



#### **Development of Automatic Speech Recognition** for the Documentation of Cook Islands Māori





Rolando Coto-Solano<sup>1</sup>, Sally Akevai Nicholas<sup>2</sup>, Samiha Datta<sup>1</sup>, Victoria Quint<sup>1</sup>, Piripi Wills<sup>3</sup>, Emma Ngakuravaru Powell<sup>4</sup>, Liam Koka'ua<sup>3</sup>, Syed Tanveer<sup>1</sup>, Isaac Feldman<sup>1</sup> 1. Dartmouth College 2. Massey University Te Kunenga Ki Pūrehuroa 3. Kōrero Rororuia 4. University of Otago Te Whare Wānanga o Ōtakou Language Resources and Evaluation Conference LREC 2022



(1) ASR for Language Documentation

(2) Cook Islands Māori

(3) Data collection and normalization

(4) Two experiments(i) Training with all speakers(ii) Training with held-out speakers

(5) Licenses and documentation workflow

Transcribing text in an Indigenous language is an extremely specialized and expensive process: It can take up to 50hrs to transcribe one hour (Shi et al. 2021).

Using ASR would accelerate transcription (Prud'hommeaux et al. 2021) and help in the creation of linguistic and educational materials.

### **Cook Islands Māori**





East Polynesian

Whangaunga tāta ki te Reo Māori o Aotearoa

Endangered

Indigenous to the Realm of New Zealand

14K speakers (+8K in NZ)

## **Data Collection and Normalization**





Data source: Te Vairanga Tuatua

Large Linguistically rich Under annotated Transcription bottleneck

# **Previous work on CIM NLP**

POS Tagging 92% accuracy (Random forest) http://cimpos.appspot.com/pos.jsp

#### Untrained Forced Alignment 8% error for center of words

#### First ASR Experiment

Transfer model from Dragonfly and Te Hiku Media's DeepSpeech, with 1 hour CIM audio. 31% CER ka tanu a Tere i te taro tam v det nprop acc det n



And then the pig got tangled ē e tāpeka'ia te puaka

e tāpaka te pu'ka

## **Data Collection and Normalization**

Recordings transcribed by Ake, Piripi Wills, Emma Powell and Liam Koka'ua Two main challenges:

-Transforming transcriptions for ASR algorithms

-Accounting for variation (e.g. code-switching)

A V						
I	00:01:04.000	00:01:05.000	00:01:06.000	00:01:07.000	00:01:08.000	00:01:09.00
default						
Speaker 1 Māori Tr	Kua tuku tā rātou k	upenga,	ē kia pōpōiri ak	e, kua mou tā rāto	u ika	
[150]						



Speaker 1 Māori Transcrip	ion Ana Andrew	29.126	31.067	1.941	I runga i te 'enua ko Tupuaki,
Speaker 1 Māori Transcrip	ion Ana Andrew	31.635	32.731	1.096	i te tuātau ta'ito,
Speaker 1 Māori Transcrip	ion Ana Andrew	33.202	37.468	4.266	tē no'o ra tēta'i māpū māro'iro'i, ko Rū tōna ingoa.
Speaker 1 Māori Transcrip	ion Ana Andrew	38.356	39.477	1.121	Kāre ia i te ariki,
Speaker 1 Māori Transcrip	ion Ana Andrew	39.932	42.371	2.439	ē kāre katoa aia i te tamaiti nā te ariki,
Speaker 1 Māori Transcrip	ion Ana Andrew	42.617	43.383	0.766	inārā,

#### **Data Collection and Normalization**

5033 files: 237 minutes (~4 hrs) Median: 6 words (29 chars), 2.3 seconds

10 speakers, 4 islands (Rarotonga, Tongareva, Ma'uke, 'Atiu)

A V						
I	00:01:04.000	00:01:05.000	00:01:06.000	00:01:07.000	00:01:08.000	00:01:09.00
default						
Speaker 1 Māori Tr	Kua tuku tā rātou k	upenga,	ē kia pōpōiri ak	e, kua mou tā rāto	u ika	



Speaker 1 Māori	Transcription Ar	na Andrew	29.126	31.067	1.941	I runga i te 'enua ko Tupuaki,
Speaker 1 Māori	Transcription Ar	na Andrew	31.635	32.731	1.096	i te tuātau ta'ito,
Speaker 1 Māori	Transcription Ar	na Andrew	33.202	37.468	4.266	tē no'o ra tēta'i māpū māro'iro'i, ko Rū tōna ingoa.
Speaker 1 Mãori	Transcription Ar	na Andrew	38.356	39.477	1.121	Kāre ia i te ariki,
Speaker 1 Māori	Transcription Ar	na Andrew	39.932	42.371	2.439	ē kāre katoa aia i te tamaiti nā te ariki,
Speaker 1 Māori	Transcription Ar	na Andrew	42.617	43.383	0.766	inārā,

### **Experiment 1: ASR with all speakers**

Training with three ASR systems:

- Kaldi (Povey et al. 2011)
- DeepSpeech (Hannun et al. 2014)
- Wav2Vec2 / XLSR (Conneau et al. 2020, Baevski et al. 2020)

Training/validation/test splits: 80% (4027 files), 10% (503 files), 10% (503 files)

Random permutations of all the data for all the speakers x20 runs

# **Experiment 1: ASR with all speakers**

DeepSpeech

Kaldi Wav2Vec2

Cook Islands Māori ASR Error rate by type of training (approx. 4 hrs of data)



	WER	CER
Kaldi	$\textbf{17.9} \pm \textbf{1.7}$	$7.5\pm0.8$
DeepSpeech	$41.1\pm2.0$	$21.9\pm1.6$
Wav2Vec2	$22.9\pm2.0$	$\textbf{6.1} \pm \textbf{0.6}$

### **Experiment 1: ASR with all speakers**

English	One day I was just sitting in my car		
Target	i tēta'i rā tē no'o 'ua ara au i roto i tōku motoka	WER	CER
Kaldi	ki tēta'i rā tē no'o 'ua ara 'oki i roto i tōku motoka	15	9
DeepSpeech	i tēta'i a te no'o ara i roto i tōku motoka	31	18
Wav2Vec2	i tēta'i rā tē no'o 'ua ara au i roto i tōku moutakā	8	5
English	I was sure that it was the pig who had rooted (it up)		
Target	kua kite ra 'oki au ē nā te puaka i ketu	WER	CER
Kaldi	kua kite rā 'oki au e nā te puaka i ketu	18	5
DeepSpeech	kite rāi koe i nā te puaka i ki	55	38
Wav2Vec2	kua kite rā 'aki au ē nā te puaka i kit	27	10
English	Absolutely, it will get mixed up		
Target	āe 'oki ka iroiro atu	WER	CER
Kaldi	'aere ka'iro i roa atu	80	50
DeepSpeech	āe ki ka'iro 'oki roa te	100	50
Wav2Vec2	āe 'oki kā'iro'i roa atu	40	23

### **Experiment 2: Held-Out Speakers**

We split the sets to ensure that there were speakers who were never seen during training/validation.

Partition 1Training and validation sets:T1, T3, M1, M2, B, J(90% of files)Test set:A, K, T2, R(10% of files)

#### Partition 5

 Training and validation sets:
 A, B, J, K, T2, R, T3, M1, M2, T1
 (70% of files)

 Test set:
 J
 (30% of files)

Random permutations within the train/val sets. x5 times

#### **Experiment 2: Held-Out Speakers**

Partition	Train-Validation-Test	WFR	CER	Test	% total	% total
1 artition	Splits (#files and %)	WLK	CLK	speaker(s)	files	time
1	4036 - 504 - 493	$32.9\pm0.9$	$8.4\pm0.2$	А	3.7	3.4
	80% - 10% - 10%			Κ	3.6	4.5
				T2	2	4.5
				R	0.5	1.0
2	4007 - 500 - 526	$40.1\pm1.9$	$11.0\pm0.5$	T3	6.9	7.6
	80% - 10% - 10%			M2	3.4	7.2
3	3849 - 481 - 703	$64.5\pm3.1$	$24.5\pm1.0$	M1	14.0	8.0
	76% - 10% - 14%					
4	3769 - 419 - 845	$25.0\pm0.0$	$5.9\pm0.3$	В	17.0	18.5
	75% - 8% - 17%					
5	3268 - 408 - 1357	$50.0\pm0.0$	$16.4\pm0.5$	J	27	30
	65% - 8% - 27%					
6	3532 - 392 - 1109	$65.9 \pm 1.9$	$23.0\pm0.2$	T1	22	15
	70% - 8% - 22%			7		
Average		$46.4\pm15.6$	$14.9\pm7.2$			

#### 

#### Partition 5

Meaning: Target: Inference: *From morning till night.* mei te pōpongi mai e pō mei te pupongi mai ēpo (CER=16, WER=50)

#### Partition 1

Meaning: When we die we die, when we live we live. Target: mē mate tātou kua mate mē ora kua ora Inference: mē mati tātou kua mate me ora kua ra (CER=8, WER=33)



Deep Learning systems can indeed work with extremely low-resource languages.

We need to test if different transcriptions will have an effect (e.g. should long vowels be  $\bar{a}$  or ax?)

CIM has relatively few phonemes (9~10 cons, 10 vowels) and few affixes. This might be enhancing the results.

## Analysis

We need recordings from more islands.

We need to generate a virtuous cycle where transcriptions create revitalization materials and create further training materials for the ASR.



Model and data are freely available.

Kaitiakitanga License:

- Non-commercial use allowed
- In consultation with Indigenous community

Objective: Maintain data sovereignty.

#### Conclusions

# (1) We built an ASR with CER~6 for Cook Islands Māori

(2) Deep Learning ASR methods are approaching a point of usability for low-resource languages.