

# 论文分享

ZiXi Yan

2022/06/10

# TiBERT: Tibetan Pre-trained Language Model

- To better express the semantic information of Tibetan and reduce the problem of OOV, this paper uses the unigram language model of Sentencepiece to segment Tibetan words and constructs a vocabulary that can cover 99.95% of the words in the corpus.
- To further promote the development of various downstream tasks of Tibetan natural language processing, this paper collected a large-scale Tibetan dataset and trained the monolingual Tibetan pre-trained language model named TiBERT.
- To evaluate the performance of TiBERT, this paper conducts comparative experiments on the two downstream tasks of text classification and question generation. The experimental results show that the TiBERT is effective.

# TiBERT: Tibetan Pre-trained Language Model

- Tibetan data from 21 Tibetan websites including Tibet People's Network and Qinghai Provincial People's Government Network
- The data contains knowledge in various fields such as current affairs, economy, technology, society, law, sports, life, nature, culture, geography, art, military, education, history, and people

# TiBERT: Tibetan Pre-trained Language Model

TABLE IV  
PERFORMANCES ON DOCUMENT CLASSIFICATION

Model	Accuracy(%)	Macro-Precision(%)	Macro-Recall(%)	Macro-F1(%)
CNN(syllable)	61.51	59.39	56.65	57.34
CINO-large	-	-	-	68.6
Transformer	28.63	41.21	28.63	28.79
TextCNN	61.71	61.65	61.71	61.53
DPCNN	62.91	63.61	62.91	61.17
TextRCNN	63.67	64.37	63.67	62.81
TiBERT	<b>71.04</b>	<b>71.20</b>	<b>71.04</b>	<b>70.94</b>
TiBERT+CNN	70.39	70.54	70.39	70.23

# TiBERT: Tibetan Pre-trained Language Model

	Gold questions	S2S+ATT+CP	TiBERT
1	<p>རང་བྱུང་ཁམས་ཀྱི་སྐྱེག་ཕྱིའི་འོད་ཀྱི་འབྱུང་ཁུངས་ནི་ཅི་ ཞིག་ཡིན།</p> <p>What is the source of ultraviolet light in nature?</p>	<p>རང་བྱུང་ཁམས་ཀྱི་སྐྱེག་ཕྱིའི་འོད་ཀྱི་འབྱུང་ཁུངས་ནི་གང་ ཡིན།</p> <p>(What is the source of ultraviolet light in nature?)</p>	<p>རང་བྱུང་ཁམས་ཀྱི་སྐྱེག་ཕྱིའི་འོད་ཀྱི་འབྱུང་ཁུངས་ནི་གང་ཡི ན།</p> <p>(What is the source of ultraviolet light in nature?)</p>
2	<p>མངའ་ཁོངས་ཀྱི་ཆ་སྟོམས་མཚོ་ངོས་མཐོ་ཚད་སྤེལ་ག་ཚོད་ ཡིན།</p> <p>(What is the average elevation of the territory in meters?)</p>	<p>མངའ་ཁོངས་ཀྱི་ཆ་སྟོམས་མཚོ་ངོས་མཚོ་ངོས་མཐོ་ཚད་ ག་ཚོད་ ཡིན།</p> <p>(What is the average <b>elevation elevation</b> of the territory?)</p>	<p>མངའ་ཁོངས་ཀྱི་ཆ་སྟོམས་མཚོ་ངོས་མཐོ་ཚད་སྤེལ་ག་ཚོད་ཡོད།</p> <p>(What is the average elevation of the territory in meters?)</p>
3	<p>དགོན་གོང་མ་གོང་རྩལ་སྲིད་གཞུང་གི་སྟོ་རྩལ་ཕྱོགས་སུ་ ཉ་ལམ་སྤེལ་ག་ཚོད་ཡོད།</p> <p>(How many kilometers southwest of Gongma town government?)</p>	<p>དགོན་གོང་མ་གོང་རྩལ་སྲིད་གཞུང་གི་སྟོ་རྩལ་ཕྱོགས་ ཞུང་གི་སྟོ་རྩལ་ཕྱོགས་སུ་ཉ་ལམ་སྤེལ་ག་ཚོད་ཡོད།</p> <p>(How many centimeters is the <b>southwest southwest</b> of Gongma town government?)</p>	<p>དགོན་གོང་མ་གོང་རྩལ་སྲིད་གཞུང་གི་སྟོ་རྩལ་ཕྱོགས་སུ་ཉ་ལ མ་སྤེལ་ཅི་ཙམ་མཚམས་ན་ཡོད།</p> <p>(How many kilometers southwest of Gongma town government?)</p>
4	<p>ངལ་ཚོལ་རུས་ཤུགས་ཅན་གྱི་མི་ག་ཚོད་ཡོད།</p> <p>(How many people have the ability of work?)</p>	<p><b>གོང་ཚོ་ཡོངས་ལ་སྤེལ་ག་ཚོད་ཡོད།</b></p> <p>(How many villages are <b>there in the whole town?</b>)</p>	<p>ངལ་ཚོལ་རུས་ཤུགས་ཅན་གྱི་མི་ག་ཚོད་གི་མི་ག་ཚོད་ཡོད།</p> <p>(How <b>many many</b> people have the ability of work?)</p>
5	<p>དུས་ནམ་ཞིག་ལ་རྒྱུ་ཞེང་གྲང་རྩེ་ལ་བསྐྱར།</p> <p>(When will it be changed to Red Star Commune?)</p>	<p>དུས་ནམ་ཞིག་ལ་གྲང་རྩེ་གྲང་རྩེ་གྲང་རྩེ་ལ་བསྐྱར།</p> <p>(When will it be changed to a <b>commune commune</b> <b>commune?</b>)</p>	<p>སྤྱི་ལོ་<b>1966</b>ལོར་རྒྱུ་ཞེང་གྲང་རྩེ་ལ་བསྐྱར།</p> <p>(Changed to Red Star Commune in <b>1966</b>)</p>

# Quantifying Language Variation Acoustically with Few Resources

- 10 words (armen: 'arms' ,deeg: 'dough' ,draden: 'wires' , duiven: 'pigeons' ,naalden: 'needles' ,ogen: 'eyes' , pijpen: 'pipes' ,tangen: 'pliers' ,volk: 'people' ,vuur: 'fire' )
- pronounced in 106 locations in the Netherlands. On average, the duration of these 10 words is only 6.3 seconds for each location.

# Quantifying Language Variation Acoustically with Few Resources

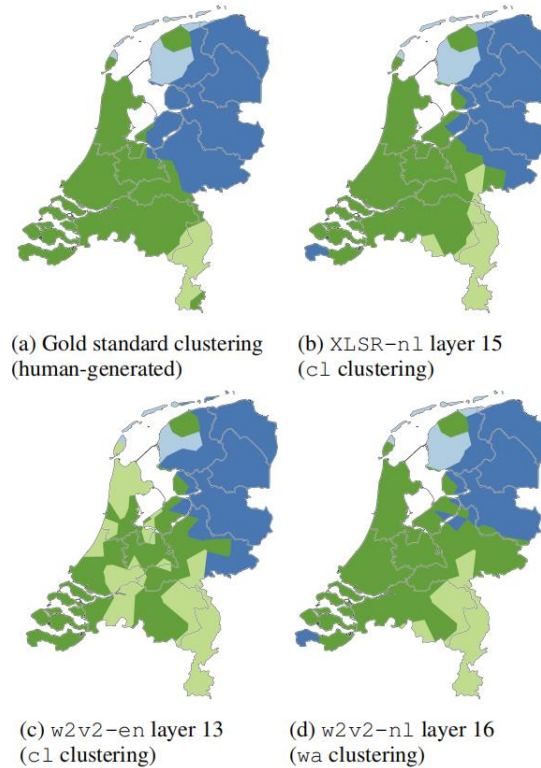


Figure 1: Cluster maps visualizing four clusters on the map of the Netherlands. Separate clusters are indicated by the different colours.

# Automatic Pronunciation Assessment using Self-Supervised Speech Representation Learning

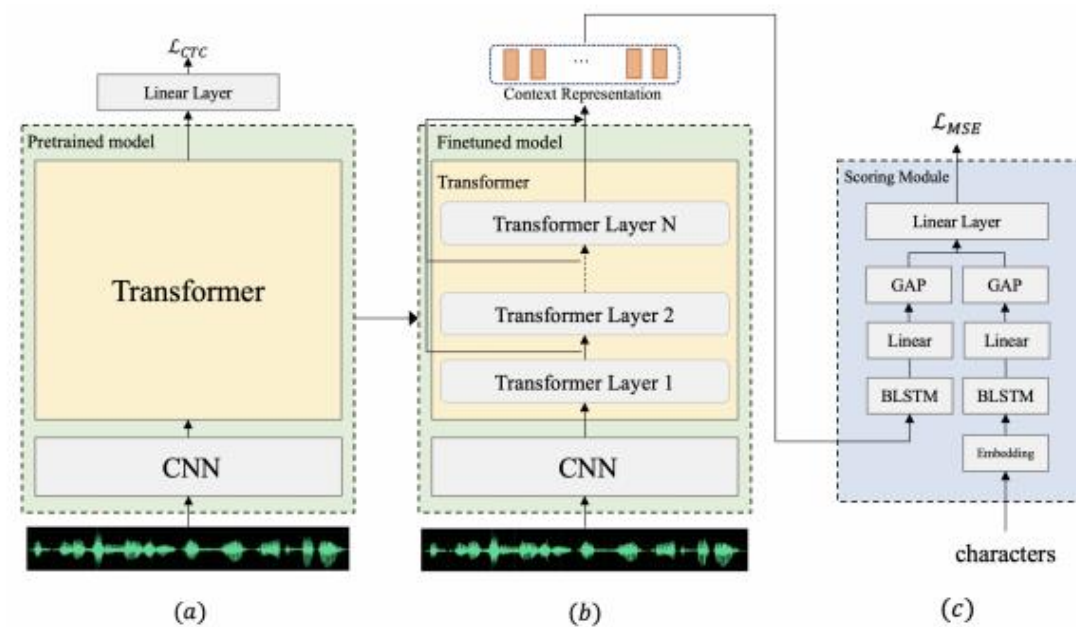


Figure 1: Overall procedure of the proposed method for automatic pronunciation assessment based on SSL models.



# Automatic Pronunciation Assessment using Self-Supervised Speech Representation Learning

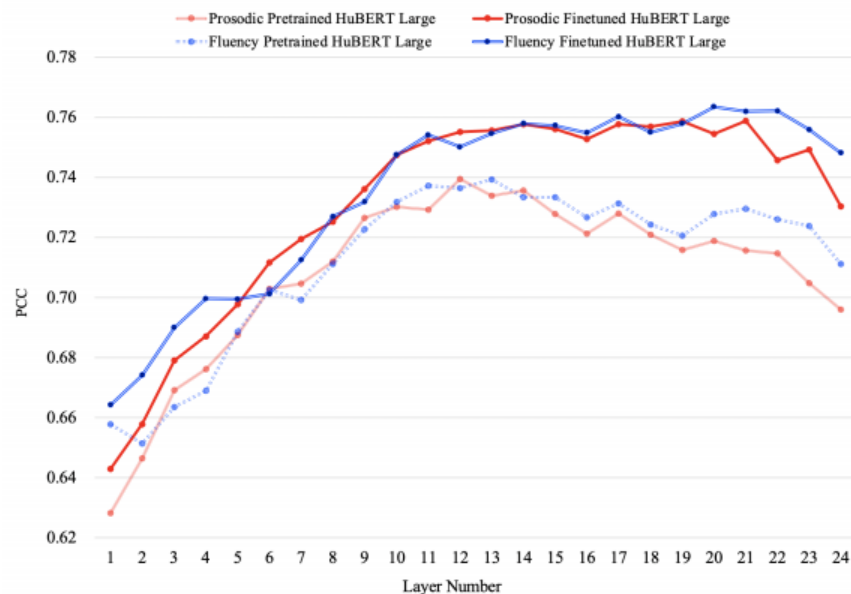


Figure 3: PCCs when using the hidden states of different transformer layers of the fine-tuned HuBERT Large model as input to the BLSTM scoring module.

Feature	KESL	Speechocean762	
	Holistic	Fluency	Prosodic
Local	0.56	0.60	0.62
Layer 20	0.81	0.76	0.76
All Layers (Proposed)	<b>0.82</b>	<b>0.78</b>	<b>0.77</b>

Table 2: Comparison of the performance of local representation of a convolutional layer, the contextual representation, and layer-wise contextual representation of the transformer layers in HuBERT Large model.

# CATEGORICAL REPARAMETERIZATION WITH GUMBEL-SOFTMAX

The Gumbel-Max trick (Gumbel, 1954; Maddison et al., 2014) provides a simple and efficient way to draw samples  $z$  from a categorical distribution with class probabilities  $\pi$ :

$$z = \text{one\_hot} \left( \arg \max_i [g_i + \log \pi_i] \right) \quad (1)$$

where  $g_1 \dots g_k$  are i.i.d samples drawn from  $\text{Gumbel}(0, 1)$ . We use the softmax function as a continuous, differentiable approximation to  $\arg \max$ , and generate  $k$ -dimensional sample vectors  $y \in \Delta^{k-1}$  where

$$y_i = \frac{\exp((\log(\pi_i) + g_i)/\tau)}{\sum_{j=1}^k \exp((\log(\pi_j) + g_j)/\tau)} \quad \text{for } i = 1, \dots, k. \quad (2)$$

The density of the Gumbel-Softmax distribution (derived in Appendix B) is:

$$p_{\pi, \tau}(y_1, \dots, y_k) = \Gamma(k) \tau^{k-1} \left( \sum_{i=1}^k \pi_i / y_i^\tau \right)^{-k} \prod_{i=1}^k (\pi_i / y_i^{\tau+1}) \quad (3)$$

This distribution was independently discovered by Maddison et al. (2016), where it is referred to as the concrete distribution. As the softmax temperature  $\tau$  approaches 0, samples from the Gumbel-Softmax distribution become one-hot and the Gumbel-Softmax distribution becomes identical to the categorical distribution  $p(z)$ .

# CATEGORICAL REPARAMETERIZATION WITH GUMBEL-SOFTMAX

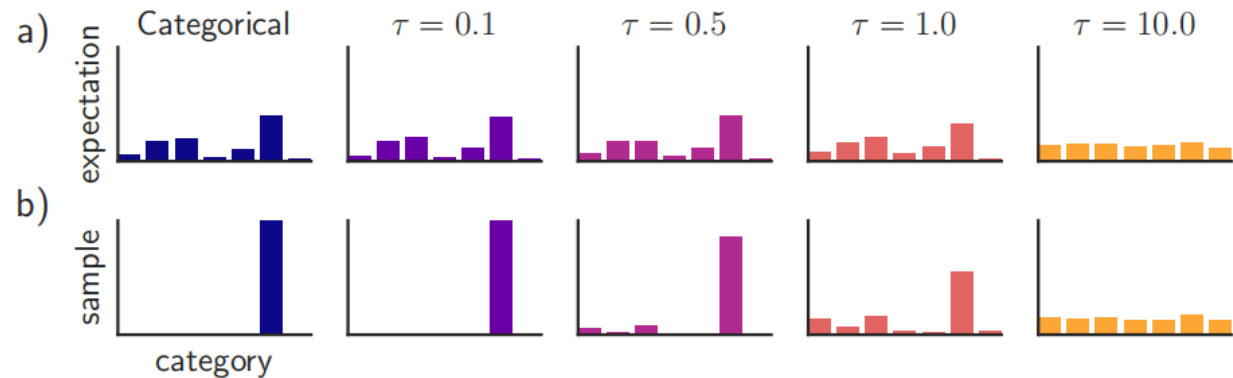


Figure 1: The Gumbel-Softmax distribution interpolates between discrete one-hot-encoded categorical distributions and continuous categorical densities. (a) For low temperatures ( $\tau = 0.1, \tau = 0.5$ ), the expected value of a Gumbel-Softmax random variable approaches the expected value of a categorical random variable with the same logits. As the temperature increases ( $\tau = 1.0, \tau = 10.0$ ), the expected value converges to a uniform distribution over the categories. (b) Samples from Gumbel-Softmax distributions are identical to samples from a categorical distribution as  $\tau \rightarrow 0$ . At higher temperatures, Gumbel-Softmax samples are no longer one-hot, and become uniform as  $\tau \rightarrow \infty$ .

# CATEGORICAL REPARAMETERIZATION WITH GUMBEL-SOFTMAX

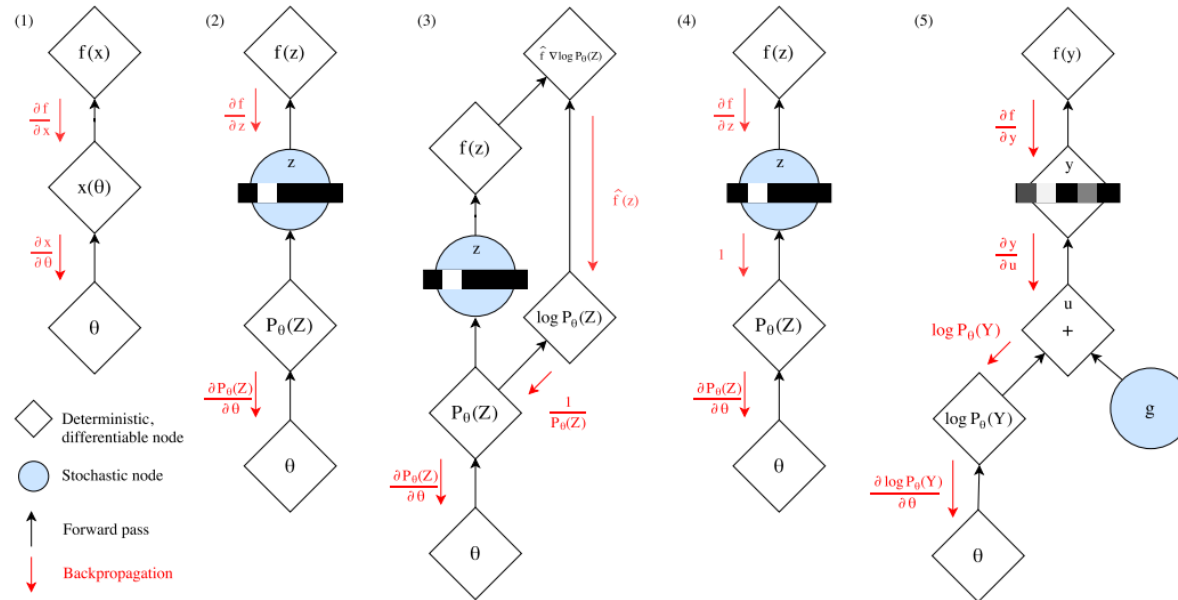


Figure 2: Gradient estimation in stochastic computation graphs. (1)  $\nabla_{\theta} f(x)$  can be computed via backpropagation if  $x(\theta)$  is deterministic and differentiable. (2) The presence of stochastic node  $z$  precludes backpropagation as the sampler function does not have a well-defined gradient. (3) The score function estimator and its variants (NVIL, DARN, MuProp, VIMCO) obtain an unbiased estimate of  $\nabla_{\theta} f(x)$  by backpropagating along a surrogate loss  $\hat{f} \log p_{\theta}(z)$ , where  $\hat{f} = f(x) - b$  and  $b$  is a baseline for variance reduction. (4) The Straight-Through estimator, developed primarily for Bernoulli variables, approximates  $\nabla_{\theta} z \approx 1$ . (5) Gumbel-Softmax is a path derivative estimator for a continuous distribution  $y$  that approximates  $z$ . Reparameterization allows gradients to flow from  $f(y)$  to  $\theta$ .  $y$  can be annealed to one-hot categorical variables over the course of training.

Thanks