

Acquisition Sketch Project Meeting 5:

Data Processing Part 1 – Setup, Transcribers, Segmenting, Tiers

6th March 2024

Agenda:

- Selecting a setup
- Selecting and training transcribers
- Segmenting the data
- Tiers for transcribing and annotation

Selecting a setup

Two most common options

ELAN plus Toolbox or FieldWorks

- Typically used in language documentation
- Shows sound waves and video on screen
- Tiers can be hierarchically connected
- Best option for annotating visual information – e.g. gesture
- <https://archive.mpi.nl/tla/elan>

CHAT transcription system and CLAN analysis tools in CHILDES

- Typically used in child language acquisition
- Shows sound waves and video on screen
- Tiers not hierarchically connected
- Relatively easy to learn and use
- Dedicated set of analysis tools specific to child language research (CLAN)
- <https://chilDES.talkbank.org>

ELAN

ELAN is an annotation tool for audio and video recordings.

Screenshot 1

Nr	Annotation	Begin TL...	End Time	Duration
6	oh	00:00:08	00:00:09	00:00:00
7	a square	00:00:09	00:00:10	00:00:00
8	oh	00:00:13	00:00:13	00:00:00
9	what's that ?	00:00:15	00:00:15	00:00:00
10	the blue cross	00:00:15	00:00:16	00:00:00
11	and the ... [+RQ]	00:00:17	00:00:17	00:00:00
12	oh dear	00:00:17	00:00:18	00:00:00
13	where's the basket gone ?	00:00:20	00:00:21	00:00:00
14	where's it gone ?	00:00:22	00:00:22	00:00:00

A sample from the [ACLEW project](#).

Screenshot 2

Nr	Annotation	Begin TL...	End Time	Duration
1	[a] guess!	00:00:00	00:00:01	00:00:00
2	oh I do [and] know what I can say ?	00:00:01	00:00:02	00:00:00
3	I don't think you've seen it yet.	00:00:03	00:00:04	00:00:00
4	it's in that bag.	00:00:05	00:00:06	00:00:00
5	it's in the bag.	00:00:08	00:00:09	00:00:00
6	Lego	00:00:12	00:00:13	00:00:01
7	have you got it ?	00:00:14	00:00:14	00:00:00
8	Lego	00:00:18	00:00:18	00:00:00
9	yeah [I] you build [] a big massive tower ?	00:00:28	00:00:28	00:00:00
10	you can't ?	00:00:30	00:00:30	00:00:00
11	why ?	00:00:30	00:00:31	00:00:00
12	what [] what if you [do you] mean there's no.	00:00:34	00:00:34	00:00:00

A sample from the [ACLEW project](#).

Description: With **ELAN** a user can add an unlimited number of textual annotations to audio and/or video recordings. An annotation can be a sentence, word or gloss, a comment, translation or a description of any feature observed in the media. Annotations can be created on multiple layers, called *tiers*. Tiers can be hierarchically interconnected. An annotation can either be time-aligned to the media or it can refer to other existing annotations. The content of annotations consists of Unicode text and annotation documents are stored in an XML format (EAF).

Main features and characteristics of ELAN:

- provides several ways to view the annotations, each view is connected to and synchronized with the media timeline
- supports creation of multiple tiers and tier hierarchies
- supports Controlled Vocabularies
- allows linking of up to 4 video files with an annotation document
- media support
 - builds on existing, native media frameworks, like Windows Media Player, QuickTime or VLC
 - support for audio and video formats depends on operating system, high performance media playback can usually be achieved
- technical
 - written in the Java programming language
 - distributions available for Windows, macOS and Linux
 - open source, the sources are available under a GPL 3 license



ELAN

ELAN 6.0 - LAMP_20130502_WF_01.eaf

File Edit Annotation Tier Type Search View Options Window Help

po@Acacia

Nr	Annotation	Begin Time	End Time	Duration
117	***	00:28:10.246	00:28:11.496	00:00:01.250
118	***	00:28:15.625	00:28:16.750	00:00:01.125
119	* ngarra nhinhi	00:29:41.444	00:29:43.762	00:00:02.318
120	***	00:29:57.817	00:29:59.885	00:00:02.068
121	***	00:30:09.454	00:30:10.033	00:00:00.579
122	kura nangkal nukun *	00:30:19.034	00:30:21.216	00:00:02.182
123	kura ngay	00:30:23.102	00:30:24.238	00:00:01.136
124	nart ne	00:30:25.897	00:30:26.829	00:00:00.932

00:13:35.622 Selection: 00:13:32.796 - 00:13:34.489 1693

Selection Mode Loop Mode

LAMP_201... 00:13:25.000 00:13:26.000 00:13:27.000 00:13:28.000 00:13:29.000 00:13:30.000 00:13:31.000 00:13:32.000 00:13:33.000 00:13:34.000 00:13:35.000 00:13:36.000

wakay mananthi warda
that's it, we've got nothing left!

po@Emily [466]
tf@Emily [440]
tn@Emily [60]
po@Acacia [361]
tf@Acacia [332]
tn@Acacia [37]
po@Tania [834]
tf@Tania [816]
tn@Tania [551]
notes [63]

mamba nhinhi w | i i | dhu-wa kanamkaykay | thangkunu thanamka
ok now your turn | There you go, he's calling out | say to him 'what are

bird call

Murrinhpatha, data from LAMP project

CHAT

CHILDES



Child Language Data
Exchange System

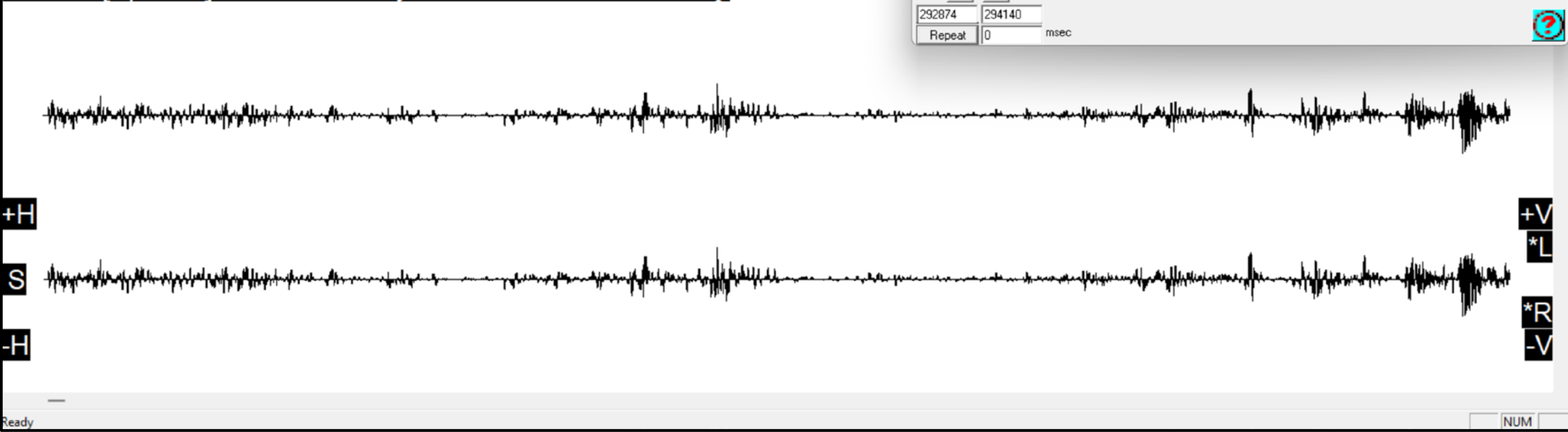
CHILDES is the child language component of the TalkBank system.

System	Database	Programs
<p><u>**Ground Rules**</u></p> <p>Contributing New Data</p> <p>IRB Principles</p>	<p><u>**Index to Corpora**</u></p> <p>Browsable Database</p> <p>TalkBankDB database search</p> <p>Hints on Downloading</p>	<p>CLAN</p> <p>XML creator and XML Schema</p> <p>Related Software</p>
Links	Teaching	Manuals
<p>Other Child Language sites</p> <p>Research based on CHILDES</p> <p>Child Language Diaries Data</p>	<p>Topics in Language Acquisition</p> <p>Teaching Resources</p> <p>YouTube Examples</p> <p>Bibliographies</p>	<p>CHAT Transcription Manual</p> <p>CLAN Program Manual</p> <p>Tutorial Screencasts</p> <p>Overviews, Other Languages</p>
Contact	Phonology and Fonts	Morphsyntax
<p>Brian MacWhinney : homepage</p> <p>How to subscribe to Mailing Lists</p>	<p>Phon and PhonBank</p> <p>Unicode and IPA for Mac</p> <p>Unicode and IPA for Windows</p>	<p>Universal Dependencies</p> <p>MOR manual</p>
Media, CA	Resources	Versions
<p>CA analysis</p> <p>Digitized video</p> <p>Digitized audio</p>	<p>Building a New Corpus</p> <p>CCT Computerized Comprehension</p> <p>LEAT Assessment Tool</p>	<p>Derived Corpora and Counts</p> <p>XML version of the database</p> <p>Database Versioning</p>

CHAT

```
File Edit View Tiers Mode Window Help
[Icons]
311 *FAT: i-lla-gi-s ulgwmaḍ ? •
312 %mor: 3SM#v|lla&PERF~prep|gi~pro:obj|s&3S n|ulgwmaḍ ?
313 %pho: il:agis ulgwmaḍɿ
314 *FAT: managu ad-t t-zra inna-m ? •
315 %mor: pro:int|managu comp|ad~pro:obj|t&3SM 3SF#v|zra&PERF
316 pro:poss:kin|inna-2SF?
317 %pho: manuk at:i tɿzra in:am
318 *CHI: nkki zri-y-t . •
319 %mor: pro:per|nkki&1S v|zri&PERF-1S~pro:obj|t&3SM.
320 %pho: nk:i zlixt
321 %mod: nk:i zɿrixt
322 *CHI: s-all-n-inu i-lla ammas n-yi-nn i-lla-sul wayya ammas n-yi-nn
```

26oct23[E|CHAT] 311 : Press any KEY or click mouse to stop



Paused 02:45

165254 Save

292874 294140

Repeat 0 msec

Berber, data from Abdellah Elouatiq

Helpful questions to ask

What setup do you already know?

What setup do you have support for (e.g. colleague who can help you)?

What setup is common in your subfield?

Is annotating video important for you (e.g. gesture)?

Is annotating phonetics important for you?

Do you need/want hierarchically connected tiers?

What programs/procedures will you use to analyze the data?

Where do you want to archive the data and what requirements do they have?

Selecting and training transcribers

Selecting transcribers

Typical choices

- You
- Caregivers or family members (alone or with you)
- Other community members (alone or with you)
- Maybe get help from older children and/or target child

Typical criteria

- Fluent speaker of the language
- Able to understand the child's speech
- Knowledge of the child's daily context
- Able to write the language – ideally using a computer
- Sufficient time and patience to complete the job (could be divided in parts)
- Acceptable to the family (ethical considerations)

Examples of transcribers

Inuktitut

- Students from the community studying in Montreal, working at researcher's office

Cree

- Mother of one target child who served as interlocutor for all recordings
- IPA transcription – students at the university trained in Cree phonology

Mongolian

- Researcher and Master's student in psychology who wanted to gain experience

Sesotho

- Mother of each child working together with the researcher

Chintang

- Community members who had studied linguistics

Supporting transcribers

Provide guidelines

- What to transcribe – child speech, caregiver speech, others?
- How to divide intonation units, utterances, etc.
- What to do with repetition, uninterpretable speech, etc.
- Transcribe what child **actually** says, not what you think they intended (that goes on another tier)
- Emphasize need for detail, quality, consistency, etc. as much as is reasonable

Amount of support

- Range from sitting together to fully independent work – depends on your situation
- Brainstorm ideas for motivation

Building in “quality control” if reasonable

Why?

- Increase transcriber confidence and knowledge
- Journals may request evidence of transcriber “reliability”

Hold training sessions

- All transcribers transcribe same short portion and compare (e.g. 2 minutes)
- Repeat a few times till cross-transcriber consistency is high
- Discuss and find solutions for differences
- Add solutions to guidelines document

Make plan to calculate transcriber agreement

- Pairs of transcribers blind-transcribe 10% of data
- Calculate Cohen’s D or % agreement

Segmenting the data

Select portions to transcribe

Scan each recording to select the portion to transcribe

Continuous chunk of 30 minutes

Maximum amount of child and caregiver speech

Avoid excessive crying, adult-only talk, distracting background noise, etc.

Need to listen to / watch recording – finding dense parts in sound file can be misleading

Prepare selected portion for transcription

Prepare audio (or video) files for your selected setup

Integrate audio (or video) file into the file for transcription

Prepare templates of tiers or have workflow in place to automatically generate them

If the transcriber is computer-savvy and has lots of time:

- Start transcribing, segmenting into units as you go

If the transcriber is less computer-savvy and/or has limited time:

- Segment small portions of speech within your setup so the transcriber can just go to each portion directly and not spend time finding the parts to transcribe

Figure 2a. Initial segmentation in the Qaqet corpus (for first-pass transcription).

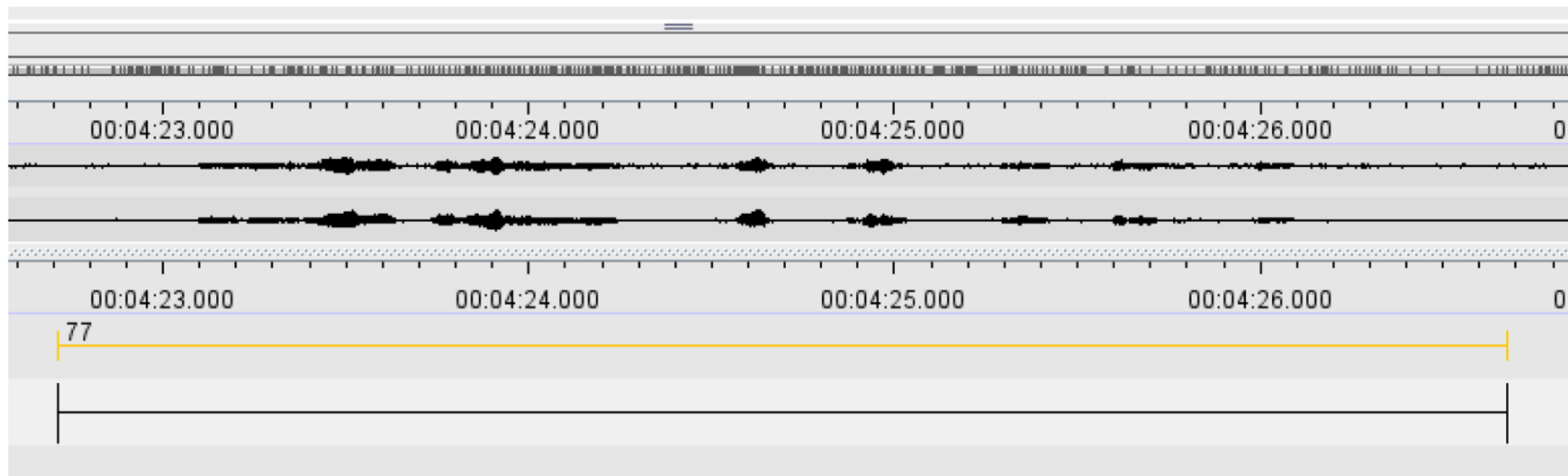
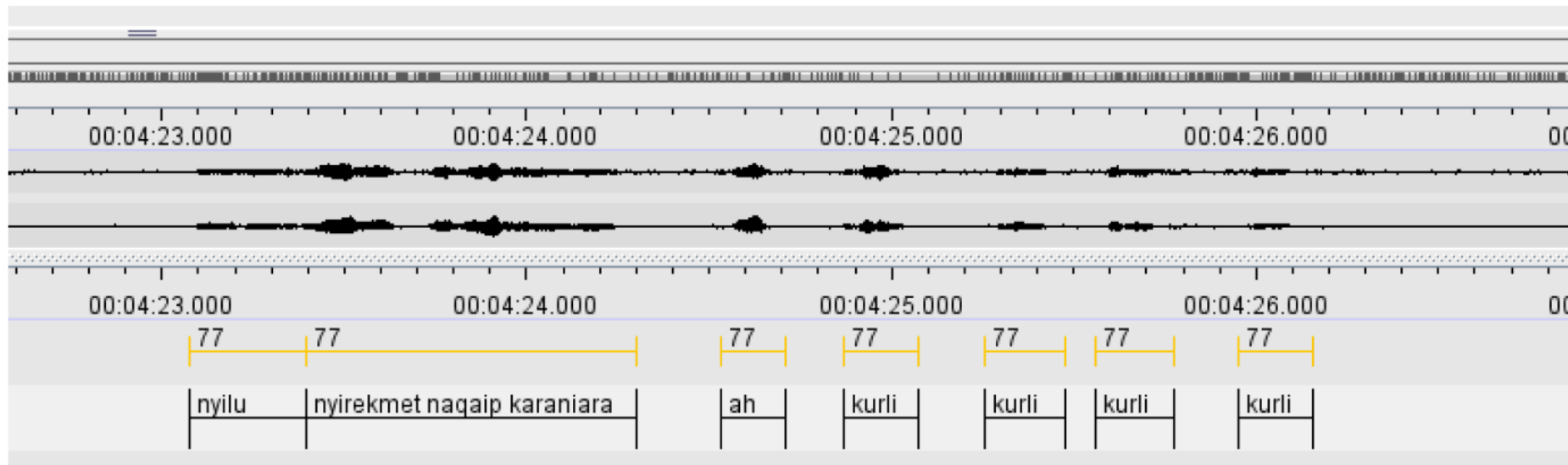


Figure 2b. Final segmentation in the Qaqet corpus (intonation units).



Select unit of segmentation to facilitate analysis

Intonation unit (typical in language documentation – Himmelmann 2006b)

- Segment of speech that occurs with a single prosodic contour
- Likely to begin with a brief pause and to end in a clause-final intonation contour
- Not necessarily identical with a syntactic unit, but often corresponds to a clause or a group of closely related clauses
- Can be as small as a single word or as long as a full sentence or more

Utterance (typical in child language – MacWhinney 2021)

- Continuous piece of speech, by one person, before or after which there is silence on the part of the person
- Typically a syntactic unit
- Can be as small as a single word or as long as a full sentence
- See CHAT Manual, chapter 9

Helpful tips

Child language is often repetitive

- Plan how to deal with repetitions – e.g. transcribe as one unit, decide depending on intonation

Child language is often hard to interpret

- Use video context for clues
- Listen to preceding and following utterances for clues
- Plan strategies to reduce frustration – e.g. maximum number of listens or amount of time before moving on, symbol/code to mark uninterpretable utterance or utterance to come back to later

If transcription resources are limited

- Think carefully about ways to reduce time and cognitive load – e.g. pre-segmentation, "removing" repetitions from first-pass transcription, providing motivation in procedure / setting / etc.

Tiers for transcribing and annotation

Potential tiers

Minimum

- Transcription in orthography
- Translation into national language, English, etc.
- Morphemic breakdown, gloss, part of speech

Ideas for optional tiers

- Interpretation of child utterance / likely target utterance
- Errors / child-like productions / differences from likely target
- Addressee
- Notes / comments / situational context
- Phonetic transcription
- Annotation for structures of interest

Helpful tips

Create a tier “template” for each utterance (cut/paste or automatically generated)

First pass = just the orthographic transcription and translation

- this will likely take at least 1 hour per 5 minutes of recording
- maybe also include interpretation of child utterance, addressee, and notes

Gradually add information for other tiers as you need them

Only use the tiers/information you need for your own questions and analyses

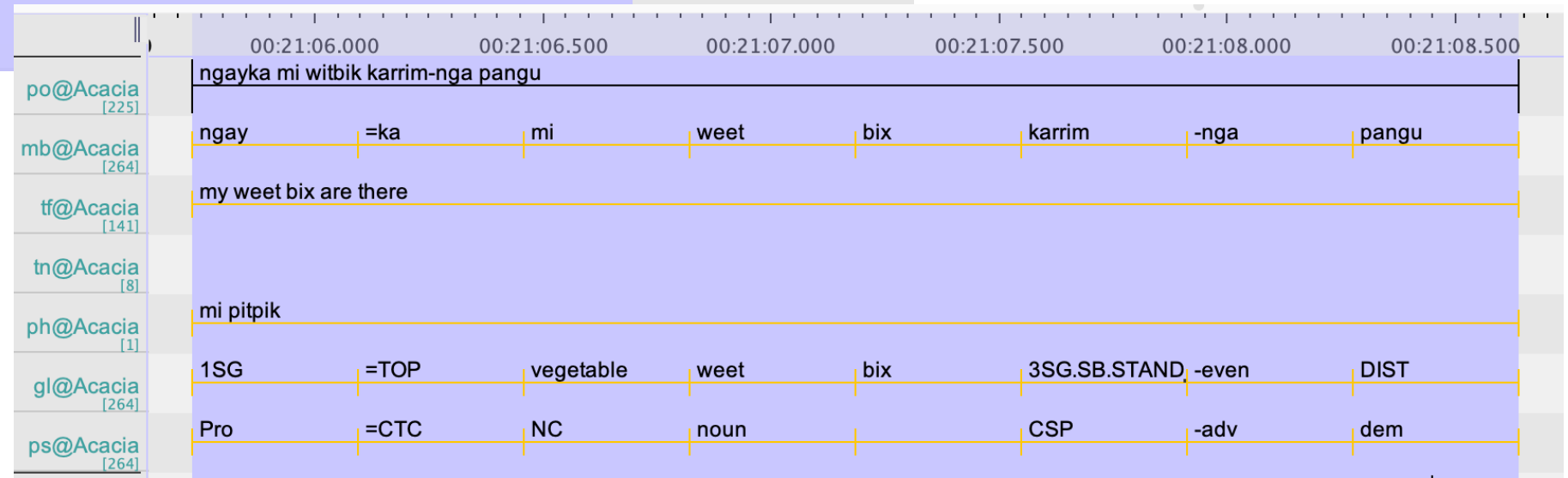
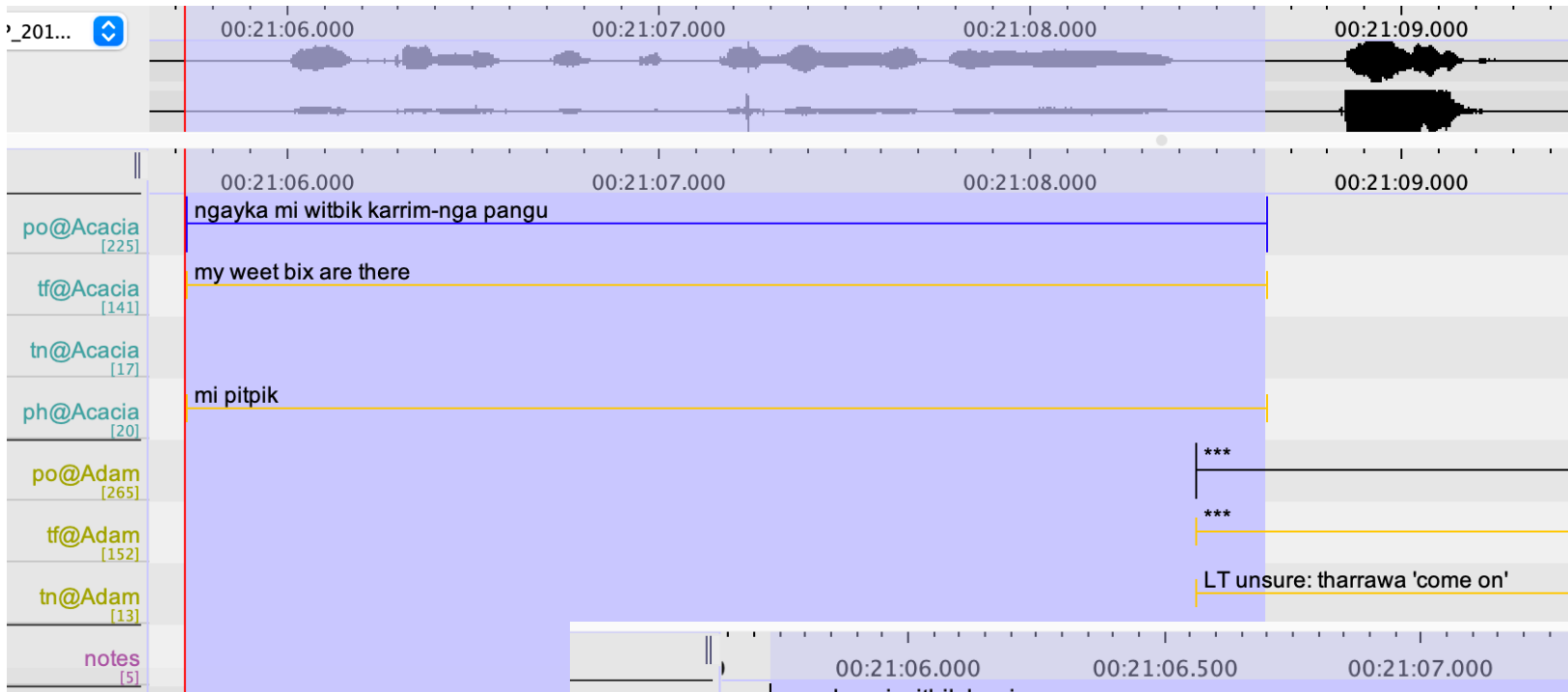
Do minimal annotations for the entire five hours with more detailed annotations for parts of the five hours and/or specific phenomena

Different people can do different tiers depending on their expertise

- e.g. mother for transcription, RA for morphemic breakdown and gloss

Murrinhpatha in ELAN + Toolbox

po	practical orthography
tf	free translation (English)
tn	transcription notes
ph	phonological notes (focus child)
addr	addressee (adult speech)
mb	morphemic breakdown
gl	gloss
ps	part of speech
notes	notes
***	unclear speech



Murrinhpatha, data from LAMP project

Murrinhpatha in ELAN + Fieldworks

po@Emily [466]		thangkunu *-aykay-yu
tf@Emily [440]		what (?are you) calling out for?
tn@Emily [63]		PROMPTED, similar number of syllables & melody as adult's
ph@Emily [37]		thangkunu yayaykayyo
po@Tania [834]	thangkunu thanamkaykay-ya nana	
tf@Tania [816]	say to him 'what are you calling out for?'	
tn@Tania [59]	'him' is the bird	
addr@Tania [233]	CDS: Emily	
notes [67]	They are talking to a bird (pulupulu) in the tree. Pulupulu is also the nickname and totem of Emily's older brother	

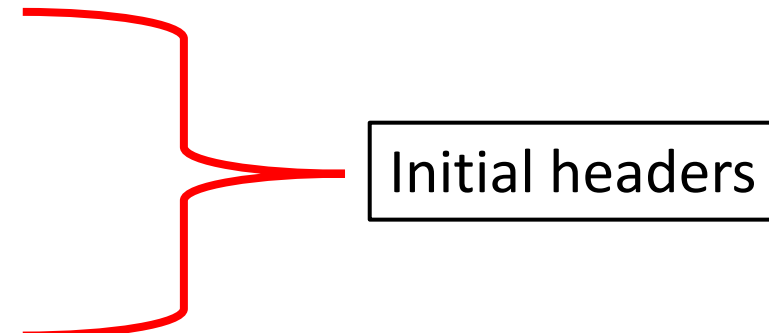
po	practical orthography
tf	free translation (English)
tn	transcription notes
ph	phonological notes (child)
addr	addressee (adults)
notes	notes
*	unclear speech

[605]	thangkunu	thanam	kay	-ya	na	-na
Tania_morph-cf-m [3572]	why	2SG.BE(4).NFUT	call.out	INTJ	2SG.SAY/DO(8).FUT	3SG.M.IO
Tania_morph-gls-e [3573]	INT	CSP	Cov	CTC	CSP	IO
Tania_morph-msa- [3571]	thangkunu	thanam	kaykay	-ya	na	-na
Tania_morph-txt-m [3627]	stem	stem	stem	suffix	stem	suffix
Tania_morph-type [3571]	say to him 'what are you calling out for?'					
Tania_phrase-gls-e [836]	232					
Tania_phrase-segn						

Mongolian in CHAT

```
@Begin
@Languages: mon
@Participants: MOT Mother, CHI Target_Child
@ID: mon|Sketch-Mongolian|MOT|25;05|female|||Mother|||
@ID: mon|Sketch-Mongolian|CHI|3;0.|female|||Target_Child|||
@Media: SAM001, audio
*CHI:  Ügüi. •187_5157•
*MOT:  Zovoogoogüi baina. •4075_5296•
*CHI:  Naadkhyg chini khen gedeg yum. •8963_9704•
*MOT:  Melkhii gedeg. •10445_10944•
*CHI:  Melkhii gedeg yum uu? •11416_12184•
*MOT:  Mmkhn melkhii gedeg. •12192_14424•
*CHI:  Yoooyo. •18203_18440•
*MOT:  Öör bambaruush baigaa yum uu? •18461_19603•
*CHI:  Baikhgüi. •19488_20624•
```

12Sep22[E][CHAT] 32



- Audio and Video Time Marks:
 - 4075_5296• time alignment
- Transcribing directly from an audio file using sonic mode in CLAN
- Translation, morphemic breakdown, gloss and other tiers will be done later

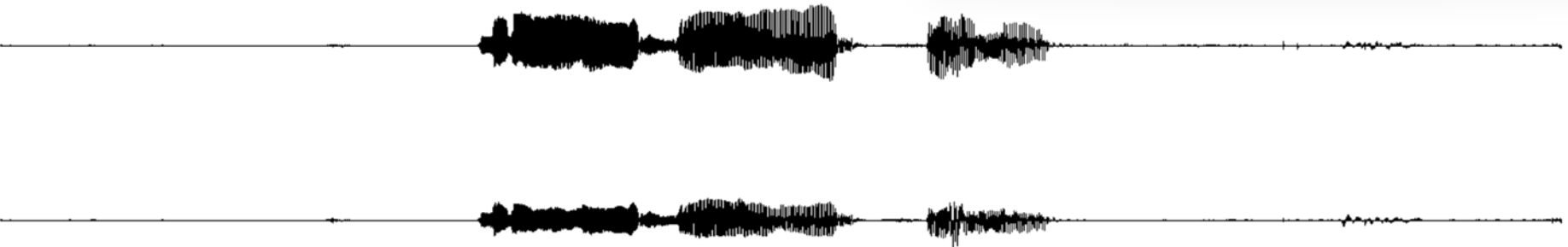
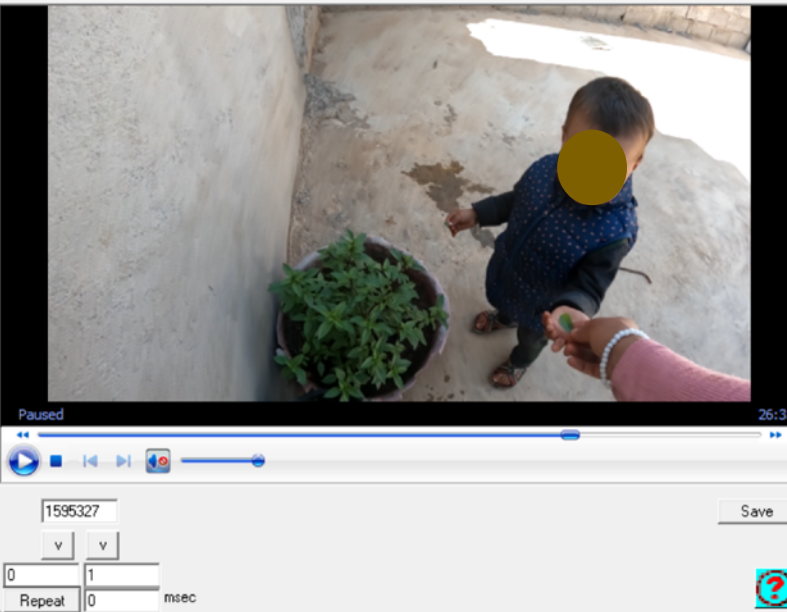
Mongolian, data from Dorjderem Byambasuren

Berber in CHAT

File Edit View Tiers Mode Window Help

86 *MOT: ffi-iyi-d şşkk^war ffi-iyi şşkk^war γ-yi . •
87 %mor: v|ffi&IMP~pro:dat|iyi&1S~dir:prox|d n|şşkk^war
88 v|ffi&IMP~pro:dat|iyi&1S n|şşkk^war prep|γ\$adv:loc|yi .
89 %pho: f:ij:id s:^sk:^war f:ij:id s:^sk:^war Ɂ:i
90 *CHI: huh ? •
91 *MOT: ffi şşkk^war γi-d ad asi-γ zzit . •
92 %mor: v|ffi&IMP n|şşkk^war adv:loc|yi~det:dem:prox|d comp|ad v|as
93 n|zzit .
94 %pho: f:i s:^sk:^war Ɂid ad asiɁ z:it
95 *CHI: huh ? •
96 *SIS: ffi-t-inn ffi-t-inn . •
97 %mor: v|ffi&IMP~pro:obj|t&3SM~dir:dist|inn

26oct23|E|CHAT| 10



+H +V
S *L
-H *R
-V

Ready NUM

- *MOT mother
- *CHI child
- *SIS sister
- %mor morphemic breakdown and gloss
- %pho phonetic transcription

Inuktitut in CHAT

```
*LIZ: anisaarta.  
%eng: Let's go out.  
%mor: VR|ani^go_out+VV|ADV|saag^quickly+VI|ta^IMP_1pS .  
%cod: $VER $ITR:vr  
%arg: $YRF $STR $INT $SUB:nu1:1:hum:crx:lcn:gun:cxp:ugi:rxx:gsn:prn $WOR:1 $MOR:3  
%err:  
%tim: 01:02:07  
%snd: "mae14.sd" 3734737 3736477  
%add: DOM; JIN  
%sit: they try to go out the window  
%com:
```

%eng English translation
%mor morpheme breakdown and gloss
%cod codes for verb types
%arg codes for argument structure
%err errors

%tim time on tape
%snd link to audio sound
%add addressee
%sit situational context
%com comments

Cree in Phon extension of CHAT including Praat

The screenshot displays the CHAT software interface with the following components:

- Media Player:** Shows a video of two people sitting on a couch. Their faces are obscured by white boxes.
- Speech Analysis:** Features a Praat spectrogram of the audio segment. The time axis ranges from 003:41.724 to 003:42.980. The frequency axis is labeled 'Hz' and has a 500 Hz marker. A blue spectrogram shows the vocal formants.
- Record Data:** Displays linguistic analysis for record 95, speaker Daisy. The analysis includes:
 - Orthography: [mâu=tâh âhtut=ikiniwi=ch-h]
 - Morpheme Meaning: [this=loc do=passive.3=0.pl]
 - Morpheme Type: [p,dem+G.pxl=loc IC.initial=passive=CIN]
 - Translation: [We do it like this.]
 - Segment: [003:41.724-003:43.347]
 - IPA Actual: [ondædɛdɪgɪnʊfʰ]
 - IPA Target: [maʊdatotɪgɛnoʃʰ]
 - Actual Morphology: [on=dæ dɛd=ɪgɪnʊ=fʰ]
 - Target Morphology: ['maʊ=da tot=ɪgɛ'no=fʰ]
 - Passives Coding: [B; nSAP; Act; T I; Conj; do; ihtutim]
 - Passives Commentary: [She is talking about scooping them with the spoon. I don't think it is a peeler but a spoon. scooping or dishing: â ūtîh=î=kiniwi=ch=h 'to dish or scoop']
 - Notes: [She is talking about how to peel potatoes. From previous record.]
- Footer:** Shows the session ID 'B1/01-Daisy-03_08_10', a warning icon for '86 warnings', and the tier information 'Tier: Passives Coding Group: 1 Character: 35'.

Cree,
data
from
CCLAS
project

Cree – data from Phon section of Talkbank

```
CHI: îtuht~â~u . ▶  
%pho: thəhesaw  
%mod: itʃ'daw  
%xactmor: thəhes~a~w  
%xmodmor: itʃ'd~a~w  
%xmormea: go~vai.fin~3.sg  
%xmortyp: initial~vai.fin~IIN  
%xtrans: is going
```

*CHI	child; orthographic transcription
%pho	IPA transcription - actual
%mod	IPA transcription - target
%xactmor	morphemic breakdown - actual
%xmodmor	morphemic breakdown - target
%xmormea	gloss of morpheme meaning
%xmortyp	gloss of morpheme function
%xtrans	English translation

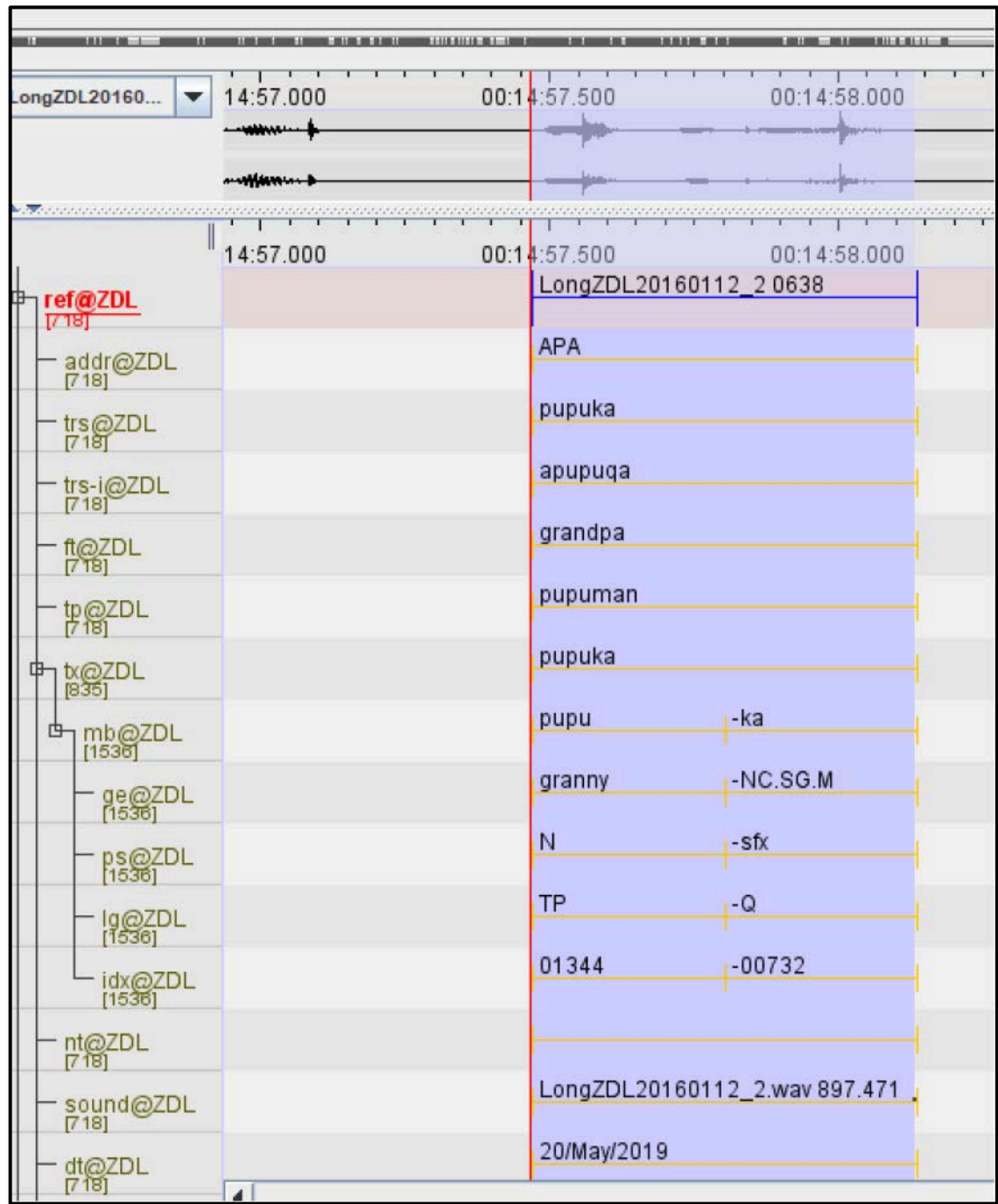
Ani 3;01.18, CCLAS database

<https://phon.talkbank.org/access/Other/Cree/CCLAS.html>

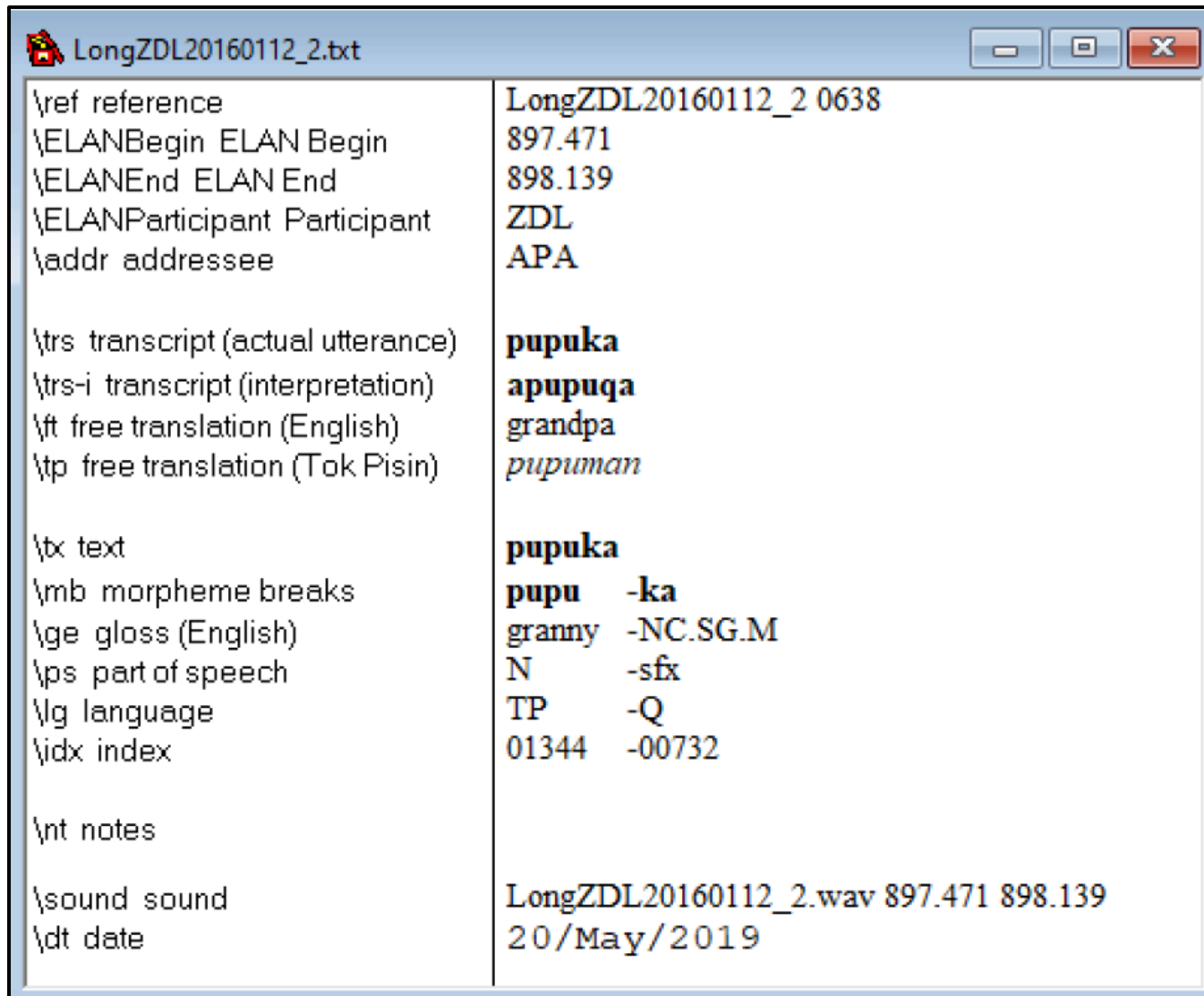
<https://sla.talkbank.org/TBB/phon/Other/Cree/CCLAS/Ani>

Qaqet in ELAN

ref	reference for data file
addr	addressee
trs	transcription
trs-i	interpretation of child utterance
ft	free translation (English)
tp	free translation (Tok Pisin)
tx	text
mb	morphemic breakdown
ge	gloss
ps	part of speech
lg	language
idx	index
nt	notes
sound	link to sound file
dt	date



Qaqet in Toolbox



\ref reference	LongZDL20160112_2 0638
\ELANBegin ELAN Begin	897.471
\ELANEnd ELAN End	898.139
\ELANParticipant Participant	ZDL
\addr addressee	APA
\trs transcript (actual utterance)	pupuka
\trs-i transcript (interpretation)	apupuqa
\ft free translation (English)	grandpa
\tp free translation (Tok Pisin)	<i>pupuman</i>
\tx text	pupuka
\mb morpheme breaks	pupu -ka
\ge gloss (English)	granny -NC.SG.M
\ps part of speech	N -sfx
\g language	TP -Q
\idx index	01344 -00732
\nt notes	
\sound sound	LongZDL20160112_2.wav 897.471 898.139
\dt date	20/May/2019

Qaqet in CHAT

*ZDL:	pupuka
%int:	apupuqa
%eng:	grandpa
%tkp:	pupuman
%mor:	n pupu=grandparent&TP+nc ka=sg&m&Q
%pos:	n+sfx
%add:	APA
%not:	
%snd:	LongZDL20160112_2.wav 897.471 898.139

%int	interpretation
%eng	English translation
%tkp	Tok Pisin translation
%mor	morphemic breakdown and gloss

%pos	part of speech
%add	addressee
%not	notes
%snd	link to audio file

Questions and Discussion

Good luck in data processing!