

Language Models for Code-switch Detection of te reo Māori and English in a Low-resource Setting

Anonymous ACL submission

Abstract

Te reo Māori, New Zealand’s only indigenous language, is code-switched with English. Most Māori speakers are bilingual, and the use of Māori is increasing in New Zealand English. Unfortunately, due to the minimal availability of resources, including digital data, Māori is under-represented in technological advances. Cloud-based systems such as Google and Azure support Māori language detection. However, we provide experimental evidence to show that the accuracy of such systems is low when detecting Māori. Hence, with the support of Māori community, we collect Māori and bilingual data to use natural language processing (NLP) to improve Māori language detection. We train bilingual sub-word embeddings and provide evidence to show that our bilingual embeddings improve overall accuracy compared to the publicly-available monolingual embeddings. This improvement has been verified for various NLP tasks using three bilingual databases containing formal transcripts and informal social media data. We also show that BiLSTM with bilingual sub-word embeddings outperforms large-scale contextual language models such as BERT on down streaming tasks of detecting Māori language. The best accuracy of 87% was obtained using BiLSTM with bilingual embeddings for detecting code-switch points of bilingual sentences.

1 Introduction

Te reo Māori (referred to as Māori) is New Zealand’s only indigenous language, spoken by 4.5% of the total population of 5 million. Māori speakers are generally bilingual, and code-switching between Māori and English is common. Māori revitalisation efforts have increased Māori use in the otherwise English-speaking country. Hence, detecting Māori language and code-switch instances is a prerequisite to analysing language data. As Māori and English both use the Roman script, currently annotations are done manually, making the process time-consuming, slowing down

research and technology development. Consider the following sentences:

- (a) Pērā anō i ngō mate kua hinga atu i te motu.
- (b) I want to give no offence to my mate Willie Jackson, but once a week hardly qualifies as the significant Māori voice.

where green indicates Māori, red is used to indicate that the word has same spelling in Māori and English, and the remaining are English. Based on expert knowledge, we know the word mate in the sentence (a) is Māori, while sentence (b) is English. In this research, we focus on two primary tasks:

Task 1: Language Detection (LD) - detecting Māori language words from input text.

Task 2: Code-switch Detection (CS) - detecting Māori to English or English to Māori code-switch points from input text.

There is limited Māori-only and Māori-English bilingual data available. We collected data by seeking feedback from the Māori community, where data-sharing is based on trust. As researchers, we remain guardians of the data, ensuring data sovereignty (Stats, 2020). Hence, all the resources shared from this study are bound by the Kaitiakitanga license (Te-Hiku-Media). This paper presents one of the first research to use advances in NLP to detect Māori and code-switching. There are no existing models using NLP techniques for code-switch detection. The cloud-based services Google and Azure are the only options available for language detection. This paper’s contributions are:

1. Evaluation of detecting Māori using cloud-based services such as Google and Azure.
2. Pre-training Māori-English bilingual, and Māori-only monolingual sub-word embeddings using the collection of data. Experiments using three different bilingual data for various NLP tasks show that bilingual embeddings outperform monolingual embeddings.
3. Providing evidence to show large scale language models such as Bidirectional Encoder

085 Representations from Transformers (BERT)
 086 are outperformed by BiLSTM with non-
 087 contextual sub-word bilingual embeddings for
 088 low-resourced language such as Māori.

- 089 4. Providing baseline results for detecting low-
 090 resourced Māori and code-switch between
 091 Māori-English language pair.

092 2 Te reo Māori (The Māori Language)

093 Māori is a Polynesian language belonging to the
 094 Austronesian family. Phonologically, Māori has ten
 095 consonants /p t k m n ŋ f r w h/. The Māori vowel
 096 system is described by five short vowels /i e a o
 097 u/ (Bauer et al., 1993). Orthographically, there is
 098 mostly a one-to-one mapping of a Māori phoneme
 099 to a grapheme, except two digraphs, ‘wh’, which is
 100 /f/, and ‘ng’ which is /ŋ/. In modern orthography,
 101 long vowels are denoted with a macron (e.g. ā).
 102 Long vowels are denoted in modern orthography
 103 with a macron (e.g. ā). In older text, they are
 104 sometimes expressed as double vowels (e.g. aa),
 105 with an umlaut (e.g. ä), or ignored completely. In
 106 addition, there is some regional variation in the
 107 way words are spelt (e.g. Aorangi vs Aoraki). This
 108 contrasts with English, which has a non-phonemic
 109 orthography. The Māori syllable structure consists
 110 of a nucleus, which may be occupied by a vowel (or
 111 a diphthong), and an optional onset (syllable start)
 112 occupied by a single consonant. Hence, consonant
 113 clusters are not present in Māori (Harlow, 2007).

114 3 Related Work

115 Research using NLP to tasks relating to Māori
 116 is relatively young. Examples include statistical
 117 machine translation for Māori-English pair (Mo-
 118 haghagh et al., 2014), and the inclusion of Māori
 119 language detection and translation using cloud ser-
 120 vices Google and Azure. Keegan (2017) (Keegan,
 121 2017) indicates that although the growth of cloud
 122 services for Māori translations is welcoming, due
 123 to the minimal availability of digitised Māori data,
 124 the resulting output is inaccurate. Google also ac-
 125 knowledges that for low-resource languages, the
 126 quality of language detection and automatic ma-
 127 chine translation is far from perfect (Blog)

128 We present the first research that uses deep
 129 learning techniques to detect code-switch between
 130 Māori and English. Hence, except for the above-
 131 mentioned cloud-services, we are limited by the
 132 availability of baseline systems for Māori language
 133 detection and Māori-English code-switch detection.

Name and Database	# Words
Māori only	
D1: Te Taka Database* (Keegan, 2021)	9,862,131
D2: Nga Mahi corpus (James et al., 2020)	81,036
D3: Māori Wikipedia	431,280
D4a: LMC Corpus (LMC)	5,486,328
<i>Total size of Māori-only database = 92 MB</i>	
Māori and English	
D4b: LMC Corpus (LMC)	7,197,059
D5: Niupepa (Māori Newspapers) (Niupepa)	5,050,988
D6: Twitter Corpus*(Trye et al., 2019)	48,289,375
<i>Total size of bilingual data = 0.4 GB</i>	

Table 1: Māori-English Words (MEW) database. ‘*’ indicates private collections of data.

We use approaches that were inspired by the literature on other language pairs. Examples include XNLI cross-lingual classification benchmark (Conneau et al., 2018) where the bidirectional Long short-term memory (BiLSTM) model was used across several low resource languages, including Swahili and Urdu; and code-switch detection using BiLSTM and Character-LSTM for language pair English-Hindi (Lal et al., 2019; Mukherjee et al., 2019). XNLI benchmark uses fastText common-crawl embeddings (denoted as E300 in this paper) and aligns it with the MUSE library. Comparison among deep learning models shows that adding background information through sub-word pre-trained embeddings trained using fastText and in the form of lexicons improves the overall performance of deep neural networks on databases of low-resource languages (Adouane et al., 2018).

Transformers such as BERT is the state-of-the-art in many NLP tasks, including language detection, name entity recognition, and machine translation (Devlin et al., 2019; Conneau et al., 2020). There are many large scale multilingual models, such as XLM-R (Conneau et al., 2020) and multilingual BERT (mBERT) (Devlin et al., 2019) trained on more than 100 languages. Research shows that for languages that are under-sampled during training, the effectiveness of large scale multilingual models such as mBERT are sub-optimal (Wu and Dredze, 2020; Wang et al., 2020). In comparison to the contextual representations like BERT, embeddings with sub-word representation are more data-efficient when data availability is limited (Wu and Dredze, 2020). Furthermore, Muller et al. (2021) (Muller et al., 2021) provides evidence to show that many under-sampled or unseen languages during training –such as Maltese or Narabizi– code-mixed with French perform worse when using mBERT compared to an RNN with

173	non-contextual dependency parsing baseline. It has	223
174	been shown that for such unseen or under-sampled	224
175	languages, there is a need to further train or fine-	225
176	tune directly with available raw data in the unseen	
177	target languages (Muller et al., 2021).	
178	4 Databases	226
179	Due to the low-resource nature of the Māori, there	227
180	is no single extensive database. We collected text	228
181	data from different sources to form the Māori-	229
182	English Words (MEW) database, as summarised	230
183	in Table 1. MEW database contains legal context,	231
184	stories, social media posts and newspaper articles.	232
185	The unlabelled MEW database is used to pre-train	233
186	bilingual and Māori-only monolingual embeddings.	234
187	We use three labelled databases for experiments:	235
188	Hansard database, MLT corpus, and RMT corpus.	236
189	Details of these databases are provided in Table 2.	237
190	Hansard database contains the New Zealand Par-	238
191	liament debates from 2003 onwards. Together	239
192	with experts in Māori (Media), we have labelled	240
193	the Hansard database, where English or Māori la-	241
194	bels are assigned using linguistic rules and manual	242
195	checking. Each sentence in the databases is marked	243
196	as Māori, English or bilingual. Each word of each	244
197	sentence is labelled as Māori or English. The re-	245
198	sulting data includes 102,559 bilingual, 1,909,876	
199	English-only and 8,826 Māori-only sentences.	
200	Labelled Māori Loanword Twitter (MLT) corpus	
201	is a small database, where each tweet is labelled as	
202	‘relevant’ and ‘irrelevant’, based on the presence of	
203	a pre-determined set of Māori loanwords in a given	
204	tweet. Given detecting Māori language in tweets	
205	is a prerequisite to this task, we consider this task	
206	also as a Māori LD task. Reo Māori Twitter (RMT)	
207	corpus contains tweets, where at least 80% of text	
208	is in Māori. RMT corpus provides a list of 879,000	
209	Māori words across the tweets. We use this corpus	
210	also for LD task where the aim is to detect the	
211	Māori words identified by the researchers.	
212	5 Language Models and Classifiers	246
213	This section provides details of the language mod-	247
214	els and classifiers we used. We evaluate the perfor-	248
215	mance of cloud-based language detection systems	249
216	from Google and Azure for Māori. We represent	250
217	text as bag-of-words and sub-word embeddings us-	251
218	ing fastText. We use logistic regression and multi-	252
219	nomial naive Bayes as baseline classifiers for lan-	
220	guage detection. We also use neural networks such	
221	as RNNs and CNNs to train and evaluate language	
222	detection and code-switch detection tasks. Further-	
	more, we perform transfer learning of pre-trained	223
	transformer models, BERT and mBERT, for the	224
	down streaming task of language detection.	225
	5.1 Cloud-based Online Tools	226
	Google Translate (Google) and Microsoft Azure	227
	Cognitive Services language detection (Microsoft)	228
	are two popular cloud-based online tools that can	229
	detect multiple languages. Google supports 108	230
	languages, including New Zealand English and	231
	Māori. Google’s RNN-based GNMT model (Wu	232
	et al., 2016) showed significant improvements in	233
	enabling translations to cover many languages, in-	234
	cluding low-resourced languages. Google recently	235
	replaced the GNMT model with a hybrid model	236
	(transformer encoder and RNN decoder). This	237
	model has shown significant improvements to the	238
	other machine translation systems. Azure’s cog-	239
	nitive services can translate 100+ languages, in-	240
	cluding Māori. Azure’s early-stage neural net-	241
	work model (Xiong et al., 2017) included a CNN-	242
	BiLSTM architecture. Recently, Azure has com-	243
	combined several machine learning algorithms and	244
	neural networks to provide various cognitive services.	245
	5.2 Bag of words	246
	Bag of words (BOW) is an effective method (Gold-	247
	berg, 2017; Joulin et al., 2016b) to represent text	248
	as a sparse vector, where the order of words in a	249
	document is not considered. The number of occur-	250
	rences of a word or a binary value indicating that	251
	the word is present in the document is stored.	252
	5.3 Word Embeddings	253
	For language processing tasks, continuous word	254
	representations such as word embeddings trained	255
	on large unlabelled databases facilitate effective	256
	representation learning (Bojanowski et al., 2017;	257
	Joulin et al., 2016a). Here, we use fastText (Bo-	258
	janowski et al., 2017) to learn word embeddings,	259
	as novel words not present during training can	260
	also be represented using fastText-based embed-	261
	dings. This can be beneficial for a low-resource	262
	setting. FastText supports two word embeddings	263
	models: continuous bag-of-words (CBOW) and	264
	Skip-grams (Mikolov et al., 2013). The CBOW	265
	predicts the specific word from the source context.	266
	Skip-gram predicts the source context from the	267
	specific word. The embeddings in this research	268
	are trained to the specifications of Wikipedia and	269
	common crawl fastText models (Grave et al., 2018)	270
	(referred to as E300) for both CBOW and Skip-	271

Data	# Sentences	# Words	Text	Labels	Task
Hansard data (Hansard)	2,021,261	36,757,230	formal	word-level & sentence level language labels	LD, CS
MLT corpus (Trye et al., 2019)	2,500	50,000	informal	tweet level labels: relevance/irrelevance	LD
RMT corpus (Trye et al., 2022)	79,018	1,000,000	informal	Māori words are identified and labelled	LD

Table 2: Databases used for experimental evaluations. LD: Language Detection, CS: Code-Switch Detection.

Embeddings Model	Data	Size
Monolingual Embeddings		
E300 (Grave et al., 2018)	Downloaded	7GB
Māori-300/300SG	D1 - D4a	3GB
Bilingual Embeddings		
Model-Māori-Eng-300/300SG	D1 - D6	3GB

Table 3: Outline of fastText pre-trained 300 dimensional embeddings. The MEW database (Table 1) was used for training. ‘SG’: Skipgram model, otherwise it is CBOW.

gram¹. E300 uses the CBOW method, character n-grams of length 5, window of size 5, 10 negative samples per positive sample with 300 dimensions. The learning rate is 0.05. Table 3 provides details of our bilingual embeddings, which are available to on request, including E300 details for comparison.

5.4 Baseline Classifiers

We use multinomial naive Bayes (John and Langley, 1995) and logistic regression (LR) (Cox, 1958) to classify text features represented by BOW and static word embeddings. LR is a statistical model used to analyse databases where independent variables determine an outcome. Naive Bayes (John and Langley, 1995) is an easy to build supervised learning algorithm, which applies Bayes’ theorem with the “naive” assumption of independence.

5.5 Convolutional Neural Network (CNN)

CNN for text (Kim, 2014) combines one-dimensional convolutions with a max-over-time pooling layer and a fully connected layer. If $x_{i:i+j}$ is a concatenation of words from a sentence, each word, x_i, x_{i+1}, \dots is mapped to its k -dimensional embeddings using word embeddings. A new feature is produced using convolution. Max-over-time pooling is applied over the feature map to capture the most important feature value. The final prediction is made by computing a weighted combination of the pooled values and applying Softmax function.

5.6 Recurrent Neural Networks (RNN)

RNNs (Rumelhart et al., 1986) are designed to handle sequential data, such as text, where the data contains complex temporal dependencies and

¹Embeddings trained on a 4 core Intel i7-6700K CPU @ 4.00GHz with 64GB of RAM. Average time: <30 minutes.

hidden information. Long Short Term Memory networks (LSTM) (Hochreiter and Schmidhuber, 1997) are modified RNNs designed to overcome the issue of vanishing gradient with RNNs. LSTM consists of a gating mechanism, input gate, forget gate, and output gate, ensuring a constant error flow and avoiding long-term dependency problems. The memory in LSTM is stored in an internal state, and the three gates play a vital role in deciding which information be included, added or removed from the memory. Over time, the memory cells learn which information is essential based on the weights. Bidirectional RNNs are widely used extensions where the input sequence is fed from beginning to end and from end to beginning. For BiLSTM (Grave et al., 2018), given there are two LSTM layers, the hidden layer output is split into two - for forward and backwards passes over the input.

5.7 Transformers

BERT (Devlin et al., 2019) is one of the early transformer models that apply bidirectional training of encoders (Vaswani et al., 2017) to language modelling. The 12-layer BERT-base model with a hidden size of 768, 12 self-attention heads, 110M parameter neural network architecture was pre-trained from scratch on BookCorpus and English Wikipedia. The mBERT-base (Devlin et al., 2019) model uses the same pre-training objective as BERT-base and is pre-trained with Wikipedia text of 104 languages with most articles. In this research, we use BERT and mBERT to refer to BERT-base and mBERT-base.

6 Experimental Setup

We experiment with various language models and classifiers for two main tasks: language detection (LD) and code-switch detection (CS). Our ultimate goal is to find a combination of language modelling and NLP techniques to improve the overall accuracy of LD and CS tasks. We use three databases to evaluate these tasks with details provided in Table 2. We use the Hansard database sentences as the primary data for training and testing. All three datasets were pre-processed by lower-casing and using regular expressions to remove punctuation us-

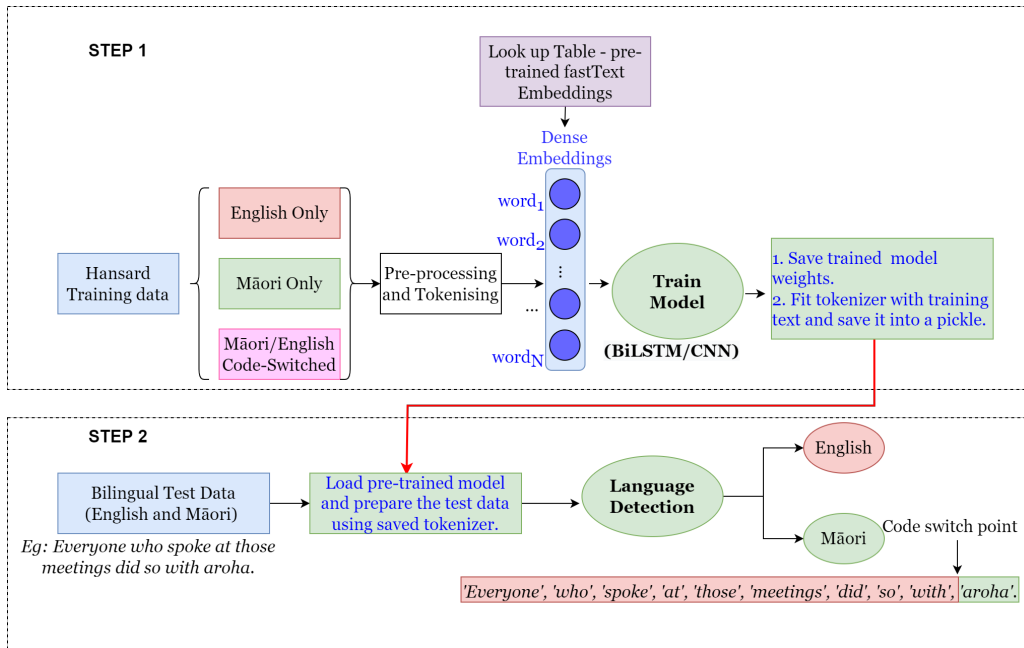


Figure 1: Code-switch detection using neural networks. Example shows ‘English’ words {Everyone, who, spoke, at, those, meetings, did, so, with} are detected as ‘English’ and ‘aroaha’ detected as ‘Māori’.

ing Python 3.9 library with Pandas data frame. All experimental results are obtained from a random seeds training-testing scheme; 70% of the shuffled data is used for training, with 10% for validation and 20% for testing, and averaged over three runs. The variation of these three independent runs is within a range of ± 0.015 .

To represent text we use fastText to pre-train embeddings (see Table 3) and BOW. An overview of code-switch detection using trained models such as BiLSTM and CNN is presented in Figure 1. This diagram is an example to demonstrate the system we used for end-to-end code-switch detection using neural networks. Step 1 includes training and evaluating a neural network. We use the training set of the Hansard database to train the model and use validation loss as the stopping condition to avoid over-fitting. In step 2, we load the trained model and detect languages at the word level on testing data. Once the language detection is done, the points in the sentence where the language labels switch from Māori to English or from English to Māori are marked as code-switch points.

Neural network models presented in this research are implemented using Keras/Tensorflow. Adam (Kingma and Ba, 2015), an adaptive learning rate optimisation algorithm, is used as the optimiser for neural networks. Softmax activation function is used in the output layer of the network. We use a combination of dropout (Srivastava et al., 2014), with a rate of 0.5, and early stopping (Zhang et al.,

2017) to avoid over-fitting. We use a maximum length of 250 tokens for BiLSTM and CNN, and padding for sentences with less than the maximum length. The embeddings layer is with a dimension of 300. The hidden units of BiLSTM are 128, and the hidden units of one-dimensional convolutions are 128. For both CNN and BiLSTM, categorical cross-entropy is used as the loss function.

We also perform transfer learning of pre-trained transformers, BERT and mBERT on the downstream task of language detection. We use batch size of 16, maximum sequence length of 256 and learning rate of $1e-5$. For both BERT and mBERT, the loss and accuracy were reported at each epoch. For both BERT and mBERT, the model converges fast, needing an average of 5 epochs per run.

All evaluations were done using Sklearn metrics (Scikit-Learn). Evaluations using baseline classifiers such as multilingual naive Bayes and LR with BOW and static features from embeddings require CPU only² machines and are very quick to train and evaluate. Neural networks require GPU devices³ for efficient training and testing. The average training time for CNN was 150-180 minutes, and BiLSTM was 300-360 minutes, while BERT and mBERT required 240 minutes per epoch being trained for an average of 5 epochs. The testing time for trained deep learning models is rapid, requiring a few minutes. The code used in this research is

² 4 core Intel i7-6700K CPU @ 4.00GHz with 64GB of RAM.

³ 12 core Intel(R) Xeon(R) W-2133 CPU @ 3.60GHz, GV100GL

made available⁴.

We present overall macro-F1 score and weighted-F1 score to provide different insights (Toftrup et al., 2021; Khanuja et al., 2020). We also provide label F1-score where needed. Macro-F1 provides average per-language results and is equally important to all languages. The weighted-F1 score considers the popularity of the languages in the data set.

The Nemenyi posthoc test (95% confidence level) identifies statistical differences between learning methods. Critical Difference (CD) plots show the average ranking of individual F1 scores obtained using various language models. The lower the rank, the better the model is. The difference in average ranking is statistically significant if there is no bold line connecting the two settings.

7 Experimental Results

The results are presented for the language detection (LD) tasks and code-switch detection (CS) tasks. The language detection task is a crucial first step for detecting code-switching (Rijhwani et al., 2017; Barman et al., 2014). First, we present the results of the language detection tasks using the three databases (Table 2), followed by the results of the code switch task using the Hansard database. As indicated in the experimental setup, all experimental results are obtained from a random seeds training-testing scheme and averaged over three runs. The variation of these three independent runs is within a range of ± 0.015 .

7.1 Task 1: Language Detection

7.1.1 Cloud-based Online Tools

To analyse the effectiveness of using Google Translate and Azure services to detect Māori (and English), we experimented with the test set of the Hansard database. Google Translate detected 99.7% of the words, and Azure detected 97.8% of the words correctly. Figure 2 presents pie charts of the resulting language detection for ‘Māori’ word (i.e. the gold-standard labels for the words is ‘Māori’). For Māori words, Google Translate detected with an accuracy of 65.2%, and Azure detected with an accuracy of 52%. Although the accuracy of Google Translate was better than Azure, the error rate of both services are too high for Māori language detection. In addition, apart from wrongly

⁴Pre-trained bilingual and monolingual embeddings are available for researchers on request. Experimental details, model implementations, and trained language models are available for researchers, all bound by the Kaitiakitanga license: <https://github.com/MaoriEnglish-Codeswitch/MaoriEnglish-CodeSwitch-Detection>

Model	Data	Results
	Multi-class	Macro-F1
Multinomial NB (BOW)	Hansard	0.887
LR (BOW)	Hansard	0.913
LR (Eng300)	Hansard	0.831
LR (Māori-Eng-300)	Hansard	0.853
LR (Māori-Eng-300SG)	Hansard	0.859
Binary		F1-score
LR (Eng300)	MLT corpus	0.833
LR (Māori-300SG)	MLT corpus	0.812
LR (Māori-Eng-300)	MLT corpus	0.849
LR (Māori-Eng-300SG)	MLT corpus	0.846

Table 4: Macro-F1 scores and F1-scores for the validation set of Hansard database and labelled MLT corpus respectively, where BOW or sentence level features are used to represent text. **Bold**: best results for each task.

detecting Māori words as English, around 14-21% of the words were classified as various other languages by both cloud services.

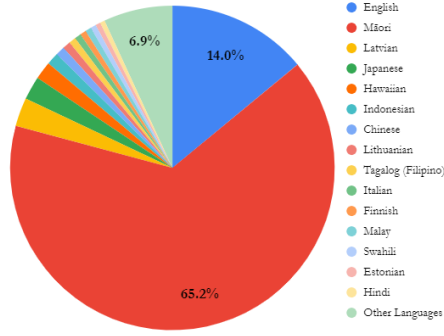
7.1.2 Baseline Classifiers

LD task using the Hansard database is a multi-class classification problem at the sentence level (classes: Māori, English or Code-Switched sentence). The LD task using MLT corpus is a binary classification problem of relevant/irrelevant tweets based on the usage of the Māori loanwords. Table 4 presents overall macro-F1 and F1 scores for the LD task using Hansard database and MLT corpus, respectively, where BOW and static word embeddings at the sentence level (or tweet level) are used to represent the text. We obtain embeddings for each sentence by computing the vector sum of the embeddings for each word in the sentence. This vector sum is then normalised to have length one, to ensure that sentences of different lengths have representations of similar magnitudes. The bilingual embeddings perform better than monolingual embeddings for both Hansard and MLT corpus. However, BOW outperforms static embeddings feature representation for LR.

7.1.3 Neural Networks

After evaluating the performance of baseline classifiers, we further proceed with LD task using neural networks. As the size of the labelled MLT corpus is small, it is insufficient for training and evaluating neural networks. Table 5 presents macro-F1 and weighted-F1 scores obtained using the validation set of the Hansard database for performance comparison across language models. The macro-F1 score is an unweighted average score of all the classes. In comparison, weighted-F1 scores

Google's Language detection for Māori Words



Azure Text Analytics Detect Language for Māori Words

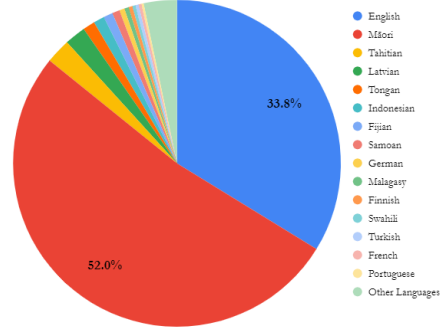


Figure 2: Pie Chart of the languages detected by Google (left) and Azure (right) at word level for the test set of the Hansard Database. The gold-standard label for all the words used here is ‘Māori’.

Model	Macro-F1	Weighted-F1
Monolingual Embeddings		
CNN (E300)	0.946	0.985
CNN (Māori-300)	0.905	0.986
CNN (Māori-300SG)	0.914	0.990
BiLSTM (E300)	0.943	0.996
BiLSTM (Māori-300)	0.926	0.995
BiLSTM (Māori-300SG)	0.940	0.995
Bilingual Embeddings		
CNN (Māori-Eng-300)	0.963	0.995
CNN (Māori-Eng-300SG)	0.969	0.996
BiLSTM (Māori-Eng-300)	0.984	0.997
BiLSTM (Māori-Eng-300SG)	0.989	0.997
Contextual Embeddings		
BERT-base	0.931	0.988
mBERT-base	0.946	0.991

Table 5: Comparison of results for the Hansard database (validation set) with various models. **Bold**: best results.

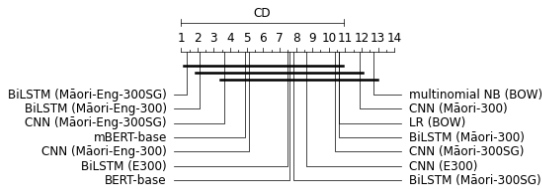


Figure 3: Critical difference plots identifying statistical differences between models presented in Tables 4 & 5.

are higher than macro-F1 scores across the models. The imbalanced distribution in the data, where labels are predominantly English, is reflected in the scores where the minority classes penalise the macro-F1 scores. Bilingual embeddings (Māori-Eng-300) consistently perform better than monolingual embeddings. BiLSTM with Māori-Eng-300SG embeddings are the best across all models, including BERT-base and mBERT-base. Skip-gram models are better than CBOW. In comparison, English-only embeddings E300 outperform Māori-only monolingual embeddings. One possible explanation for this is the lack of training data for Māori-only embeddings compared to E300.

Model	Training data	Testing data	Accuracy (Māori)
Google	Wikipedia	RMT	68.2%
BiLSTM (E300)	Hansard	RMT	56.6%
BiLSTM (Māori-Eng-300)	Hansard	RMT	85.4%
BiLSTM (Māori-Eng-300SG)	Hansard	RMT	85.6%

Table 6: Accuracy of Māori words detection in RMT corpus using Hansard-based trained models (Table 5).

Figure 3 presents critical difference plots across the models presented in Table 5 and BOW representation presented in Table 4. BiLSTM (Māori-Eng-300SG) has the lowest rank, and multinomial naïve Bayes (BOW) has the highest rank with no bold line connecting the two, indicating the difference in average ranking is statistically significant. Bold lines are connecting BiLSTM (Māori-Eng-300SG) with mBERT and BERT-base in the CD-plot, indicating that the difference in average ranking is not statistically significant. A 4-6 % improvement was observed between BERT/mBERT and BiLSTM (Māori-Eng-300SG).

To further evaluate the language models, we used the models trained with the Hansard data to detect Māori words in RMT corpus. Table 6 presents the accuracy of the detection. We also present the accuracy of Māori language detection using Google Translate. Evidently, BiLSTM with Māori-Eng-300SG embeddings model trained on the training set of the Hansard database has the best accuracy. As observed with other databases, the accuracy of the bilingual embeddings is higher than the monolingual embeddings. However, the accuracy of BiLSTM with E300 embeddings is considerably lower than other models, including Google. One possible reason is the lack of vocabulary in E300 for the informal language used in RMT data (Tweets).

7.1.4 In Summary

The results suggest that the bilingual embeddings perform better than monolingual embeddings for

Model	CS: Accuracy
CNN (E300)	35%
BiLSTM (E300)	83%
BiLSTM (Māori-Eng-300)	67%
BiLSTM (Māori-Eng-300SG)	87%

Table 7: Accuracy of code-switch detection in the Hansard data (bilingual sentences of the test set) using the trained models, as shown in Figure 1.

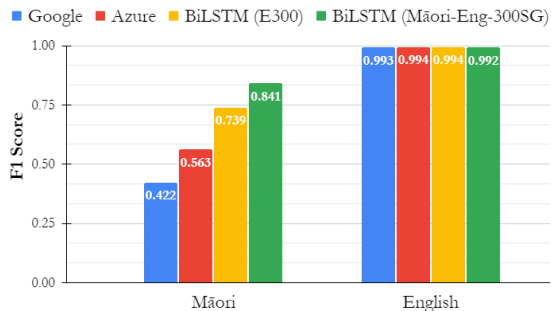


Figure 4: F1-scores for Māori and English calculated at the word level for the Hansard database.

the LD task. This finding was verified across the Hansard database (Tables 4, 5) and the MLT corpus (Table 4). Further evidence is provided in Māori words detection using RMT corpus (Table 6). We also observed that the bilingual embeddings outperformed the pre-trained contextual embeddings. One possible reason for this finding is the lack of vocabulary in BERT alike models as we did not perform any further training using Māori data. As emphasised before, the Māori data availability and access is the biggest limitation to this research. Among the experimented models for LD task, BiLSTM with Māori-Eng-300SG performed the best.

7.2 Task 2: Code-Switch Detection

For evaluation of the code-switch detection between Māori-English pair, we require word-level labels and hence, only use the test set of the Hansard database. We use selected trained models presented in Section 7.1, and identify the code-switch point (see Figure 1). Figure 4 presents word-level F1 scores of Māori and English for CS task. For English words, all systems perform equally well. However, for Māori, cloud-based systems perform poorly, and BiLSTM with bilingual embeddings shows a substantial improvement in F1 score, as observed before. Furthermore, Table 7 presents the accuracy of detecting the code-switch points of the test set of the Hansard database. Among the reported results, CNN with E300 performed poorly, and BiLSTM with Māori-Eng-300SG outperformed the other models.

8 Discussions and Conclusions

This research is the first attempt to use advances in NLP in two tasks - low-resourced Māori language detection and Māori-English code-switch detection. Our experiments show that the accuracy of existing cloud-based systems to detect Māori language is very low. We collect data in collaboration with Māori researchers for training and evaluations. Experiments obtained across tasks using three databases show that our bilingual embeddings outperformed English-only embeddings trained on large databases. Among the models tested in this research, BiLSTM with bilingual embeddings trained using the Skip-gram model is the best for both tasks. We provide evidence to show BERT-base used on the down-streaming task of language detection – where Māori is under-represented or unseen by the model vocabulary– is not always the best solution (as also observed by (Wu and Dredze, 2020; Wang et al., 2020)). For most low-resourced languages, including Māori, the Wikipedia data is significantly smaller than English, resulting in a fewer vocabulary. Due to lack of resources, continuous training or training from scratch of models such as BERT-base is not possible. For future work, it is a possibility to use ideas such as Extend M-BERT (Wang et al., 2020) and explore the possibility of using more efficient pre-training techniques to improve the accuracy of BERT like models for language detection of low-resource languages such as Māori. In addition, hybrid models using handcrafted rules based on the phonotactic differences between the languages and deep learning-based methods are a pathway for future work. It is vital to point out that the availability of digitised Māori and bilingual data is limited, which restricts the ability to train large language models. In addition, considering this is the first deep learning-based research in this area, comparison with published work is not possible. We overcome these limitations by respecting the available data and data sovereignty for this research, and we use the available cloud services as the baseline existing systems for comparisons. The study reported here is a much-wanted contribution to Māori language technology development. Word embeddings developed in this research are available to other researchers on request, bound by the Kaitiakitanga license.

Acknowledgements

Thanks to all researchers who shared their data collections with us, and to the xx Fund for support.

References

New Zealand Parliament Hansard. Hansard reports. 670

New Zealand. [https://www.parliament.](https://www.parliament.nz/en/pb/hansard-debates/rhr/) 671

[nz/en/pb/hansard-debates/rhr/](https://www.parliament.nz/en/pb/hansard-debates/rhr/). 672

Ray Harlow. 2007. *Maori: A linguistic introduction*. 673

Cambridge University Press. 674

Sepp Hochreiter and Jürgen Schmidhuber. 1997. Long 675

short-term memory. *Neural computation*, 9(8):1735– 676

1780. 677

Jesin James, Isabella Shields, Rebekah Berriman, Pe- 678

ter J. Keegan, and Catherine I. Watson. 2020. Devel- 679

oping Resources for Te Reo Māori Text To Speech 680

Synthesis System. In *Proc. Sojka P., Kopeček I., Pala* 681

K., Horák A. (eds) Text, Speech, and Dialogue, Lec- 682

ture Notes in Computer Science, pages 294–302. 683

George H John and Pat Langley. 1995. Estimating con- 684

tinuous distributions in bayesian classifiers. In *Proc.* 685

Conference on Uncertainty in Artificial Intelligence, 686

pages 338–345, San Mateo. Morgan Kaufmann. 687

Armand Joulin, Edouard Grave, Piotr Bojanowski, 688

Matthijs Douze, Hérve Jégou, and Tomas Mikolov. 689

2016a. Fasttext. zip: Compressing text classification 690

models. *arXiv preprint arXiv:1612.03651*. 691

Armand Joulin, Edouard Grave, Piotr Bojanowski, and 692

Tomas Mikolov. 2016b. Bag of tricks for efficient 693

text classification. *arXiv preprint arXiv:1607.01759*. 694

Te Taka Keegan. 2021. Private collection of te reo Māori 695

text data. *The Univeristy of Waikato, New Zealand*. 696

Te Taka Adrian Gregory Keegan. 2017. Machine trans- 697

lation for te reo māori. 698

Simran Khanuja, Sandipan Dandapat, Anirudh Srini- 699

vasan, Sunayana Sitaram, and Monojit Choudhury. 700

2020. [GLUECoS: An evaluation benchmark for](#) 701

[code-switched NLP](#). In *Proc. of the 58th Annual* 702

Meeting of the Association for Computational Lin- 703

guistics, pages 3575–3585, Online. Association for 704

Computational Linguistics. 705

Yoon Kim. 2014. Convolutional neural networks for 706

sentence classification. In *Proc. of the 2014 Con-* 707

ference on Empirical Methods in Natural Language 708

Processing (EMNLP), pages 1746–1751. Association 709

for Computational Linguistics. 710

Diederik P Kingma and Jimmy Ba. 2015. Adam: A 711

method for stochastic optimization. In *Proc. Inter-* 712

national Conference on Learning Representations 713

(ICLR), pages 1–13. 714

Yash Kumar Lal, Vaibhav Kumar, Mrinal Dhar, Manish 715

Shrivastava, and Philipp Koehn. 2019. De-mixing 716

sentiment from code-mixed text. In *Proc. of the 57th* 717

Annual Meeting of the Association for Computational 718

Linguistics: Student Research Workshop, pages 371– 719

377. 720

Wafia Adouane, Jean-Philippe Bernardy, and Simon 618
Dobnik. 2018. Improving neural network perfor- 619
mance by injecting background knowledge: Detect- 620
ing code-switching and borrowing in algerian texts. 621
In *Proc. of the Third Workshop on Computational Ap-* 622
proaches to Linguistic Code-Switching, pages 20–28. 623

Utsab Barman, Amitava Das, Joachim Wagner, and Jen- 624
nifer Foster. 2014. Code mixing: A challenge for 625
language identification in the language of social me- 626
dia. In *Proc. of the first workshop on computational* 627
approaches to code switching, pages 13–23. 628

Winifred Bauer, William Parker, and Te Kareongawai 629
Evans. 1993. *Māori*. London: Routledge. 630

Google AI Blog. Google ai blog: Recent advances in 631
google translate. 632

Piotr Bojanowski, Edouard Grave, Armand Joulin, and 633
Tomas Mikolov. 2017. Enriching word vectors with 634
subword information. *Transactions of the Associa-* 635
tion for Computational Linguistics, 5:135–146. 636

Alexis Conneau, Kartikay Khandelwal, Naman Goyal, 637
Vishrav Chaudhary, Guillaume Wenzek, Francisco 638
Guzmán, Edouard Grave, Myle Ott, Luke Zettle- 639
moyer, and Veselin Stoyanov. 2020. Unsupervised 640
cross-lingual representation learning at scale. In *Proc.* 641
ACL Conference. 642

Alexis Conneau, Guillaume Lample, Ruty Rinott, Adina 643
Williams, Samuel R Bowman, Holger Schwenk, and 644
Veselin Stoyanov. 2018. XNLI: Evaluating Cross- 645
lingual Sentence Representations. In *Proc. EMNLP*, 646
pages 2475–2485. 647

David R Cox. 1958. The regression analysis of binary 648
sequences. *Journal of the Royal Statistical Society:* 649
Series B (Methodological), 20(2):215–232. 650

Jacob Devlin, Ming-Wei Chang, Kenton Lee, and 651
Kristina Toutanova. 2019. [BERT: Pre-training of](#) 652
[deep bidirectional transformers for language under-](#) 653
[standing](#). In *Proc. of the 2019 Conference of the* 654
North American Chapter of the Association for Com- 655
putational Linguistics: Human Language Technolo- 656
gies, Volume 1 (Long and Short Papers), pages 4171– 657
4186, Minneapolis, Minnesota. Association for Com- 658
putational Linguistics. 659

Yoav Goldberg. 2017. Neural network methods for 660
natural language processing. *Synthesis Lectures on* 661
Human Language Technologies, 10(1):1–309. 662

Google. Google translate. [https://translate.](https://translate.google.com/) 663
[google.com/](https://translate.google.com/). 664

Edouard Grave, Piotr Bojanowski, Prakhar Gupta, Ar- 665
mand Joulin, and Tomas Mikolov. 2018. Learning 666
word vectors for 157 languages. In *Proc. of the Inter-* 667
national Conference on Language Resources and 668
Evaluation, pages 3483–3487. 669

721	LMC. Legal Māori corpus. <i>Victoria University of Wellington, New Zealand</i> . http://nzetc.victoria.ac.nz/tm/scholarly/tei-legalMaoriCorpus.html .	775
722		776
723		777
724		
725	Media. Xx. Name avoided for anonymity.	
726	Microsoft. Azure translator. https://azure.microsoft.com/en-us/services/cognitive-services/translator/ .	
727		
728		
729	Tomas Mikolov, Kai Chen, Greg Corrado, and Jeffrey Dean. 2013. Efficient estimation of word representations in vector space. <i>arXiv preprint arXiv:1301.3781</i> .	
730		
731		
732		
733	Mahsa Mohaghegh, Michael McCauley, and Mehdi Mohammadi. 2014. Māori-english machine translation. <i>NZCSRSC New Zealand Computer Science Research Student Conference-Canterbury University. Unitec Research Bank</i> .	
734		
735		
736		
737		
738	Siddhartha Mukherjee, Vinuthkumar Prasan, Anish Nediyanath, Manan Shah, and Nikhil Kumar. 2019. Robust deep learning based sentiment classification of code-mixed text. In <i>Proc. of the 16th International Conference on Natural Language Processing</i> , pages 124–129.	
739		
740		
741		
742		
743		
744	Benjamin Muller, Antonios Anastasopoulos, Benoît Sagot, and Djamé Seddah. 2021. When being unseen from mBERT is just the beginning: Handling new languages with multilingual language models . In <i>Proc. of the 2021 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies</i> , pages 448–462, Online. Association for Computational Linguistics.	
745		
746		
747		
748		
749		
750		
751		
752		
753	Niuepepa. Maori newspapers - new zealand digital library. <i>Ministry of Education, New Zealand</i> . http://www.nzdl.org/cgi-bin/library.cgi?a=p&p=about&c=niuepepa .	
754		
755		
756		
757	Shruti Rijhwani, Royal Sequiera, Monojit Choudhury, Kalika Bali, and Chandra Shekhar Maddila. 2017. Estimating code-switching on twitter with a novel generalized word-level language detection technique. In <i>Proc. of the 55th annual meeting of the association for computational linguistics (volume 1: long papers)</i> , pages 1971–1982.	
758		
759		
760		
761		
762		
763		
764	David E Rumelhart, Geoffrey E Hinton, and Ronald J Williams. 1986. Learning representations by back-propagating errors. <i>Nature</i> , 323(6088):533–536.	
765		
766		
767	Scikit-Learn. Sklearn’s principal component analysis. https://scikit-learn.org/stable/modules/generated/ , obtained 10 Nov 2021.	
768		
769		
770	Nitish Srivastava, Geoffrey Hinton, Alex Krizhevsky, Ilya Sutskever, and Ruslan Salakhutdinov. 2014. Dropout: A simple way to prevent neural networks from overfitting . <i>Journal of Machine Learning Research</i> , 15(56):1929–1958.	
771		
772		
773		
774		
	N. Z. Stats. 2020. Ngā tikanga paihere: a framework guiding ethical and culturally appropriate data use. <i>Guidelines</i> , page 8.	775
		776
		777
	Te-Hiku-Media. Kaitiakitanga license. https://github.com/TeHikuMedia/Kaitiakitanga-License , accessed 10 Dec 2021.	778
		779
		780
		781
	Mads Tofttrup, Søren Asger Sørensen, Manuel R. Ciosici, and Ira Assent. 2021. A reproduction of apple’s bi-directional LSTM models for language identification in short strings . In <i>Proc. of the 16th Conference of the European Chapter of the Association for Computational Linguistics: Student Research Workshop</i> , pages 36–42, Online. Association for Computational Linguistics.	782
		783
		784
		785
		786
		787
		788
		789
	David Trye, Andreea S Calude, Felipe Bravo-Marquez, and Te Taka Adrian Gregory Keegan. 2019. Māori loanwords: a corpus of new zealand english tweets . In <i>Proc. of the 57th Annual Meeting of the Association for Computational Linguistics: Student Research Workshop, Florence, Italy: Association for Computational Linguistics.</i> , pages 136–142.	790
		791
		792
		793
		794
		795
		796
	David Trye, Te Taka Keegan, Paora Mato, and Mark Apperley. 2022. Harnessing indigenous tweets: The reo māori twitter corpus. <i>Language Resources and Evaluation</i> . (in press).	797
		798
		799
		800
	Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz Kaiser, and Illia Polosukhin. 2017. Attention is all you need. <i>Advances in neural information processing systems</i> , 30:5998–6008.	801
		802
		803
		804
		805
	Zihan Wang, K Karthikeyan, Stephen Mayhew, and Dan Roth. 2020. Extending multilingual bert to low-resource languages. In <i>Proc. of the 2020 Conference on Empirical Methods in Natural Language Processing: Findings</i> , pages 2649–2656.	806
		807
		808
		809
		810
	Shijie Wu and Mark Dredze. 2020. Are all languages created equal in multilingual BERT? In <i>Proc. of the 5th Workshop on Representation Learning for NLP</i> , pages 120–130, Online. Association for Computational Linguistics.	811
		812
		813
		814
		815
	Yonghui Wu, Mike Schuster, Zhifeng Chen, Quoc V Le, Mohammad Norouzi, Wolfgang Macherey, Maxim Krikun, Yuan Cao, Qin Gao, Klaus Macherey, et al. 2016. Google’s neural machine translation system: Bridging the gap between human and machine translation. <i>arXiv preprint arXiv:1609.08144</i> .	816
		817
		818
		819
		820
		821
	Wayne Xiong, Jasha Droppo, Xuedong Huang, Frank Seide, Mike Seltzer, Andreas Stolcke, Dong Yu, and Geoffrey Zweig. 2017. The microsoft 2016 conversational speech recognition system. In <i>Proc. Acoustics, Speech and Signal Processing (ICASSP), 2017 IEEE International Conference on</i> , pages 5255–5259. IEEE.	822
		823
		824
		825
		826
		827
		828

829 Chiyuan Zhang, Samy Bengio, Moritz Hardt, Benjamin
830 Recht, and Oriol Vinyals. 2017. [Understanding deep](#)
831 [learning requires re-thinking generalization](#). In *Proc.*
832 *International Conference on Learning Representa-*
833 *tions 2017*, pages 1–15.