

Endangered Language Documentation: Bootstrapping a Chatino speech corpus, Forced Aligner, ASR

Malgorzata cavar, Damir cavar, Hilaria cruz
University of Indiana and University of Kentucky

**10th edition of the Language
Resources and Evaluation
Conference, 27 May 2016,
Portorož (Slovenia)**

Agenda

- Work based on the NSF-funded AARDVARC project (<http://info.linguistlist.org/aardvarc/>) (NSF #[1244713](#))
- Creating the first San Juan Quiahije Eastern Chatino (CTP) Speech Corpus
- Experiment with Forced Alignment
- Why?

The Problems

- The Language Resource Bottleneck
- Large Data Collections in Archives
- The Transcription Bottleneck



From:

<http://fitforprostatesurgery.com/get-fit-get-healthy/establish-an-exercise-routine/>

Endangered Language Documentation

- Large collections of recordings over the last decades
 - Limited amount of transcribed, translated, or analyzed data
- Ethnologue lists 7,097 living languages (yesterday)
 - 1% of those languages are well-resourced.
 - Biodiversity and language diversity threatened in a similar way.
 - Large number of languages: unwritten, partially qualitatively documented, etc.
 - Growing gap between the **low 99%** and the **1% highly resourced** languages.

Language Resources

- Audio and Video from documentary linguistic work
 - The Archive of the Indigenous Languages of Latin America at UT Austin (AILLA)
 - The Alaska Native Language Archive
 - DOBES, Max Planck Institute
 - SOAS, London
 - many more
- Lack of transcription, effort
 - 50 to 100 x real-time, i.e. one hour of recording requires 50 to 100 hours of work

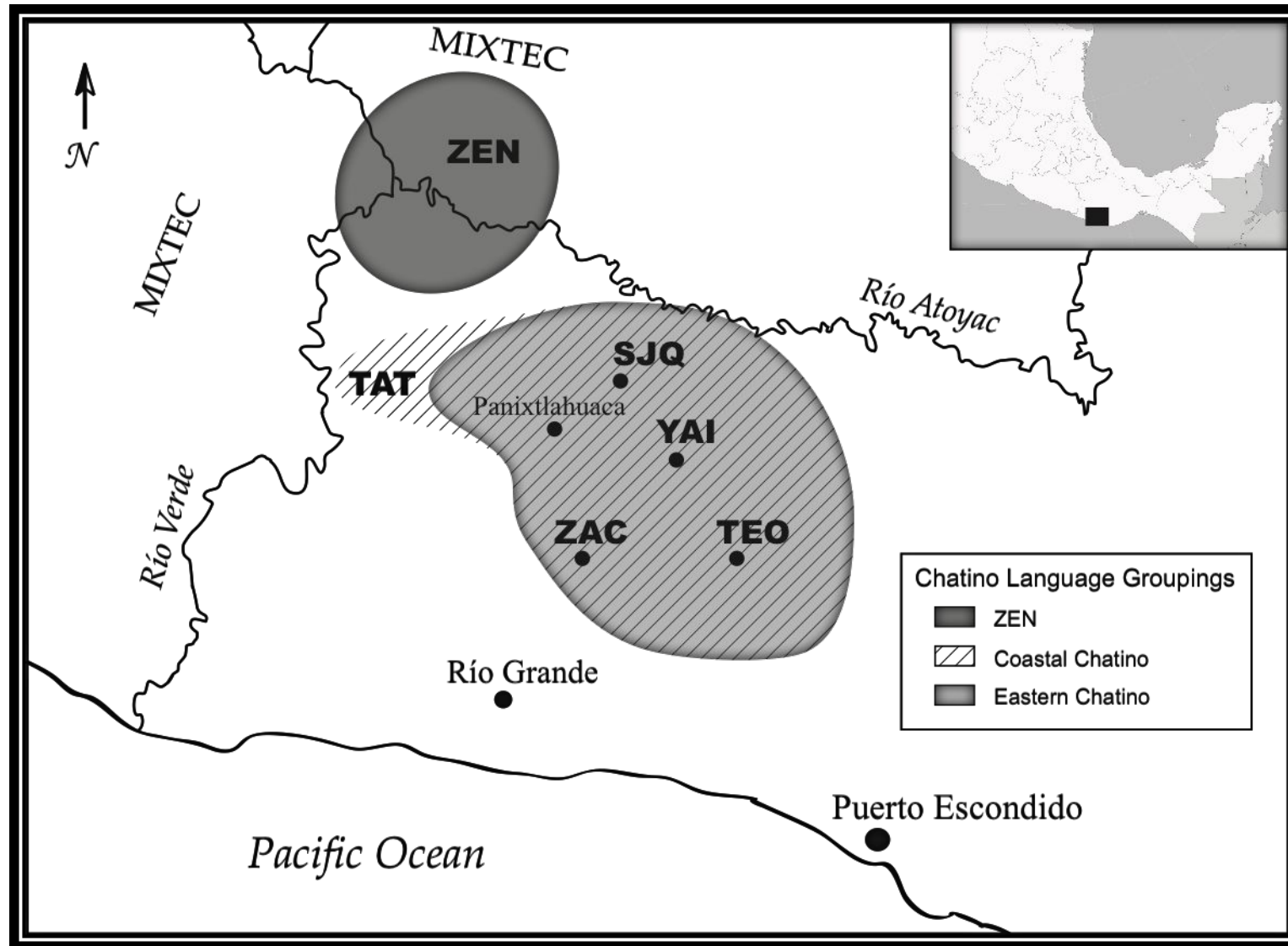
Possible Approach

- Speech and video technologies (AARDVARC project)
 - Forced Alignment to speed up the production of speech corpora
 - Automatic Speech Recognition (ASR) for the recordings could speed up the transcription task.
 - Low number of speakers and longer spoken sequences = less signal variation.
- Task to estimate effort for:
 - Initial speech corpus for Forced Aligner (approx. 2 to 5 hours).
 - Effective corpus size for ASR.

The Chatino languages spoken in Oaxaca, Mexico



Chatino is a group of three languages



ELAN - 2015-03-21_SJQ_EZ-ME_LYSH.eaf

File Edit Annotation Tier Type Search View Options Window Help

Grid Text Subtitles Lexicon Comments Recognizers Metadata Controls

Elena_Chantino

Nr	Annotation	Begin Time	End Time	Duration
1	Na-f kwan-h nya-j ran-f ntyqo-h chaq-f i-j	00:00:00.000	00:00:01.270	00:00:01.270
2	kwiq-j la-e, kwiq-j la-e, ntqen-h ndyan-l wan-a nky-a ntqoh chaq-f Nt...	00:01:01.515	00:01:06.400	00:00:04.885
3	ne-c janq-g in-h, ne-c no-e nkyqan-j re-h qya-j ne-c ntyqo-h chaq-f, n...	00:01:06.590	00:01:11.170	00:00:04.580
4	kwan-h chaq-f no-a, no-a, ne-h jin-c, ye-g qa-j, kyqan-j qa-k chaq-f n...	00:01:11.495	00:01:16.965	00:00:05.470
5	no-a nkan-j an-h krensy-a, na-h an chaq-f, na-h an son-k, na-h an jan...	00:01:17.155	00:01:21.595	00:00:04.440
6	ntqen-g janq-g ti-c nyi-e janq-h, ntqen-g janq-g nyi-h qo-j ji-l qna-g, n...	00:01:22.250	00:01:29.840	00:00:07.590
7	ne-h jin-c, kwa-f sten-h ne-c, ja-a la-l qa-j ran-f ndywenq-h qo-e janq-h	00:01:30.450	00:01:32.900	00:00:02.450
8	Kyqan-j, kwiq-j twen-f no-a wa-c yqan-l, si-k naq-g ti-c ndwa-b qa-k c...	00:01:33.095	00:01:40.835	00:00:07.740
9	Qo-e kwan-h nya-j qan-e nkwa-e ran-f janq-h ntyqo-h chaq-f kanq-g...	00:01:41.445	00:01:47.160	00:00:05.715
10	Wa-c qne-e X ne-c ti-h xka-l yaq-a, wa-c qne-e ne-c ntyqo-h chaq-f k...	00:01:47.360	00:01:51.330	00:00:03.970
11	wa-c ntqan-e qan-e ntqan-e ten-j qo-e wa-f ne-c jin-h in-h	00:01:52.100	00:01:54.895	00:00:02.795
12	Jan-a jan-f	00:02:13.080	00:02:13.930	00:00:00.850

00:13:16.635 Selection: 00:13:16.635 - 00:13:19.350 2715

2015-03-21...

00:13:57.000 00:13:58.000 00:13:59.000 00:14:00.000 00:14:01.000 00:14:02.000 00:14:03.000

Maestra-Epifania [3]

Maestra-Epifania [177]

Elena_Chantino [208]

Elena_Spanish [2]

Elena_English [204]

Elena-relative_Cha [43]

Elena-relative_Eng [40]

kwan-h nya-j ska-a janq-g ja-a nkwa-c tykwiq-e qo-e ska-a janq-a nkwa-c tyqan-e. Ja-a nkwa-c tykwiq-e janq-g janq-h in-h qo-e ja-a nkwa-c tykwiq-e janq-a janq-h. Qo xka-i janq-a in-j qan-a janq-a janq-h.

~there were two of them. One of them could not speak and could not walk and the other one was able to walk, but could not speak

The shown scene was recorded by Lynn Hou.

This is how Chatino sounds



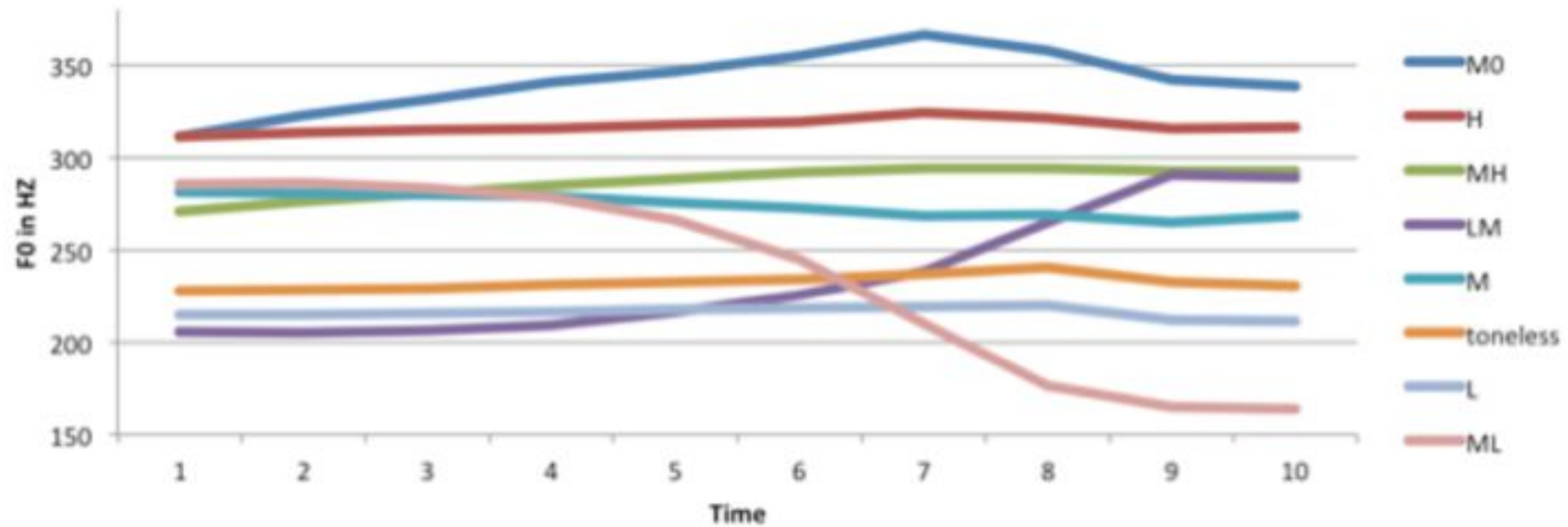
Margarita Balthazar Garcia, Cieneguilla.

Typological features of the Chatino Languages

- Highly tonal (Especially Eastern Chatino)
- Strongly head-marking
- Head-initial syntax
- VSO word order
- Alienable vs. inalienable possession
- complex verbal inflection classes
- Vigesimal numeral system
- no plural marking on nouns
- Excl vs. Incl. pronouns

TONES

Time Normalized F0 of Basic Chatino Tones



Representation of tone



	Number	Tonal group	gloss
Level tones	ska1+0	skaK	sugar
	kla1	klaE	Weaving loom
	kla2	klaC	water pool
	kla3	klaF	dream
	kla4	klaA	fish, old, star
descending	kla24	klaJ	twenty
	tqwa14	tqwaB	cold
	tyuq04	tyuqM	term of endearment
ascending	kla42	klaG	you will arrive
	sqen32	sqenI	scorpion
	kla20	klaH	you will sing
	xkwan40	xkwanL	I will throw it



кта^E 'foráneo (foreign)'

кта^C 'harina (flour)'

кта^F 'chepil (chepil *Crotalaria longirostrata*)'

кта^A 'tabaco, se bañará (tabacco, s/he will take a shower)'

кта^G 'lo sembraras, te bañaras (you will take a shower)'

кта^H 'machucar (to bruise)'

кта^I 'se sembrará (it will be planted)'

Кта^B 'ganado (livestock)'

Approach

- Identification of text for reading and recording
- Recording initial speech corpus of 5 hours
- Transcription and time alignment
- Training a Forced Aligner
 - Prosodylab Aligner (PLA, Python module using HTK)
 - ELAN2split (<https://bitbucket.org/dcavar/elan2split>) corpus creation for PLA
 - Pronunciation dictionary (tokens plus tokenized phonetic transcription)
 - Espeak for Praat: Language model for TXT2Speech
- Extending the speech corpus

Outcome

- Initial speech corpus:
 - Approx. 5 hours speech transcribed and time aligned, PoS-tagged and translated
 - Initial annotation in ELAN (time-alignment correction in Praat)
- Workload (ignoring previous investment in text, transcription schema)
 - Ford Assembly Line approach: bootstrapping initial corpus, FA-training, ...
 - 4 to 6 person weeks
 - Estimate \$ 24,000 – 50,000
 - If we would do that for 3,500 languages: < \$ 84 mil.

Resources

- GORILLA site (<http://gorilla.linguistlist.org/>)
 - Audio, ELAN and Praat transcription/annotation files
 - Corpus licenses: CC BY-SA, i.e. free for commercial use (donation-ware and copyright free resources)
 - Code and software: Apache 2.0 licensed, i.e. free for commercial use
 - Every resource comes with a paper to cite
 - LLOD-linked
 - CLARIN-linked
 - ...

Resources

- Internships at LINGUIST List and Indiana University:
 - Work on LL resources, but also: corpus creation, speech and language technologies, qualitative and quantitative language related research
 - Corpora created with colleagues, students, native speakers and community members:
 - Chatino, Burmese, Turkic (Baharlu, Khorasan, Iran), Croatian, Russian, Spanish, ...
 - Welcoming students from all over the world!

WaC xqweF
qwanJ

Hvala!