# **CORI:** CJKV Benchmark with Romanization Integration - A step towards Cross-lingual Transfer Beyond Textual Scripts

### Hoang Nguyen<sup>1</sup>, Chenwei Zhang<sup>2</sup>, <u>Ye Liu<sup>3</sup></u>, Natalie Parde<sup>1</sup>, Eugene Rohrbaugh<sup>4</sup>, Philip S. Yu<sup>1</sup>

<sup>1</sup> University of Illinois at Chicago <sup>2</sup> Amazon <sup>3</sup> Salesforce Research

<sup>3</sup> Harrisburg University of Science and Technology





# Overview

- Background
- Preliminary Study
- Dataset Construction
- Framework
- Evaluation
- Conclusion

# BACKGROUND

# Background



Performance Comparison between languages on downstream task when fine-tuning on EN



 Performance on degrades language with different script



4

### Q1: Should EN always be chosen as source language?

# Q2: How to overcome performance gap between languages with different textual writing scripts

Results are obtained from Yang et al., 2022. Enhancing Cross-lingual Transfer via Manifold Mixup. ICLR 2022

# PRELIMINARY STUDY

# Background

### **Observation:**

When ZH is leveraged as source language:

- Higher average performance among target JKV languages
- Lower STD among JKV

 $\Rightarrow$  ZH is a better source language than EN when considering target JKV languages



#### Average performance across JKV languages when using different source languages

### Source language matters on downstream task target language performance

# Background

### **Observation:**

- Representations of parallel inputs between JA and ZH are more aligned than between EN and ZH
- JA-ZH has higher language contact
- ⇒ ZH is a better source language than EN when considering target JKV languages



n Figure 1: Representation visualization of parallel sentences between source language and target language (EN-JA (left) and ZH-JA (right)) of the fine-tuned XLM-R model on parallel PAWSX test set. The significant overlapping representation of ZH-JA demonstrates the lower representation discrepancy and higher language contact between source and target language.

Closely-related benchmark dataset is needed CJKV



**CORI Dataset Construction** 

# DATASET CONSTRUCTION

## **Dataset Construction**

### **<u>Challenge 1:</u>** Language Availability

Table 4: Details of CORI benchmark dataset. (X) and ( $\checkmark$ ) denote unavailable and available data existent in XTREME benchmark respectively. MT, SEG, ROM correspond to Machine Translation, Presegmentation, Romanization processing respectively as described in Section 5. Y, N denote if the corresponding preprocessing is conducted for the dataset or not.

		MT	SEG	ROM	ZH		JA	KO	VI
					Train	Dev	Test	Test	Test
Sent-level	PAWSX	Y	Y	Y	49.4k ( <b>X</b> )	2k (🗸)	2k ( <b>x</b> )	2k (✔)	2k ( <b>X</b> )
	XNLI	Y	Y	Y	392.7k (🗸)	2.49k (🗸)	5.01k ( <b>x</b> )	5.01k ( <b>x</b> )	5.01k (🗸)
Token-level	UDPOS	N	Ν	Y	10k (🗸)	1.6k (🗸)	2.5k (✔)	4.7k(✓)	0.8k (🗸 )
	PANX	N	Y	Y	20k (🗸)	10k (🗸)	10k (✔)	10k (🗸)	10k (🗸)
Question Answering	XQuAD	Y	Y	Y	80.1k ( <b>x</b> )	8.87k ( <b>x</b> )	1.19k ( <b>X</b> )	1.19k ( <b>X</b> )	1.19k (✔)
	MLQA	Y	Y	Y	80.1k ( <b>X</b> )	8.87k ( <b>x</b> )	11.24 k ( <b>X</b> )	11.24k ( <b>x</b> )	11.24k (🗸)

### Multilingual Datasets have been created unequally across CJKV languages

## **Dataset Construction**

### **Challenge 2: Pre-segmentation**



### Inconsistent pre-segmentation across CJKV languages existent in XTREME

### Dataset Construction <u>Challenge 3:</u> Romanization

Table 1: Orthographic and Romanized representations (abbreviated as Ortho and Roman) of a sample sentence across CJKV languages where **colored segments** denote the corresponding semantic **words** defined in Section 5. (.) denotes the specific name of Romanization system of the respective language. The first sentence for each language denotes the currently preprocessed XTREME benchmark dataset.

Language	Input Type	Sample input sentence								
EN	Ortho	He was a scholar in Metaphysical Literature , Theology and Classical sciences .								
	Ortho	他是形而上学文学、神学和古典科学方面的学者。								
ZH (source)	Ortho (seg)	他 // 是 // 形而上学 // 文学 // 、// 神学 // 和 // 古典 // 科学 // 方面 // 的 // 学者 // 。								
	Roman (Pinyin)	tā // shì // xíngérshàngxué // wénxué // 、 // shénxué // hé / gǔdiǎn // kēxué // fāngmiàn // de // xuézhě // 。								
	Ortho	Ông là một học giả về Văn học Siêu hình , Thần học và Khoa học Cổ điển .								
VI (target)	Ortho (seg)	Ông // là // một // học giả // về // Văn học // Siêu hình // , // Thần học // và // Khoa học // Cổ điển // .								
	Roman	Ông // là // một // học giả // về // Văn học // Siêu hình // , // Thần học // và <mark>/</mark> Khoa học /- Cổ điển // .								
	Ortho	彼は形而上学文学、神学、古典科学の学者でした。								
JA (target)	Ortho (seg)	彼 // は // 形而上学 // 文学 //、神学 、// 古典 // 科学 //の// 学者 //でし // た // 。								
	Roman (Romaji)	kare // ha // keiji ue gaku // bungaku // ,// shingaku // ,// koten // kagaku // no // gakusha // deshi // ta // .								
	Ortho	그는형이상학문학, 신학및고전과학의학자이었습니다.								
KO (target)	Ortho (seg)	그〃는〃 형이상학 〃 문학 〃,〃 <mark>신학</mark> 〃및〃 고전 〃 과학 〃 의학자 〃 <u>이</u> 〃었〃습니다〃.								
	Roman (Romaja)	eu // neun // hyeongisanghak // munhak // ,// sinhak // mit // gojeon // gwahak // uihakja // i // eot // seupnida // .								

Romanization captures linguistic contact beyond textual scripts, providing beneficial signals for cross-lingual transfer

## **Dataset Construction**



## **Dataset Construction**

Table 4: Details of CORI benchmark dataset. (X) and ( $\checkmark$ ) denote unavailable and available data existent in XTREME benchmark respectively. MT, SEG, ROM correspond to Machine Translation, Presegmentation, Romanization processing respectively as described in Section 5. Y, N denote if the corresponding preprocessing is conducted for the dataset or not.

		MT	SEG	ROM	ZH		JA	КО	VI
					Train	Dev	Test	Test	Test
Sent-level	PAWSX	Y	Y	Y	49.4k ( <b>X</b> )	2k (🗸)	2k ( <b>X</b> )	2k (🗸)	2k ( <b>x</b> )
	XNLI	Y	Y	Υ	392.7k (🗸)	2.49k (🗸)	5.01k ( <b>X</b> )	5.01k ( <b>x</b> )	5.01k (🗸)
Token-level	UDPOS	N	Ν	Y	10k (🗸)	1.6k (🗸)	2.5k (✔)	4.7k(✔)	0.8k (✔)
	PANX	N	Y	Y	20k (🗸)	10k (🗸)	10k (🗸)	10k (🗸)	10k (🖌)
Question Answering	XQuAD	Y	Y	Y	80.1k ( <b>X</b> )	8.87k ( <b>x</b> )	1.19k ( <b>X</b> )	1.19k ( <b>X</b> )	1.19k (✔)
	MLQA	Y	Y	Υ	80.1k ( <b>X</b> )	8.87k ( <b>x</b> )	11.24 k ( <b>x</b> )	11.24k ( <b>X</b> )	11.24k (🗸)

CORI addresses the presented challenges from multilingual XTREME benchmark

- 1. Language Availability
- 2. Pre-segmentation
- 3. Romanization

# FRAMEWORK

## Framework: Overview



# **EVALUATION**

# Evaluation

#### **Observation:**

- Higher performance on CORI than XTREME benchmark dataset
- No difference in performance on UDPOS task performance since no MT or SEG is applied to improve quality of the dataset.

90 80 70 60 Performance 50 XTREME 40 CORI 30 20 10 0 PANX XquAD PAWSX XNLL UDPOS MLQA Tasks

#### JKV's Performance Comparison between CORI and XTREME (original) datasets

### Proposed preprocessing steps for CORI are effective

# **Evaluation**

### **Observation:**

- Romanized transcription enhances the textual representation across CJKV languages
- Leading to improvements of downstream multi-level tasks across target JKV languages

Romanization provides
additional helpful signals for cross-lingual transfer









Average JKV performance across multiple-level NLU tasks



# Romanization is an essential addendum, not replacement, for the orthographic representation

### **Observation:**

Evaluation

- Romanization provides helpful information for downstream tasks for target JKV languages
- Relying purely on Romanization is not sufficient for cross-lingual transfer on text-based LMs

#### Empirical Study on the impact of Romanization on representative multi-level NLU tasks



# CONCLUSION

# CONCLUSION

- Choice of source language is essential to downstream task performance for target languages
- <u>CORI:</u> CJKV-specific dataset addresses the limitations of current multilingual benchmark datasets
- **Romanization**, one type of phonemic signals, is valuable for cross-lingual transfer beyond the limitations of textual scripts

# Thank you for your attendance