madeleine Target Child, MOT Mother, OBS Martine Observer, UNI Unidentified @ID: fra|Paris-Corpus Madeleine|CHI|2;01.02|female|||Target Child|| @ID: fra|Paris-Corpus Madeleine|MOT||female|||Mother|| @ID: fra|Paris-Corpus\_Madeleine|OBS||female|||Observer|| @ID: fra|Paris-Corpus\_Madeleine|UNI||||Unidentified|| @Birth of CHI: 14-APR-2005 @Media: MADELEINE-13-2 01 02, video @Date: 16-MAY-2007 @Time Duration: 00:00:00-01:02:30 @Location: Madeleine's home @Comment: code"'r - Stéphanie Caët @G: dans la cuisine \*CHI: 0. %act: se dirige vers la table \* Lignes \*MOT: et ta table elle est propre ? • \*CHI: oui. • principales %pho: wi \*OBS: délicieux c(e) café (.) merci beaucoup ! • % Lignes %add: OBS s'adresse à MOT secondaires %sit: CHI nettoie sa table Balises de temps \*CHI: moi i(e) nettoie . • %pho: mwa z etwa \*CHI: moi je nettoie moi . • %pho: mwa 39 letwa mwa \*MOT: tu nettoies ? · 04nov10[E|CHAT] \* 13

29/79

CLAN File Edit Font Size/Style Tiers Mode Windows Help

Lire une transcription de la base de données



🧯 CLAN Fil	e Edit	Font	Size/Style	Tiers	Mode	Windows	Help	
$\Theta \Theta \Theta$			/Users	/cp/Des	ktop/Par	is/madeleine	/MADELEI	NE-13-2_01_02-6
@Begin								
@Languages: fra	а							
@Participants: CH	II Madelein	ne Targe	t_Child, MOT	Mother,	OBS Ma	artine Observe	er, UNI Uni	identified
@ID: fra Paris-Co	orpus_Mad	leleine C	HI 2;01.02 fe	male   Ta	arget_Ch	ild		
@ID: fra Paris-Co	orpus_Mad	leleine N	10T  female	Mother				
@ID: fra Paris-Co	orpus_Mad	leleine C	)BS  female	Observe	r			
@ID: fra Paris-Co	orpus_Mad	leleine L	JNI     Unident	ified				
@Birth of CHI: 14-	APR-2005	04.00						
@Media: MADELE	INE-13-2_0	01_02, v	laeo					
@Time Duration: 0	0.007	1.02.30						
@Location: Mar	teleine's h	ome						
@Comment: co	de"r - Sté	onie	Caët	_				
@G: dans la cuisin	e	Priorito		Date		مامم		lines de terre
*CHI: 0.	53			rou	ram	cherie	es pa	ilses de tem
%act: se dirige vers	s la table		(		M	lode>F	Txnar	nd bullets
*MOT: et ta table	elle est pro	opre ? •1	102		1.0			
*CHI: oui . •10303	38_10324	15•				O	u Esc	c+a
%pho: wi			_					
*OBS: délicieux c	(e) café (.)	merci t	ucoup ! •10	32415_1	1035696			
%add: OBS s'adres	sse à MOT							
%sit: CHI nettoie s	sa table							
*CHI: moi j(e) nett	toie . •1035	5697_1	895•					
%pho: mwa 3 etw	a Ia mai . 10	027005	1040277-					
% pho: mwa za lat	ie mor . •10	021992	1040377•					
*MOL: tu nettoies	2 •104037	7 1041	498.					
04poy10[EICHAT]	+ 40	1_1041	400-					
	- 13							

#### Lire une transcription de la base de données







Lire une transcription de la base de données


## Dans cette présentation ...

#### 1. Le Projet CHILDES

### 2. Le logiciel CLAN

1. Lire et analyser des transcriptions de la base de données CHILDES

#### 2. Transcrire

- 3. Aligner la transcription avec la vidéo ou le son
- 4. Faire des analyses automatiques sur les transcriptions



#### 2.2 Transcrire avec CLAN La feuille de texte

Ś	CLAN	File	Edit	Font	Size/Style	Tiers	Mode	Windows	Help	
00	0				newfile.c	ha				
1										C



Le logiciel CLAN



. . .



### 2.2 Transcrire avec CLAN Les en-têtes (méta-données)







\*CHI: maman ə@fs veux ça s'i(I)+te+plaît [=! pleure]. •

 Pour les interactions adulte-enfant : découpage en énoncés => MLU

Parfois problématique ...

\*CHI: moi je veux pas jouer +...

\*CHI: parce que moi je veux lire un livre.

OU

\*CHI: moi je veux pas jouer (.) parce que moi je veux lire un livre.



\*CHI: maman ə@fs veux ça s'i(I)+te+plaît [=! pleure]. •

- Pour les interactions adulte-enfant : découpage en énoncés
- Transcription en orthographe => analyses
   Problématique aussi …
   \*CHI: tombé / tombée / tombés / tomber.
   %pho: tõbe
   Mais parfois nécessaire …





\*CHI: tigre. %pho: ti:I**ʁ** 

\*CHI: tigre. %pho: ti:g

\*CHI: un tigre. %pho: ε̃ ti:lg

\*CHI: faut trouver du tigre. %pho: fo kruve dy ti:ʁ



\*CHI: maman ə@fs veux ça s'i(I)+te+plaît [=! pleure]. •

- Pour les interactions adulte-enfant : découpage en énoncés
- Transcription en orthographe
- Transcription adaptée à l'oral et à CLAN

   (.) pauses
   (élisions)
   mots+composés
   +< chevauchements</li>



\*CHI: maman ə@fs veux ça s'i(I)+te+plaît [=! pleure]. •

- Pour les interactions adulte-enfant : découpage en énoncés
- Transcription en orthographe
- Transcription adaptée à l'oral et à CLAN
- Adaptations au langage de l'enfant 
   @@fs veux un tiktik@c



#### 2.2 Transcrire au format CHAT Les lignes secondaires : %cod

\*CHI: maman ə@fs veux ça s'i(I)+te+plaît [=! pleure]. •

%pho: mamã ə vø sa sitəplε %act: se dirige vers l'étagère %xpnt: CHI pointe du chocolat %sit: MOT le suit Transcription phonétique Transcription du non verbal

Lignes secondaires liées à la ligne principale

Lignes secondaires adaptables selon sujet de recherche





## Dans cette présentation ...

#### 1. Le Projet CHILDES

### 2. Le logiciel CLAN

- 1. Lire et analyser des transcriptions de la base de données CHILDES
- 2. Transcrire
- 3. Aligner la transcription avec la vidéo ou le son
- 4. Faire des analyses automatiques sur les transcriptions



- 1. Ouvrir CLAN. Le fenêtre de commandes s'ouvre (sinon : ctrl/ **é**+d)
- 2. Dans la fenêtre de commandes, spécifier le dossier dans lequel se trouve la vidéo à transcrire et spécifier la grammaire

00	Commands	
working /Use	rs/cp/Desktop/Paris/made	leine
output		
lib /App	lications/CLAN/fr	
mor lib /App	lications/CLAN/fr	
GLAN		Help
First-		
Recall 04no	v10	Run

3. Fermer la fenêtre de commandes

Le logiciel CLAN





- 4. Ouvrir un nouveau document (File>New)
- 5. L'enregistrer aussitôt dans le dossier où se trouve la vidéo (Save as)
- 6. Remplir les en-têtes





- Lancer la vidéo et le processus d'alignement Mode>transcribe sound or movie (ou F5)
- 8. Appuyer sur Espace pour découper en énoncés et insérer les balises de temps







- 9. Réécouter chaque balise (Mode>Play bullet Media ou F4 ou ctrl/ c+clic)
- 10. Insérer le nom du locuteur (ctrl+1 ou ctrl+2 ...)
- 11. Et transcrire les énoncés juste devant les balises et le non verbal sous les balises



43/79



12. Pour réécouter l'ensemble : Mode>Continuous Playback ou (Esc+8)13. That's all !





## Dans cette présentation ...

#### 1. Le Projet CHILDES

### 2. Le logiciel CLAN

- 1. Lire et analyser des transcriptions de la base de données CHILDES
- 2. Transcrire
- 3. Aligner la transcription avec la vidéo ou le son
- 4. Faire des analyses automatiques sur les transcriptions



#### 2.3 Analyses automatiques Les commandes de base

Nombre total d'énoncés Nombre total de morphèmes ou de mots MLU : Longueur Moyenne des énoncés



MLU

Nombre total de mots et de mots différents Type/Token ratio Fréquences d'items



Combinaisons de critères de recherche Fréquence des combinaisons



Analyse morphosyntaxique désambiguisée







#### 2.3 Analyses automatiques La fenêtre de commandes

Ouvrir la fenêtre de commandes : Windows > Commands Ou ctrl/ **#**+d



Le logiciel CLAN

46/79



#### 2.3 Analyses automatiques Structure d'une commande simple







#### 2.3 Analyses automatiques MLU : longueur moyenne des énoncés

000	CLAN Output	t	
> MLU +t*CHI -t%mor -syy -sy mlu +t*CHI -t%mor -syy -sxx Mon Nov 15 11:29:08 2010	∝	ppel de la commande	0
mlu (04-Nov-2010) is conduct ONLY speaker main tiers ma	ting analyses on: tching: *CHI;	appel du nom de fichier	
From file MLU for Speaker: *CHI	/Paris/madeleine/MAD	ELEINE-13-2_01_02.cha>	
MLU (xxx and yyy are EXCLU	DED from the utterance	e and morpheme counts):	
Number of: utterances =	526, morphemes = 15	511	
Ratio of morphemes over	r utterances = $2.873$	MLU	
Standard deviation = 1.	574		



#### 2.3 Analyses automatiques MLU : longueur moyenne des énoncés

#### Après report des résultats dans Excel ...



Le logiciel CLAN

50/79



#### 2.3 Analyses automatiques MLU : longueur moyenne des énoncés

#### Permet de comparer avec d'autres enfants ...



Le logiciel CLAN



#### **2.3 Analyses automatiques** FREQ 1 : liste des mots produits par l'enfant



Le logiciel CLAN

52/79



#### **2.3 Analyses automatiques** FREQ 1 : liste des mots produits par l'enfant

<pre>&gt; freq +t*CHI +o @ freq +t*CHI +o @</pre> Rappe	l de la commande
freq (04-Nov-2010) is conducting analyses on:	appel du nom de fichier
UNLT speaker main tiers matching: "CHI;	
From file <td>ADELEINE-13-2_01_02-mor.cha&gt;</td>	ADELEINE-13-2_01_02-mor.cha>
85 c'est	>
56 oui	
53 ça	
43 Veux	
36 non	
32 le	
29 la	
27 on	
27 pas	
T Voudrais	
331 Total number of different word types used	
1745 Total number of words (tokens)	
0.190 Type/Token ratio	

Le logiciel CLAN



	000	CLAN Output	
	> FREQ +t*CHI +s"je" MADELEINE*.cha		<u> </u>
The !	freq +t*CHI +sje MADELEINE*.cha		
	Mon Nov 15 11:41:04 2010		
112-	freq (04-Nov-2010) is conducting analy	ses on:	
	ONLY speaker main tiers matching: *CH	-11;	
	*******	**	
ſ	From file <madeleine-11-1_10_07-cor.< td=""><td>.cha&gt;</td><td></td></madeleine-11-1_10_07-cor.<>	.cha>	
	1 je		
Ź			
	1 Total number of different word typ	bes used	
	1 Total number of words (tokens)		
C	1.000 Type/Token ratio		
	From file <madeleine-12-1_11_13.cna< td=""><td>&gt;</td><td></td></madeleine-12-1_11_13.cna<>	>	
	0 Total number of different word typ	nes used	
	0 Total number of words (tokens)		
	From file <madeleine-13-2 01="" 02-bis.<="" td=""><td>.cha&gt;</td><td></td></madeleine-13-2>	.cha>	
	0 Total number of different word typ	bes used	
	0 Total number of words (tokens)		
	From file <madeleine-13-2_01_02-mon< td=""><td>r-pst.cha&gt;</td><td></td></madeleine-13-2_01_02-mon<>	r-pst.cha>	
	19 je		
$\boldsymbol{\boldsymbol{1}}$		>	
	1 Total number of different word typ	bes used	
	19 Total number of words (tokens)	J	55/79

.....



Calcule les « je » et les « tu »

Le logiciel CLAN



#### 2.3 Analyses automatiques FREQ 2 : fréquences d'items

000	CLAN Output				
> freq +t*CHI +sje +stu @ freq +t*CHI +sje +stu @ Wed Nov 17 12:29:18 2010 freq (04-Nov-2010) is conducting analyses on: ONLY speaker main tiers matching: *CHI;		Madeleine			
From file 19 je 6 tu	MADELEINE-13-2_01_02.cha>	la Tiarr Mada Windows Holn			
2. Total number of different word tensor word	CLAN File Edit Font Size/Sty	put			
25 Total number of words (tokens) 0.080 Type/Token ratio	> freq +t*MOT +sje +stu @ freq +t*MOT +sje +stu @ Wed Nov 17 17:05:04 2010 freq (04-Nov-2010) is conducting analyses on: ONLY speaker main tiers matching: *MOT;				
From file 49 je 157 tu					
La mère de Madeleine <sup>2</sup> Total number of different word types used <sup>2</sup> Total number of words (tokens) 0.010 Type/Token ratio					





#### 2.3 Analyses automatiques FREQ 2 : fréquence d'items

#### Après report des résultats dans Excel ...



Le logiciel CLAN



#### 2.3 Analyses automatiques FREQ 2 : fréquence d'items

#### Permet de comparer la mère et l'enfant ...





#### 2.3 Analyses automatiques COMBO 1 : contexte de production





#### 2.3 Analyses automatiques COMBO 1 : contexte de production





Le logiciel CLAN



#### 2.3 Analyses automatiques COMBO 2 : recherche double



#### **CLAN** Output

\*\*\* File "/Users/cp/Documents/corpus/MADELEINE/MADELEINE-13-2\_01\_02.cha": line 1989.

\*CHI: voilà [>] . •

\*MOT: <il en a bien> [<] besoin .

\*CHI: <(1)je v(eux)> [/] (2)je veux pas le chandail . •

\*CHI: maman prend le chandail . •

\*MOT: bah oui ben j(e) vais le mettre à la poupée . •

\*\*\* File "/Users/cp/Documents/corpus/MADELEINE/MADELEINE-13-2\_01\_02.cha": line 2300.

\*MOT: le col c'est important qu' i(l) soit bien r(e)passé . •

\*MOT: 0 [=! imite bruit de la vapeur] . •

\*CHI: (1)t(u) as mis des [/] des +/. •

\*MOT: voilà ça y+est . •

\*CHI: t(u) as fini ? •

Triple clic sur nom de fichier : Accès au passage dans la transcription

Le logiciel CLAN

63/79



Le logiciel CLAN

64/79


#### 2.3 Analyses automatiques COMBO 3 : collocations

$\Theta \Theta$	CLAN Output	
> Combo +t*CHI +s"moi^je" @		
(moi^je))		
combo +t*CHI +smoi^je @		
Wed Nov 17 12:45:32 2010		
combo (04-Nov-2010) is conducting analyses on:		
ONLY speaker main tiers matching: *CHI;		
From file <td>ELEINE-13-2_01_02.cha&gt;</td> <td></td>	ELEINE-13-2_01_02.cha>	
*** File "/Users/co/Documents/cornus/MADELEINE/MADE	EINE 13.2 01 02 cha": line 275	
CHI: <(1)moi (1)ie your enlayer ton manteau>[>] •	LEHNE-13-2_01_02.cna . me 275.	
ern. (f)nor (f)je veux entever ton manteau> [>].		
*** File "/Users/cp/Documents/cornus/MADELEINE/MADE!	I FINE-13-2 01 02 cha": line 704	
$CHI: (D)moi \le (D)ie veux ouv(rir) > [//] zo@ifs veux fermer v$	N+ •	
erne (rynor (rynor (rynor our our fin j zo@is feux feiner y		
*** File "/Users/cp/Documents/corpus/MADELEINE/MADE!	LEINE-13-2 01 02 cha": line 793	
CHI: (Dmoi (Die veux nettover •		
*** File "/Users/cp/Documents/corpus/MADELEINE/MADE	LEINE-13-2 01 02 cha": line 821	
CHI: (Dmoi (Dife) nettoie •		
*** File "/Users/cp/Documents/corpus/MADELEINE/MADE	I FINE-13-2 01 02 cha": line 823	
CUL (De si (D) sottois (De si	LEINE-15-2_01_02.000 . Inte 025.	



sur les énoncés de l'enfant uniquement



## 2.3 Analyses automatiques MOR : analyse morphosyntaxique

$\Theta \Theta \Theta$	CLAN Output				
> MOR +t*CHI MADELE	INE-13-2_01_02.cha		0		
mor +t*CHI MADELEINE	E-13-2_01_02.cha				
Mon Nov 22 16:02:31 20	010				
mor (04-Nov-2010) is co	inducting analyses on:				
ONLY speaker main tie	ers matching: *CHI;				
and those speakers' A	ALL dependent tiers				
and ALL header tiers					
**** <mark>*********************</mark> ****	*******				
From file <madeleine< td=""><td>-13-2_01_02.cha&gt; to file <madeleine-13< td=""><td>-2_01_02.mor.cex&gt;</td><td></td><td></td><td></td></madeleine-13<></td></madeleine<>	-13-2_01_02.cha> to file <madeleine-13< td=""><td>-2_01_02.mor.cex&gt;</td><td></td><td></td><td></td></madeleine-13<>	-2_01_02.mor.cex>			
Using sf-rule: /Applicatio	ns/CLAN/fr/sf.cut.				
Using a-rules: /Applicati	ons/CLAN/fr/ar.cut.				
Using lexicon: /Applicati	ons/CLAN/fr/lex/0affix-n.cut.				
Using lexicon: /Applicati	ons/CLAN/fr/lex/0affix-v-irreg.cut.				
Using lexicon: /Applicati	ons/CLAN/fr/lex/0affix-v.cut.				
[		000	MA	DELEINE	_
Using lexicon: /Applicati	ons/CLAN/fr/lex/zero.cut.		m	0	
Loaded lexicon: 44572				ų	
Using c-rules: /Application	ons/CLAN/fr/cr.cut.	Network	Name	2-1 11 13.01a	▲ Date
warning: crule "adv-me	nt" referenced but not defined	📃 main disk	MADELEINE-1	2-1_11_13.mov	Ap
warning: crule "adv-me	nt" referenced but not defined	Deckton	MADELEINE-1	3-2_01_02.cha	No
warning: crule "adv-me	nt" referenced but not defined	A cn	MADELEINE-1	3-2_01_02.mor.cex	То
warning: crule "adv-me	nt" referenced but not defined	1 CP	MADELEINE-1	3-2_01_02.mov	Jur
2600		B Documents	MADELEINE-1	4-2_02_06.cha	No
Done with file <madele< td=""><td>EINE-13-2_01_02.mor.cex&gt;</td><td>Movies</td><td>MADELEINE-1</td><td>4-2_02_06.mov</td><td>Jur</td></madele<>	EINE-13-2_01_02.mor.cex>	Movies	MADELEINE-1	4-2_02_06.mov	Jur
		& Music		5-2_03_05.cna	Oc A
		Pictures		5-2_05_05.mov	Sel
		in recures			) + + (
			67 items 40 3	A CR available	1





# 2.3 Analyses automatiques MOR : analyse morphosyntaxique

#### Fichier .mor.cha créé par CLAN :

*CHI: 0.	0
%act: CHI jette une regard à OBS puis reprend son ménage .	
*CHI: moi l@fs a yy (.) nettoyé . •	
%mor: pro moi&SING^n moi&_MASC fs I	
v:poss avoir&PRES&3SV^v:aux avoir&PRES&3SV unk yy	
vinettoyer-PPMASC .	
%pho: mwa I a tõ letaje	
*OBS: tu te vois dedans . •	
04nov10[E CHAT] * 1150	
	<b>A</b>
• = OU	
(" a " act la varba passacif avair (vup	acalayair®DDES92SV/

**OU** l'auxiliaire avoir (v:aux|avoir&PRES&3SV)







# 2.3 Analyses automatiques POST : désambiguïse l'analyse de MOR

000	CLAN Output			
> POST MADELEINE-13-	2_01_02.mor.cex			
Using file: /Applications post MADELEINE-13-2_0 Mon Nov 22 16:03:49 201 post (04-Nov-2010) is con ALL speaker tiers	s/CLAN/fr/post.db. 1_02.mor.cex 0 ducting analyses on:			
From file <madeleine-1 Done with file <madelei< td=""><td>3-2_01_02.mor.cex&gt; to file <madel NE-13-2_01_02.mor.pst.cex&gt;</madel </td><td>EINE-13-2_01_02.m</td><td>nor.pst.cex&gt;</td><td></td></madelei<></madeleine-1 	3-2_01_02.mor.cex> to file <madel NE-13-2_01_02.mor.pst.cex&gt;</madel 	EINE-13-2_01_02.m	nor.pst.cex>	
04nov10[E TEXT] 12				C
		<ul> <li>Network</li> <li>main bisk</li> <li>Desktop</li> <li>cp</li> <li>Applications</li> <li>Documents</li> <li>Movies</li> <li>Music</li> <li>Pictures</li> </ul>	Name           MADELEINE-12-1_11_13.cma           MADELEINE-12-1_11_13.mov           MADELEINE-13-2_01_02.cha           MADELEINE-13-2_01_02.mor.cex           MADELEINE-13-2_01_02.mor.pst.cex           MADELEINE-13-2_01_02.mov           MADELEINE-13-2_01_02.mov           MADELEINE-13-2_01_02.mov           MADELEINE-14-2_02_06.cha           MADELEINE-14-2_02_06.mov           MADELEINE-15-2_03_05.cha           MADELEINE-15-2_03_05.cha	Ap Ap No To Jur No Jur Oc
			68 items, 40.34 GB available	







# 2.3 Analyses automatiques POST : désambiguïse l'analyse de MOR

#### Fichier .mor.pst.cex créé par CLAN :



71/79





#### 2.3 Analyses automatiques MOR|POST et FREQ



73/79





#### 2.3 Analyses automatiques MOR|POST et COMBO

 $\Theta \Theta \Theta$ **CLAN Output** %mor: (1)pro:subjlil vltomber . \*\*\* File "MADELEINE-13-2\_01\_02.mor.pst.cex": line 381. \*CHI: (.) <i(I) tombe> [/] <i(I) tombe> [<] . • %mor: (1)pro:subjlil vltomber. \*\*\* File "MADELEINE-13-2\_01\_02.mor.pst.cex": line 711. \*CHI: on va faire un café . • %mor: (1)pro:subj|on v:mdl|aller v:mdllex|faire-INF det|un n|café . \* File "MADELEINE-13-2 01 02.mor.pst.cex": line 923. \*CHI: moi je veux nettoyer . • %mor: pro/moi (1)pro:subjlje v:mdl/vouloir v/nettoyer-INF \_\_\_\_\_ \*\*\* File "MADELEINE-13-2\_01\_02.mor.pst.cex": line 958. \*CHI: moi j(e) nettoie . • %mor: pro/moi (1)pro:subj/je v/nettoyer . \_\_\_\_\_ \*\*\* File "MADELEINE-13-2\_01\_02.mor.pst.cex": line 961. \*CHI: moi je nettoie moi . • %mor: pro|moi (1)pro:subj|je v|nettoyer pro|moi . 04nov10[E|TEXT] 38





#### 2.3 Analyses automatiques et beaucoup d'autres commandes encore ...

Command	Page	Function
CHAINS	50	Tracks sequences of interactional codes across speakers.
CHECK	53	Verifies the correct use of CHAT format.
CHIP	56	Examines parent-child repetition and expansion.
COMBO	63	Searches for complex string patterns.
COOCUR	71	Examines patterns of co-occurence between words.
DIST	71	Examines patterns of separation between speech act codes.
DSS	72	Computes the Developmental Sentence Score.
FREQ	81	Computes the frequencies of the words in a file or files.
FREQMERG	90	Combines the outputs of various runs of FREQ.
FREQPOS	90	Tracks the frequencies in various utterance positions.
GEM	91	Finds areas of text that were marked with GEM markers.
GEMFREQ	94	Computes frequencies for words inside GEM markers.
GEMLIST	94	Lists the pattern of GEM markers in a file or files.
KEYMAP	95	Lists the frequencies of codes that follow a target code.
KWAL	96	Searches for word patterns and prints the line.
MAXWD	98	Finds the longest words in a file.
MLT	100	Computes the mean length of turn.
MLU	103	Computes the mean length of utterance.
MODREP	108	Matches the child's phonology to the parental model.
PHONFREQ	111	Computes the frequency of phonemes in various positions.
RELY	112	Measures reliability across two transcriptions.
STATFREQ	113	Formats the output of FREQ for statistical analysis.
TIMEDUR	114	Uses the numbers in sonic bullets to compute overlaps.
VOCD	115	Computes the VOCD lexical diversity measure.
WDLEN	120	Computes the length of utterances in words.



### Pour aller plus loin avec CLAN ...

#### CHILDES Child Language Data Exchange System



CHILDES is the child language component of the <u>TalkBank</u> system. TalkBank is a system for sharing and studying conversational interactions.

System

Ground rules

**Guidelines for Contributors** 

**Overviews and Introductions** 

Membership list

How to subscribe to Mailing Lists

Links

The TalkBank Database

**IASCL** information

Other Child Language sites

#### **Programs and Database**

Downloadable Database

Browsable Database

The CLAN Program

Derived Corpora and Counts

Manuals

CHAT Transcription

**CLAN Programs** 

**Database Manuals** 

Workshop Novembre 2010

77/79



#### Pour aller plus loin avec CLAN et d'autres logiciels...



#### Workshop Novembre 2010

78/79



## Pour aller plus loin avec CLAN et d'autres logiciels...



Workshop Novembre 2010









Cette présentation a été créée à partir du matériel disponible sur le site Internet de CHILDES (http://childes.psy.cmu.edu/), le manuel CLAN de L. Balthasar (Laboratoire ICAR), le manuel CLAN du Projet LEONARD dirigé par A. Morgenstern, et constamment enrichie grâce aux questions des étudiants !



stephanie.caet@gmail.com



Outils et Recherches pour les Corpus d'Acquisition