

---

# Appendix

---

**Anonymous Author(s)**

Affiliation  
Address  
email

**1 A Training Corpora and Hyper-parameters****2 A.1 Training Corpora**

3 As for monolingual data, we follow Conneau et al. [2020] to build a Common-Crawl Corpus using  
4 the CCNet [Wenzek et al., 2020] tool<sup>1</sup>, which is widely used in the literature Huang et al. [2019],  
5 Luo et al. [2021], Chi et al. [2021], Wei et al. [2021]. Further, we collect parallel corpora from  
6 CCAligned El-Kishky et al. [2020], CCMATRIX Schwenk et al. [2021], WMT Akbardeh et al. [2021],  
7 and MultiUN Ziemski et al. [2016], involving 94 languages with more than 4.8 billion sentence  
8 pairs. We use the OpusFilter<sup>2</sup> tool to remove noisy bitexts, which results in 3.2 billion sentence pairs.  
9 Table 1 shows the statistics for both monolingual and parallel data. We apply subword tokenization  
10 directly on raw text data using Sentence Piece Model Kudo and Richardson [2018] without any  
11 additional preprocessing. To better support our motivation that EMMA-X can cover more languages  
12 than previous cross-lingual sentence representations, we divide Tatoeba Artetxe and Schwenk [2019]  
13 into two subsets: “Head”, containing languages usually covered in previous methods, and “Long-tail”,  
14 with other languages. We treat the 36 languages containing in XTREME Ruder et al. [2021] as head  
15 languages, which are: “**ar, he, vi, id, jv, tl, eu, ml, ta, te, af, nl, en, de, el, bn, hi, mr, ur, fa, fr, it, pt,**  
16 **es, bg, ru, ja, ka, ko, th, sw, zh, kk, tr, et, fi, hu, az, lt, pl, uk, ro**”. The remaining 76 languages in  
17 Tatoeba are treated as long-tail ones.

**18 A.2 Hyper-parameters**

19 The parameters of EMMA-X are first initialized with XLM-R, with 24 layers of Transformer [Vaswani  
20 et al., 2017] encoder, 1024 hidden states, and 16 attention heads. We set the total semantic ranks as 4.  
21 The GMM classifier is implemented as a mixture of Gaussian forms, each of which consists of a prior  
22  $\pi \in \mathbb{R}^1$ , a mean  $\mu \in \mathbb{R}^{1024}$ , and a variance  $\sigma \in \mathbb{R}^{1024}$ , all are trainable variables. We optimize the  
23 GMM classifier with Adam ( $\beta_1=0.9$ ,  $\beta_2=0.999$ ) Kingma and Ba [2015] using a batch size of 1024  
24 and a learning rate of 3e-5. For cross-lingual encoder, we apply the same training setting as MoCo He  
25 et al. [2020], with the momentum queue  $K$  to be 256 and temperature as 0.04. We set the momentum  
26 coefficient to 0.999 and use the Adam optimizer with a cosine decay learning rate whose peak is 5e-4.

**27 B FLORES-200 Dataset and Geometric Analysis****28 B.1 FLORES-200 dataset**

29 FLORES-200 Goyal et al. [2022], Costa-jussà et al. [2022] is a many-to-many multilingual bench-  
30 mark, which consists of 3001 sentences in 204 total languages. FLORES-200 sourced all sentences  
31 from English WikiMedia and translated these English sentences to 204 languages by human trans-  
32 lators. In particular, sentences in FLORES-200 have a much larger breadth of topics, for they are

<sup>1</sup>[https://github.com/facebookresearch/cc\\_net](https://github.com/facebookresearch/cc_net)

<sup>2</sup><https://github.com/Helsinki-NLP/OpusFilter>

Code	Size (GB)	Sent. (M)												
af	1.3	-	et	6.1	22.3	ja	24.2	89.2	mt	0.2	-	sq	3.0	-
am	0.7	-	eu	2.0	0.81	jv	0.2	-	my	0.9	-	sr	5.1	-
ar	20.4	72.3	fa	21.6	7.5	ka	3.4	2.0	ne	2.6	-	su	0.1	-
as	0.1	-	fi	19.2	92.8	kk	2.6	2.8	nl	15.8	66.0	sv	10.8	74.2
az	3.6	0.82	fr	46.5	331.5	km	1.0	0.84	no	3.7	-	sw	1.6	1.7
be	3.5	0.51	fy	0.2	0.13	kn	1.2	-	om	0.1	-	ta	8.2	2.79
bg	22.6	47.2	ga	0.5	-	ko	17.2	79.3	or	0.6	-	te	2.6	-
bn	7.9	7.52	gd	0.1	0.05	ku	0.4	-	pa	0.8	-	th	14.7	13.1
br	0.1	-	gl	2.9	0.77	ky	1.2	-	pl	16.8	79.7	tl	0.8	-
bs	0.1	-	gu	0.3	-	la	2.5	-	ps	0.7	-	tr	17.3	93.8
ca	10.1	14.9	ha	0.3	-	lo	0.6	-	pt	15.9	247.6	ug	0.4	-
cs	16.3	108.4	he	6.7	47.1	lt	7.2	11.0	ro	8.6	60.4	uk	9.1	0.78
cy	0.8	-	hi	20.2	3.2	lv	6.4	0.37	ru	48.1	134.9	ur	5.0	1.15
da	15.2	8.0	hr	5.4	-	mg	0.2	-	sa	0.3	-	uz	0.7	-
de	46.3	283.4	hu	9.5	55.2	mk	1.9	-	sd	0.4	-	vi	44.6	15.3
el	29.3	95.1	hy	5.5	1.7	ml	4.3	1.07	si	2.1	0.60	xh	0.1	-
en	49.7	-	id	10.6	184.6	mn	1.7	0.19	sk	4.9	-	yi	0.3	-
eo	0.9	0.18	is	1.3	-	mr	1.3	-	sl	2.8	9.8	zh	36.8	379.4
es	44.6	279.6	it	19.8	179.3	ms	3.2	2.1	so	0.4	-	-	-	-

Table 1: The statistics of CC-100 and the collected parallel corpora used for training. We report the list of 94 languages and include the size of the monolingual data (in GiB) and the number of sentence pairs (in Millions, which denotes the number of sentence pairs between the specific language and English) in parallel corpora for each language. “-” means the number of sentence pairs is less than 0.1 million.

Number of Sentences	3001
Average Words per Sentence	21
Number of Articles	842
Average Number of Sentences per Article	3.5
<b>Domain</b>	
WikiNews	309
WikiJunior	284
WikiVoyage	249
<b>Sub-Topic</b>	
Crime	155
Disasters	27
Entertainment	28
Geography	36
Health	27
Nature	17
Politics	171
Science	154
Sports	154
Travel	505
<b>Articles</b>	
<b>Sentences</b>	
993	313
1006	65
1002	68
341	86
325	67
162	45
1529	341

Table 2: Basic Statistics of FLORES-200.

33 collected from three different sources: WikiNews<sup>3</sup>, WikiJunior<sup>4</sup> and WikiVoyage<sup>5</sup>. We summa-  
 34 rize the basic statistics of all languages in FLORES-200 in Table 2. Similar to Tatoeba [Artetxe  
 35 and Schwenk, 2019], we treat English data “eng\_Latn” as retrieval labels and report the retrieval  
 36 accuracy using the same scripts as Tatoeba in XTREME [Ruder et al., 2021]. We set the 68 lan-  
 37 guages: “bel\_Cyrl, bos\_Latn, hun\_Latn, epo\_Latn, khm\_Khmr, urd\_Arab, srp\_Cyrl, jav\_Latn,  
 38 hye\_Armn, gla\_Latn, por\_Latn, lit\_Latn, bul\_Cyrl, slk\_Latn, mal\_Mlym, ita\_Latn, nno\_Latn,  
 39 mar\_Deva, hrv\_Latn, hin\_Deva, kat\_Geor, ben\_Beng, fin\_Latn, cym\_Latn, oci\_Latn, cat\_Latn,  
 40 fao\_Latn, xho\_Latn, spa\_Latn, ron\_Latn, amh\_Ethi, ces\_Latn, swe\_Latn, nld\_Latn, tat\_Cyrl,  
 41 kor\_Hang, glg\_Latn, fra\_Latn, eus\_Latn, ind\_Latn, dan\_Latn, tha\_Thai, deu\_Latn, tel\_Telu,  
 42 afr\_Latn, pol\_Latn, est\_Latn, uig\_Arab, ukr\_Cyrl, uzn\_Latn, heb\_Hebr, kaz\_Cyrl, nob\_Latn,  
 43 rus\_Cyrl, vie\_Latn, arb\_Arab, zho\_Hans, tuk\_Latn, khk\_Cyrl, jpn\_Jpan, ell\_Grek, isl\_Latn,  
 44 tam\_Taml, slv\_Latn, tur\_Latn, mkd\_Cyrl, tgl\_Latn, gle\_Latn” as “Head” languages, and the  
 45 remaining 135 languages (excluded English data) as “Long-tail” ones.

<sup>3</sup><https://en.wikinews.org/wiki/MainPage>

<sup>4</sup><https://en.wikibooks.org/wiki/Wikijunior>

<sup>5</sup>[https://en.wikivoyage.org/wiki/Main\\_Page](https://en.wikivoyage.org/wiki/Main_Page)

Task category	Task	Train	Dev	Test	Lang.	Metric	Domain
Inference	AmericasNLI	392,702	222,743	738,750	10	Accuracy	Misc.
	XNLI	392,702	2,490	5,010	15	Accuracy	Misc.
Semantic Similarity	Multi-STS	550,152+5,749	10,000+1,500	250	7	Spearman	Misc.
	WMT21QETask1	7,000	1,000	1,000	7 (11)	Pearson	News
Sentence Retrieval	LAReQA	87,599	10,579	1,190	11	mAP@20	Wikipedia
	Mewsl-X	116,093	10,252	428-1,482	11 (50)	mAP@20	News
	BUCC	-	-	1,896-14,330	5	F1	Wiki/News
	Tatoeba	-	-	1,000	36 (122)	Accuracy	Misc.
Classification	XCOPA	33,410+400	100	500	11	Accuracy	Misc.
	MultiEURLEX	55,000	5,000	5,000	23	Accuracy	Legal
	MultiARC	200,000	5,000	5,000	6	MAE	Reviews
	PAWS-X	49,401	2,000	2,000	7	Accuracy	Wiki/Quora

Table 3: Overview of XRETE tasks. For tasks that have training and dev sets in other language, we only report the number of sentences in English sets. We report the number of test examples per languages.

## 46 B.2 Three measurements in Geometric Analysis

**Invariance Measurement** implies whether the semantic distributions of all languages are similar [Abend and Rappoport, 2017]. We adopt a Gaussian form  $\mathcal{N}_l(\mu_l, \sigma_l^2)$  where  $\mu_l = \frac{\sum_{x \in l} \gamma^{(x)}}{3001}$  and  $\sigma_l^2 = \sum_{x \in l} (\gamma^{(x)} - \mu_l)(\gamma^{(x)} - \mu_l)^T$ , to approximate the semantic distribution of each language  $l$ . Further, we compute the mean averaged KL-divergence (KL-D for short) [Kullback and Leibler, 1951] among all language pairs as the overall Invariance score  $\mathcal{I}_v$  with  $L$  as the total number of languages:

$$\mathcal{I}_v = \frac{1}{L \times (L-1)} \sum_{l_1 \neq l_2} \frac{\text{KL}(\mathcal{N}_{l_1} || \mathcal{N}_{l_2}) + \text{KL}(\mathcal{N}_{l_2} || \mathcal{N}_{l_1})}{2}. \quad (1)$$

**Canonical Form Measurement** Previous works [Teller, 2000, Irwin et al., 2009] have demonstrated that a good multilingual space should distribute sentence representations based on their semantic similarities rather than language families. To measure this in quantity, we focus on Calinski-Harabasz Index (CH-I) [Caliński and Harabasz, 1974], which measures how similar an object is to its own cluster compared to other clusters. We group all semantically equivalent sentences in a cluster, which leads to 3001 clusters and each observes 204 sentences in 204 different languages. Assuming  $c_k$  and  $c$  are the centroid of cluster  $k$  and the whole dataset  $s$ , respectively. The CH-I  $\mathcal{C}_h$  is defined as:

$$\mathcal{C}_h = \left[ 204 \times \sum_{k=1}^K \|c_k - c\|^2 \right] / \left[ \sum_{k=1}^K \sum_{s \in S} \|s - c_k\|^2 \right]. \quad (2)$$

47 The higher the CH-I is, the better the semantically equivalent sentences are clustered.

**Isotropy Measurement** A high-dimensional embedding space often demonstrates poor isotropy, and deteriorates into a low-dimensional manifold that greatly limits the expressive ability of the embedding space. We adopt principal ratio (PR) [Mu and Viswanath, 2018] to measure isotropy. Let  $\mathcal{E}$  be the sentence representation matrix,  $\mathcal{V}$  be the set of the eigenvectors of  $\mathcal{E}$ , the Isotropy  $\mathcal{I}_{so}$  is

$$\mathcal{I}_{so} = \min_{v \in \mathcal{V}} \sum_{e \in \mathcal{E}} \exp(v^\top e) / \max_{v \in \mathcal{V}} \sum_{e \in \mathcal{E}} \exp(v^\top e). \quad (3)$$

48 The closer  $\mathcal{I}_{so}$  is to 1, the more isotropic the representation space is.

## 49 C XRETE: Cross-lingual Representation Transfer Evaluation

50 XRETE consists of 12 tasks that fall into four different categories. In our “translate-train-all” setting,  
51 we individually fine-tune models with English training set and its translated training sets on each  
52 task. Then we report the performance of our fine-tuned model. We give an overview in Table 3 and  
53 describe the task details as follows.

54 **XNLI** The Cross-lingual Natural Language Inference corpus Conneau et al. [2018] tasks the  
55 systems with reading two sentences and determining whether one entails the other, contradicts it,

56 or neither (neutral). A crowdsourcing-based procedure is used for collecting English examples,  
57 which are later translated into ten target languages for evaluation. Training data stays consistent  
58 with the English training data of MultiNLI Williams et al. [2018]. For evaluation, we concatenate  
59 two sentences as input and apply a new classification head to distinguish sentence relationships. We  
60 perform “translate-train-all” evaluation, where model is first fine-tuned on English training data and  
61 its translated data in other languages, then evaluated on test sets.

62 **AmericasNLI (ANLI)** The AmericasNLI Ebrahimi et al. [2022] is an extension of XNLI task to  
63 10 Indigenous languages of the Americas. All of these languages are truly low-resource languages  
64 and serve as a good testbed for zero-shot cross-lingual transferability. As Spanish is more relative to  
65 the target languages, the Spanish version of XNLI subset is translated for evaluation. For training,  
66 both English and Spanish versions of MultiNLI training data are provided. We evaluate on ANLI  
67 following the same settings as in XNLI.

68 **MultiSTS** The Multilingual Semantic Textual Similarity dataset Cer et al. [2017], Reimers and  
69 Gurevych [2020] aims to assign a semantic similarity score for a pair of sentences. The MultiSTS  
70 dataset contains 7 cross-lingual sentence pairs and 3 monolingual pairs. Stanford NLI Bowman et al.  
71 [2015] and English STS Cer et al. [2017] are provided as training sets. We report the results after first  
72 fine-tuning on English training set using a Siamese network structure [Reimers and Gurevych, 2020].  
73 Then we compute the cosine similarity between the sentence pairs and compute Spearman’s rank  
74 correlation between the predicted score and gold score following Reimers and Gurevych [2020].

75 **WMT21QETask1 (QE)** The WMT21 Quality Estimation Task 1 Sentence-level Direct Assessment  
76 Specia et al. [2021] aims at testing the translation quality and this task has been applied to test the  
77 sensitivity of language models to semantic similarity Tiyajamorn et al. [2021]. The training and  
78 evaluation sets are collected from Wikipedia by translating sentences using state-of-the-art translation  
79 models to 6 languages and annotated by professional translators. In WMT21, 4 new language pairs  
80 with no training data are given to test zero-shot cross-lingual transferability. Our evaluation setting  
81 on QE is similar to that on MultiSTS, but we report Pearson’s rank correlation [Kepler et al., 2019].

82 **LAReQA** The Language-Agnostic Retrieval Question Answering Roy et al. [2020] is a QA retrieval  
83 task where models are required to retrieve all relevant answers in different languages over a large  
84 multilingual pool. The dataset is constructed on XQuAD Artetxe et al. [2020] and a question is  
85 linked with answer sentences in different languages. The training set of SQuAD v1.1 Rajpurkar  
86 et al. [2016] is used to fine-tune the model to adapt to QA retrieval task. During evaluation, sentence  
87 embeddings are also obtained by a siamese network, and we retrieve the sentences with the highest  
88 cosine similarity as predictions.

89 **Mewsli-X** Mewsli ([Multilingual Entities in News](#), [linked](#)) requires linking an entity mention  
90 to its entry in a language-agnostic knowledge base Botha et al. [2020]. Mewsli-X Ruder et al.  
91 [2021] features 15k mentions in 11 languages. For each mention, Mewsli-X offer entity description  
92 candidate pool containing 1M candidates across 50 languages. Fine-tuning is done on a predefined  
93 set of English-only mention-entity pairs from English Wikipedia hyperlinks. Our evaluation setting is  
94 identical to LAReQA.

95 **BUCC** The second and third shared task of the workshop on Building and Using Parallel Corpora  
96 Zweigenbaum et al. [2017], Pierre Zweigenbaum and Rapp [2018] aims to examine the ability of  
97 models to detect parallel sentence pairs in a pair of monolingual corpora. The dataset provides train  
98 and test splits in 5 languages. Following XTREME Hu et al. [2020], we directly evaluate on BUCC  
99 without fine-tuning and retrieve sentences with the highest cosine similarity.

100 **Tatoeba** The goal of the Tatoeba dataset Artetxe and Schwenk [2019] is to find the nearest neighbor  
101 for each sentence in the other language according to cosine similarity and compute the error rate. The  
102 dataset consists of up to 1,000 English-aligned sentence pairs covering 122 languages. Following  
103 XTREME Hu et al. [2020], we directly evaluate on Tatoeba without fine-tuning and retrieve sentences  
104 with the highest cosine similarity.

105 **XCOPA** In the Cross-lingual Choice of Plausible Alternatives dataset Ponti et al. [2020], each  
106 XCOPA instance corresponds to a premise and two alternatives. The task formulates as a binary

Model	en	ar	bg	de	el	es	fr	hi	ru	sw	th	tr	ur	vi	zh	Avg.
XLM-R	88.6	84.5	86.7	84.6	85.2	84.7	82.0	82.5	82.6	82.4	80.6	83.1	80.3	77.3	77.2	<b>82.8</b>
INFOXLM	90.4	83.9	85.8	86.0	85.6	87.8	86.9	83.9	83.5	83.3	81.2	84.6	82.7	81.6	75.7	<b>84.2</b>
HICTL	90.6	86.8	88.2	87.4	87.0	87.4	85.0	83.9	83.3	84.8	83.1	85.7	82.8	79.7	80.9	<b>85.1</b>
ChatGPT	70.4	61.0	64.5	64.8	62.8	65.7	66.3	51.5	63.4	55.7	53.0	61.6	47.9	61.6	62.6	<b>60.9</b>
<b>EMMA-X</b>	<b>91.9</b>	<b>89.2</b>	<b>90.1</b>	<b>89.6</b>	<b>89.5</b>	<b>90.3</b>	<b>88.7</b>	<b>86.7</b>	<b>85.4</b>	<b>88.5</b>	<b>86.7</b>	<b>89.6</b>	<b>87.7</b>	<b>83.6</b>	<b>83.9</b>	<b>88.1</b>

Table 4: XNLI results (accuracy) for each language.

107 classification to predict the more plausible choice. The English COPA Gordon et al. [2012] training  
 108 set and Social IQa Sap et al. [2019] training data are used for fine-tuning, while the validation and  
 109 test sets of English COPA are translated and re-annotated into 11 languages for evaluation.

110 **MultiEURLEX** The MultiEURLEX dataset Chalkidis et al. [2021] is a legal topic classification  
 111 task which comprises 65k European Union (EU) laws in 23 official EU languages. The dataset  
 112 provides multi-granular labels per document. The dataset is split into training, development, and  
 113 test subsets chronologically, resulting in 55k training documents for 7 languages, and 5k each for  
 114 development and test subsets in all 23 languages.

115 **MultiARC** The Multilingual Amazon Reviews Corpus Keung et al. [2020] is a large-scale collection  
 116 of Amazon reviews for multilingual text classification in 6 languages. Different languages are directly  
 117 gathered from the marketplaces in different countries. The goal is to predict the reviewer’s rating on  
 118 the 5-star scale using the text of the review as input. The data is clearly split into training (200,000  
 119 reviews), development (5,000 reviews), and test sets (5,000 reviews) for each language.

120 **PAWS-X** The Cross-lingual Paraphrase Adversaries from Word Scrambling Yang et al. [2019b]  
 121 dataset requires to identify whether two sentences are paraphrases. A subset of the evaluation pairs in  
 122 English PAWS Zhang et al. [2019] are human-translated into 6 typologically distinct languages for  
 123 evaluation, while the English PAWS training set is used for training.

## 124 D Baseline Methods

125 To fairly evaluate the performance of EMMA-X, we choose XLM-R Conneau and Lample [2019]  
 126 and its several derivatives as our baselines, which contain: (1) XLM-R, which applies multilingual  
 127 MLM tasks as pre-training objectives on CCNet-100 corpus; (2) HICTL Wei et al. [2021], which  
 128 continues training on XLM-R using hierarchical contrastive learning; and (3) INFOXLM, which is  
 129 initialized with XLM-R and trains with cross-lingual contrast, multilingual MLM and TLM. Also,  
 130 we compare EMMA-X to strong sentence models: (1) S-BERT [Reimers and Gurevych, 2020],  
 131 which adopts multilingual knowledge distillation to extend monolingual sentence representations  
 132 to multilingual. We use the strongest baseline, **XLM-R ← SBERT-paraphrase**, proposed in the  
 133 original paper as a baseline. (2) LaBSE [Feng et al., 2022], which systematically combines several  
 134 best methods, including: masked language modeling, translation language modeling [Conneau and  
 135 Lample, 2019], dual encoder translation ranking [Guo et al., 2018], and additive margin softmax [Yang  
 136 et al., 2019a], to learn cross-lingual sentence representations. It filters 17B monolingual sentences  
 137 and 6B translation pairs for sentence representation learning. We take the best model, LaBSE with  
 138 Customized Vocab as our baseline. We further report the zero-shot results on Large Language Model  
 139 (LLM), ChatGPT, which is trained on a wide variety of multilingual sentences and instruction tuning  
 140 based on Reinforcement Learning with Human Feedback [Christiano et al., 2017, Ouyang et al.,  
 141 2022].

## 142 E Prompts for ChatGPT

143 In this section, we show the input prompts of ChatGPT on each task in Table 7.

Model	aym	bzd	cni	gn	hch	nah	oto	quy	shp	tar	Avg.
XLM-R	49.01	50.61	41.72	58.34	42.46	54.63	35.57	59.29	51.62	41.54	<b>48.48</b>
INFOXLM	49.87	51.29	42.41	58.83	43.07	55.25	36.14	59.87	52.20	42.12	<b>49.10</b>
HICTL	49.65	51.22	42.36	58.82	43.09	55.13	36.04	59.61	52.17	42.08	<b>49.02</b>
ChatGPT	42.0	43.6	40.8	40.4	40.0	43.8	41.1	43.1	42.0	40.0	<b>41.7</b>
<b>EMMA-X</b>	<b>51.19</b>	<b>52.50</b>	<b>43.62</b>	<b>59.88</b>	<b>44.31</b>	<b>55.44</b>	<b>39.16</b>	<b>60.14</b>	<b>52.84</b>	<b>43.10</b>	<b>50.21</b>

Table 5: AmericasNLI (ANLI) results (top-1 accuracy) across different input languages.

Model	en-ar	en-de	en-tr	en-es	en-fr	en-it	en-nl	ar-ar	en-en	es-es	Avg.
XLM-R	50.2	63.7	45.8	59.6	68.0	63.4	69.6	87.7	82.5	68.5	<b>65.9</b>
INFOXLM	81.7	80.3	79.9	79.1	80.6	83.4	81.2	86.7	87.2	81.7	<b>82.2</b>
HICTL	80.4	81.8	78.3	80.6	81.2	80.9	79.3	88.4	86.1	79.6	<b>81.6</b>
<b>EMMA-X</b>	<b>86.6</b>	<b>85.0</b>	<b>87.1</b>	<b>84.4</b>	<b>85.2</b>	<b>89.4</b>	<b>88.3</b>	<b>90.9</b>	<b>92.0</b>	<b>84.5</b>	<b>87.3</b>

Table 6: MultiSTS results (Spearman) across different input languages.

## 144 F Results of each Language

145 We show the details for tasks and all languages in Tables 4 (XNLI), 5 (AmericasNLI), 6 (MultiSTS),  
146 8 (QE), 9 (LAReQA), 10 (Mewsli-X), 11 (XCOPA), 12 (BUCC) and 13 (PAWS-X).

## 147 G Equations and Theoretical Analysis

### 148 G.1 Details of Equations

**Details of Gaussian Form  $\mathcal{N}_r$**  In EMMA-X, GMM classifier is introduced to determine the semantic rank of sentence pairs. The posterior probability  $P_{\mathcal{G}}(\cdot)$  of GMM classifier is already discussed in Eq. 5. We show the explicit calculation of Gaussian form  $\mathcal{N}_r(\gamma^{(x_i)}, \gamma^{(y_k)})$  as:

$$\mathcal{N}_r(\gamma^{(x_i)} - \gamma^{(y_k)} | \mu_r, \sigma_r) = \frac{\pi_r}{(2\pi)^{(d/2)} |\text{diag}(\sigma_r)|} \cdot e^{\left(-\frac{1}{2} [(\gamma^{(x_i)} - \gamma^{(y_k)}) - \mu_r]^T \text{diag}(\sigma_r^{-2}) [(\gamma^{(x_i)} - \gamma^{(y_k)}) - \mu_r]\right)}, \quad (4)$$

149 where  $d$  is the dimension of hidden states of  $\gamma^{(x_i)}$  and  $\gamma^{(y_k)}$ .

**Details of contrastive learning** The training objective of cross-lingual encoder in EMMA-X is the ranking InfoNCE loss. We show the explicit expansion of this loss (Eq. 7) as:

$$\begin{aligned}
\mathcal{L}_{\text{CTL}}(\mathcal{X}, \mathcal{Y}; \Theta_{\mathcal{M}}) = & -\mathbb{E}_{\mathbf{x}_i \sim \mathcal{X}} \left[ \right. \\
& \log \underbrace{\frac{\sum_{\mathbf{y}_k \sim \mathcal{Y}_{c_{\mathcal{G}}^*=1}} e^{\frac{s[\gamma(\mathbf{x}_i), \gamma(\mathbf{y}_k)]}{\tau_1}}}{\sum_{\mathbf{y}_t \sim \mathcal{Y}_{c_{\mathcal{G}}^*=1}} e^{\frac{s[\gamma(\mathbf{x}_i), \gamma(\mathbf{y}_t)]}{\tau_1}} + \sum_{\mathbf{y}_t \sim \mathcal{Y}_{c_{\mathcal{G}}^*=2}} e^{\frac{s[\gamma(\mathbf{x}_i), \gamma(\mathbf{y}_t)]}{\tau_1}} + \dots + \sum_{\mathbf{y}_t \sim \mathcal{Y}, \mathcal{Y}_{c_{\mathcal{G}}^*=4}} e^{\frac{s[\gamma(\mathbf{x}_i), \gamma(\mathbf{y}_t)]}{\tau_1}}} \}_{\ell_1} \\
& + \log \underbrace{\frac{\sum_{\mathbf{y}_k \sim \mathcal{Y}_{c_{\mathcal{G}}^*=2}} e^{\frac{s[\gamma(\mathbf{x}_i), \gamma(\mathbf{y}_k)]}{\tau_2}}}{\sum_{\mathbf{y}_t \sim \mathcal{Y}_{c_{\mathcal{G}}^*=2}} e^{\frac{s[\gamma(\mathbf{x}_i), \gamma(\mathbf{y}_t)]}{\tau_2}} + \sum_{\mathbf{y}_t \sim \mathcal{Y}_{c_{\mathcal{G}}^*=3}} e^{\frac{s[\gamma(\mathbf{x}_i), \gamma(\mathbf{y}_t)]}{\tau_2}} + \sum_{\mathbf{y}_t \sim \mathcal{Y}_{c_{\mathcal{G}}^*=4}} e^{\frac{s[\gamma(\mathbf{x}_i), \gamma(\mathbf{y}_t)]}{\tau_2}}} \}_{\ell_2} \\
& + \log \left. \underbrace{\frac{\sum_{\mathbf{y}_k \sim \mathcal{Y}_{c_{\mathcal{G}}^*=3}} e^{\frac{s[\gamma(\mathbf{x}_i), \gamma(\mathbf{y}_k)]}{\tau_3}}}{\sum_{\mathbf{y}_t \sim \mathcal{Y}_{c_{\mathcal{G}}^*=3}} e^{\frac{s[\gamma(\mathbf{x}_i), \gamma(\mathbf{y}_t)]}{\tau_3}} + \sum_{\mathbf{y}_t \sim \mathcal{Y}_{c_{\mathcal{G}}^*=4}} e^{\frac{s[\gamma(\mathbf{x}_i), \gamma(\mathbf{y}_t)]}{\tau_3}}} \right]_{\ell_3}, \tag{5}
\end{aligned}$$

where  $\tau_r$  represents the temperature term. As small temperature  $\tau$  tends to be less tolerant to similar samples, and large  $\tau$  tends to cluster similar samples together [Wang and Liu, 2021], we empirically set  $\tau_1 < \tau_2 < \tau_3 < \tau_4$ , which remains the same as Hoffmann et al. [2022].

## 153 G.2 Theoretical Analysis

154 In this section, we provide detailed proof for Eq. 14 and Eq. 15. Next, we prove the feasibility of  
155 our dual supervision. GMM classifier clusters sentence pairs in terms of Euclidean distance, while  
156 cross-lingual encoder minimizes the covariance of each semantic relation rank via cosine distance.  
157 Finally, we prove that these two metrics are actually equivalent to each other in the unit hypersphere  
158 of the embedding space.

**Proof of Eq. 14.** We provide the derivation of Eq. 14. With the assumption that  $P(\mathbf{x}_i, \mathbf{y}_k | c_{\mathcal{G}}^* = r, \Theta) \sim \mathcal{N}_r(\mathbf{x}_i - \mathbf{y}_k | \tilde{\mu}_r, \tilde{\sigma}_r)$ , we have,

$$\begin{aligned}
\sum_{\mathbf{x}_i \in \mathcal{X}} \sum_{\mathbf{y}_k \in \mathcal{Y}} \sum_{r=1}^N Q(r) \log \frac{P(\mathbf{x}_i, \mathbf{y}_k, r | \Theta)}{Q(r)} & \approx \sum_{\mathbf{x}_i \in \mathcal{X}} \sum_{\mathbf{y}_k \in \mathcal{Y}} \sum_{r=1}^N \log P(\mathbf{x}_i, \mathbf{y}_k | c_{\mathcal{G}}^* = r, \Theta) \\
& = \sum_{\mathbf{x}_i \in \mathcal{X}} \sum_{\mathbf{y}_k \in \mathcal{Y}} \sum_{r=1}^N \left( \log \left( \frac{1}{(2\pi)^{(d/2)} |\tilde{\sigma}_r|^{1/2}} \right) \right. \\
& \quad \left. + \frac{1}{2} [(\mathbf{x}_i - \mathbf{y}_k) - \tilde{\mu}_r]^T \tilde{\sigma}_r^{-1} [(\mathbf{x}_i - \mathbf{y}_k) - \tilde{\mu}_r] \right) \\
& \geq \sum_{r=1}^N \left[ \sum_{\mathbf{x}_i \in \mathcal{X}} \sum_{\mathbf{y}_k \in \mathcal{Y}} (\mathbf{x}_i - \mathbf{y}_k) \right]^2 - 2\tilde{\mu}_r \sum_{\mathbf{x}_i \in \mathcal{X}} \sum_{\mathbf{y}_k \in \mathcal{Y}} (\mathbf{x}_i - \mathbf{y}_k) + n\tilde{\mu}_r^2 \\
& = \sum_{r=1}^N n^2 \tilde{\mu}_r^2 - n\tilde{\mu}_r^2 \\
& = n(n-1) \sum_{r=1}^N \tilde{\mu}_r^2, \tag{6}
\end{aligned}$$

159 with  $n$  denoting the number of sentence pairs in semantic rank  $r$ . Here, we ignore the impact of  $\tilde{\sigma}_r$ .

**Proof of Eq. 15.** As we apply dual supervision, data in the contrastive label space also follows the distribution  $\mathcal{N}_r(\mathbf{x}_i - \mathbf{y}_k | \tilde{\mu}_r, \tilde{\sigma}_r)$ . Hence, under mild assumptions, we can get:

$$\begin{aligned}
\mathcal{L}_{\text{CTL}}^+(\mathcal{X}, \mathcal{Y}; \Theta_{\mathcal{M}}) &= \mathbb{E}_{\mathbf{x}_i \sim \mathcal{X}} \sum_{r=1}^{N-1} \log \sum_{\mathbf{y}_k \sim \mathcal{Y}_{c_{\mathcal{G}}^*=r}} e^{s[\gamma^{(\mathbf{x}_i)}, \gamma^{(\mathbf{y}_k)}]} \\
&= \sum_{\mathbf{x}_i \in \mathcal{X}} \sum_{\mathbf{y}_k \in \mathcal{Y}} \sum_{r=1}^{N-1} s(\mathbf{x}_i, \mathbf{y}_k) \\
&= \sum_{\mathbf{x}_i \in \mathcal{X}} \sum_{\mathbf{y}_k \in \mathcal{Y}} \sum_{r=1}^{N-1} \frac{(\mathbf{x}_i - \mathbf{y}_k)^2 - 2}{2} \\
&= n^2 \sum_{r=1}^{N-1} \tilde{\mu}_r^2.
\end{aligned} \tag{7}$$

Based on the definition of semantic ranks, we have  $\tilde{\mu}_1 < \tilde{\mu}_2 < \dots < \tilde{\mu}_N$ . Empirically, the number of sentence pairs in each rank  $n$  is larger than the number of semantic ranks  $N$ . Hence, it can be derived that:

$$\begin{aligned}
\mathcal{L}_{\text{CTL}}^+(\mathcal{X}, \mathcal{Y}; \Theta_{\mathcal{M}}) &= n^2 \sum_{r=1}^{N-1} \tilde{\mu}_r^2 \\
&< n^2 \sum_{r=1}^{N-1} \tilde{\mu}_r^2 + n^2 \tilde{\mu}_N^2 - n \sum_{r=1}^N \tilde{\mu}_r^2 \\
&= n(n-1) \sum_{r=1}^N \tilde{\mu}_r^2 \\
&\leq \sum_{\mathbf{x}_i \in \mathcal{X}} \sum_{\mathbf{y}_k \in \mathcal{Y}} \sum_{r=1}^N Q(r) \log \frac{P(\mathbf{x}_i, \mathbf{y}_k, r | \Theta)}{Q(r)}.
\end{aligned} \tag{8}$$

160 Therefore, we prove that minimizing the positive terms  $\mathcal{L}_{\text{CTL}}^+(\mathcal{X}, \mathcal{Y}; \Theta_{\mathcal{M}})$  in contrastive learning is  
161 equivalent to maximizing a lower bound of the likelihood in Eq. 12.

**Feasibility of Dual Supervision** According to the definition of semantic ranks, the approximated semantic rank  $c_{\mathcal{G}}^*$  from GMM classifier should satisfy the following restriction,

$$\mathbb{E}_{\mathbf{y}_k \sim \mathcal{Y}_{c_{\mathcal{G}}^*=1}} \|\gamma^{(\mathbf{x}_i)} - \gamma^{(\mathbf{y}_k)}\| < \mathbb{E}_{\mathbf{y}_k \sim \mathcal{Y}_{c_{\mathcal{G}}^*=2}} \|\gamma^{(\mathbf{x}_i)} - \gamma^{(\mathbf{y}_k)}\| < \dots < \mathbb{E}_{\mathbf{y}_k \sim \mathcal{Y}_{c_{\mathcal{G}}^*=N}} \|\gamma^{(\mathbf{x}_i)} - \gamma^{(\mathbf{y}_k)}\|. \tag{9}$$

162 Similary, the approximated semantic rank  $c_{\mathcal{M}}^*$  from cross-lingual encoder should satisfy the following  
163 restriction,

$$\mathbb{E}_{\mathbf{y}_k \sim \mathcal{Y}_{c_{\mathcal{M}}^*=1}} s[\gamma^{(\mathbf{x}_i)}, \gamma^{(\mathbf{y}_k)}] > \mathbb{E}_{\mathbf{y}_k \sim \mathcal{Y}_{c_{\mathcal{M}}^*=2}} s[\gamma^{(\mathbf{x}_i)}, \gamma^{(\mathbf{y}_k)}] > \dots > \mathbb{E}_{\mathbf{y}_k \sim \mathcal{Y}_{c_{\mathcal{M}}^*=N}} s[\gamma^{(\mathbf{x}_i)}, \gamma^{(\mathbf{y}_k)}]. \tag{10}$$

Next, we prove that these two restrictions are interchangeable with each other in a unit hypersphere. For simplicity, we consider only two ranks, but extending the explanation to more ranks is trivial. As the Euclidean distance is always larger than 0, we have:

$$\begin{aligned}
\mathbb{E}_{\mathbf{y}_k \sim \mathcal{Y}_{c_{\mathcal{G}}^*=1}} \|\gamma^{(\mathbf{x}_i)} - \gamma^{(\mathbf{y}_k)}\| &< \mathbb{E}_{\mathbf{y}_k \sim \mathcal{Y}_{c_{\mathcal{G}}^*=2}} \|\gamma^{(\mathbf{x}_i)} - \gamma^{(\mathbf{y}_k)}\| \\
\Leftrightarrow \mathbb{E}_{\mathbf{y}_k \sim \mathcal{Y}_{c_{\mathcal{G}}^*=1}} (\gamma^{(\mathbf{x}_i)} - \gamma^{(\mathbf{y}_k)})^2 &< \mathbb{E}_{\mathbf{y}_k \sim \mathcal{Y}_{c_{\mathcal{G}}^*=2}} (\gamma^{(\mathbf{x}_i)} - \gamma^{(\mathbf{y}_k)})^2 \\
\Leftrightarrow \mathbb{E}_{\mathbf{y}_k \sim \mathcal{Y}_{c_{\mathcal{G}}^*=1}} (2 - 2\gamma^{(\mathbf{x}_i)}\gamma^{(\mathbf{y}_k)}) &< \mathbb{E}_{\mathbf{y}_k \sim \mathcal{Y}_{c_{\mathcal{G}}^*=2}} (2 - 2\gamma^{(\mathbf{x}_i)}\gamma^{(\mathbf{y}_k)}) \\
\Leftrightarrow \mathbb{E}_{\mathbf{y}_k \sim \mathcal{Y}_{c_{\mathcal{G}}^*=1}} s[\gamma^{(\mathbf{x}_i)}, \gamma^{(\mathbf{y}_k)}] &> \mathbb{E}_{\mathbf{y}_k \sim \mathcal{Y}_{c_{\mathcal{G}}^*=2}} s[\gamma^{(\mathbf{x}_i)}, \gamma^{(\mathbf{y}_k)}] \\
\Leftrightarrow \mathbb{E}_{\mathbf{y}_k \sim \mathcal{Y}_{c_{\mathcal{M}}^*=1}} s[\gamma^{(\mathbf{x}_i)}, \gamma^{(\mathbf{y}_k)}] &> \mathbb{E}_{\mathbf{y}_k \sim \mathcal{Y}_{c_{\mathcal{M}}^*=2}} s[\gamma^{(\mathbf{x}_i)}, \gamma^{(\mathbf{y}_k)}].
\end{aligned} \tag{11}$$

164 From the above analyses, we can tell that the approximated semantic rank from one module can  
165 provide a reasonable supervision signal to guide the training of the other module. Hence, all sentence  
166 pairs will be uniformly distributed according to a unified ranking semantic similarity in the embedding  
167 space.

168 **References**

- 169 Omri Abend and Ari Rappoport. The state of the art in semantic representation. In *Proceedings  
170 of the 55th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long  
171 Papers)*, pages 77–89, Vancouver, Canada, July 2017. Association for Computational Linguistics.  
172 doi: 10.18653/v1/P17-1008. URL <https://aclanthology.org/P17-1008>.
- 173 Farhad Akhbardeh, Arkady Arkhangorodsky, Magdalena Biesialska, Ondřej Bojar, Rajen Chatterjee,  
174 Vishrav Chaudhary, Marta R. Costa-jussa, Cristina España-Bonet, Angela Fan, Christian Feder-  
175 mann, Markus Freitag, Yvette Graham, Roman Grundkiewicz, Barry Haddow, Leonie Harter,  
176 Kenneth Heafield, Christopher Homan, Matthias Huck, Kwabena Amponsah-Kaakyire, Jungo Ka-  
177 saki, Daniel Khashabi, Kevin Knight, Tom Kocmi, Philipp Koehn, Nicholas Lourie, Christof Monz,  
178 Makoto Morishita, Masaaki Nagata, Ajay Nagesh, Toshiaki Nakazawa, Matteo Negri, Santanu  
179 Pal, Allahsera Auguste Tapo, Marco Turchi, Valentin Vydrin, and Marcos Zampieri. Findings of  
180 the 2021 conference on machine translation (WMT21). In *Proceedings of the Sixth Conference  
181 on Machine Translation*, pages 1–88, Online, November 2021. Association for Computational  
182 Linguistics. URL <https://aclanthology.org/2021.wmt-1.1>.
- 183 Mikel Artetxe and Holger Schwenk. Massively multilingual sentence embeddings for zero-shot  
184 cross-lingual transfer and beyond. *Transactions of the Association for Computational Linguistics*,  
185 7:597–610, 2019. doi: 10.1162/tacl\_a\_00288. URL <https://aclanthology.org/Q19-1038>.
- 186 Mikel Artetxe, Sebastian Ruder, and Dani Yogatama. On the cross-lingual transferability of monolin-  
187 gual representations. In *Proceedings of the 58th Annual Meeting of the Association for Compu-  
188 tational Linguistics*, pages 4623–4637, Online, July 2020. Association for Computational Linguistics.  
189 doi: 10.18653/v1/2020.acl-main.421. URL [https://aclanthology.org/2020.acl-main.  
190 421](https://aclanthology.org/2020.acl-main.421).
- 191 Jan A. Botha, Zifei Shan, and Daniel Gillick. Entity Linking in 100 Languages. In *Proceedings of  
192 the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP)*, pages  
193 7833–7845, Online, November 2020. Association for Computational Linguistics. doi: 10.18653/  
194 v1/2020.emnlp-main.630. URL <https://aclanthology.org/2020.emnlp-main.630>.
- 195 Samuel R. Bowman, Gabor Angeli, Christopher Potts, and Christopher D. Manning. A large  
196 annotated corpus for learning natural language inference. In *Proceedings of the 2015 Conference on  
197 Empirical Methods in Natural Language Processing*, pages 632–642, Lisbon, Portugal, September  
198 2015. Association for Computational Linguistics. doi: 10.18653/v1/D15-1075. URL <https://www.aclweb.org/anthology/D15-1075>.
- 200 T. Caliński and J Harabasz. A dendrite method for cluster analysis. *Communications in Statistics*,  
201 3(1):1–27, 1974. doi: 10.1080/03610927408827101. URL <https://www.tandfonline.com/doi/abs/10.1080/03610927408827101>.
- 203 Daniel Cer, Mona Diab, Eneko Agirre, Iñigo Lopez-Gazpio, and Lucia Specia. SemEval-2017 task 1:  
204 Semantic textual similarity multilingual and crosslingual focused evaluation. In *Proceedings of the  
205 11th International Workshop on Semantic Evaluation (SemEval-2017)*, pages 1–14, Vancouver,  
206 Canada, August 2017. Association for Computational Linguistics. doi: 10.18653/v1/S17-2001.  
207 URL <https://aclanthology.org/S17-2001>.
- 208 Ilias Chalkidis, Manos Fergadiotis, and Ion Androutsopoulos. MultiEURLEX - a multi-lingual and  
209 multi-label legal document classification dataset for zero-shot cross-lingual transfer. In *Proceedings  
210 of the 2021 Conference on Empirical Methods in Natural Language Processing*, pages 6974–6996,  
211 Online and Punta Cana, Dominican Republic, November 2021. Association for Computational  
212 Linguistics. doi: 10.18653/v1/2021.emnlp-main.559. URL <https://aclanthology.org/2021.emnlp-main.559>.
- 214 Zewen Chi, Li Dong, Furu Wei, Nan Yang, Saksham Singhal, Wenhui Wang, Xia Song, Xian-Ling  
215 Mao, Heyan Huang, and Ming Zhou. InfoXLM: An information-theoretic framework for cross-  
216 lingual language model pre-training. In *Proceedings of the 2021 Conference of the North American  
217 Chapter of the Association for Computational Linguistics: Human Language Technologies*, pages  
218 3576–3588, Online, June 2021. Association for Computational Linguistics. doi: 10.18653/v1/  
219 2021.naacl-main.280. URL <https://aclanthology.org/2021.naacl-main.280>.

- 220 Paul F Christiano, Jan Leike, Tom Brown, Miljan Martic, Shane Legg, and Dario Amodei. Deep reinforcement learning from human preferences. In I. Guyon, U. Von Luxburg, S. Bengio, H. Wallach, R. Fergus, S. Vishwanathan, and R. Garnett, editors, *Advances in Neural Information Processing Systems*, volume 30. Curran Associates, Inc., 2017. URL [https://proceedings.neurips.cc/paper\\_files/paper/2017/file/d5e2c0adad503c91f91df240d0cd4e49-Paper.pdf](https://proceedings.neurips.cc/paper_files/paper/2017/file/d5e2c0adad503c91f91df240d0cd4e49-Paper.pdf).
- 225 Alexis Conneau and Guillaume Lample. Cross-lingual language model pretraining. In H. Wallach, H. Larochelle, A. Beygelzimer, F. d'Alché-Buc, E. Fox, and R. Garnett, editors, *Advances in Neural Information Processing Systems*, volume 32. Curran Associates, Inc., 2019. URL <https://proceedings.neurips.cc/paper/2019/file/c04c19c2c2474dbf5f7ac4372c5b9af1-Paper.pdf>.
- 230 Alexis Conneau, Ruty Rinott, Guillaume Lample, Adina Williams, Samuel Bowman, Holger Schwenk, and Veselin Stoyanov. XNLI: Evaluating cross-lingual sentence representations. In *Proceedings of the 2018 Conference on Empirical Methods in Natural Language Processing*, pages 2475–2485, Brussels, Belgium, October–November 2018. Association for Computational Linguistics. doi: 10.18653/v1/D18-1269. URL <https://www.aclweb.org/anthology/D18-1269>.
- 235 Alexis Conneau, Kartikay Khandelwal, Naman Goyal, Vishrav Chaudhary, Guillaume Wenzek, Francisco Guzmán, Edouard Grave, Myle Ott, Luke Zettlemoyer, and Veselin Stoyanov. Unsupervised cross-lingual representation learning at scale. In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, pages 8440–8451, Online, July 2020. Association for Computational Linguistics. doi: 10.18653/v1/2020.acl-main.747. URL <https://aclanthology.org/2020.acl-main.747>.
- 241 Marta R Costa-jussà, James Cross, Onur Çelebi, Maha Elbayad, Kenneth Heafield, Kevin Heffernan, Elahe Kalbassi, Janice Lam, Daniel Licht, Jean Maillard, et al. No language left behind: Scaling human-centered machine translation. *arXiv preprint arXiv:2207.04672*, 2022.
- 244 Abteen Ebrahimi, Manuel Mager, Arturo Oncevay, Vishrav Chaudhary, Luis Chiruzzo, Angela Fan, John Ortega, Ricardo Ramos, Annette Rios, Ivan Vladimir Meza Ruiz, Gustavo Giménez-Lugo, Elisabeth Mager, Graham Neubig, Alexis Palmer, Rolando Coto-Solano, Thang Vu, and Katharina Kann. AmericasNLI: Evaluating zero-shot natural language understanding of pretrained multilingual models in truly low-resource languages. In *Proceedings of the 60th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 6279–6299, Dublin, Ireland, May 2022. Association for Computational Linguistics. doi: 10.18653/v1/2022.acl-long.435. URL <https://aclanthology.org/2022.acl-long.435>.
- 252 Ahmed El-Kishky, Vishrav Chaudhary, Francisco Guzmán, and Philipp Koehn. CCAligned: A massive collection of cross-lingual web-document pairs. In *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP)*, pages 5960–5969, Online, November 2020. Association for Computational Linguistics. doi: 10.18653/v1/2020.emnlp-main.480. URL <https://aclanthology.org/2020.emnlp-main.480>.
- 257 Fangxiaoyu Feng, Yinfei Yang, Daniel Cer, Naveen Arivazhagan, and Wei Wang. Language-agnostic BERT sentence embedding. In *Proceedings of the 60th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 878–891, Dublin, Ireland, May 2022. Association for Computational Linguistics. doi: 10.18653/v1/2022.acl-long.62. URL <https://aclanthology.org/2022.acl-long.62>.
- 262 Andrew Gordon, Zornitsa Kozareva, and Melissa Roemmele. SemEval-2012 task 7: Choice of plausible alternatives: An evaluation of commonsense causal reasoning. In *\*SEM 2012: The First Joint Conference on Lexical and Computational Semantics – Volume 1: Proceedings of the main conference and the shared task, and Volume 2: Proceedings of the Sixth International Workshop on Semantic Evaluation (SemEval 2012)*, pages 394–398, Montréal, Canada, 7–8 June 2012. Association for Computational Linguistics. URL <https://aclanthology.org/S12-1052>.
- 268 Naman Goyal, Cynthia Gao, Vishrav Chaudhary, Peng-Jen Chen, Guillaume Wenzek, Da Ju, Sanjana Krishnan, Marc'Aurelio Ranzato, Francisco Guzmán, and Angela Fan. The Flores-101 evaluation benchmark for low-resource and multilingual machine translation. *Transactions of the Association for Computational Linguistics*, 10:522–538, 2022. doi: 10.1162/tacl\_a\_00474. URL <https://aclanthology.org/2022.tacl-1.30>.

- 273 Mandy Guo, Qinlan Shen, Yinfei Yang, Heming Ge, Daniel Cer, Gustavo Hernandez Abrego,  
 274 Keith Stevens, Noah Constant, Yun-Hsuan Sung, Brian Strope, and Ray Kurzweil. Effective  
 275 parallel corpus mining using bilingual sentence embeddings. In *Proceedings of the Third Con-*  
*276 ference on Machine Translation: Research Papers*, pages 165–176, Brussels, Belgium, Oc-  
*277 tober 2018. Association for Computational Linguistics.* doi: 10.18653/v1/W18-6317. URL  
*278* <https://aclanthology.org/W18-6317>.
- 279 Kaiming He, Haoqi Fan, Yuxin Wu, Saining Xie, and Ross Girshick. Momentum contrast for  
 280 unsupervised visual representation learning. In *2020 IEEE/CVF Conference on Computer Vision*  
*281 and Pattern Recognition (CVPR)*, pages 9726–9735. IEEE Computer Society, 2020.
- 282 David T Hoffmann, N Behrmann, J Gall, Thomas Brox, and M Noroozi. Ranking info noise  
 283 contrastive estimation: Boosting contrastive learning via ranked positives. In *AAAI Conference on*  
*284 Artificial Intelligence*, 2022.
- 285 Junjie Hu, Sebastian Ruder, Aditya Siddhant, Graham Neubig, Orhan Firat, and Melvin Johnson.  
 286 XTREME: A massively multilingual multi-task benchmark for evaluating cross-lingual generalisa-  
 287 tion. In Hal Daumé III and Aarti Singh, editors, *Proceedings of the 37th International Conference*  
*288 on Machine Learning*, volume 119 of *Proceedings of Machine Learning Research*, pages 4411–  
 289 4421. PMLR, 13–18 Jul 2020. URL <https://proceedings.mlr.press/v119/hu20b.html>.
- 290 Haoyang Huang, Yaobo Liang, Nan Duan, Ming Gong, Linjun Shou, Dixin Jiang, and Ming Zhou.  
 291 Unicoder: A universal language encoder by pre-training with multiple cross-lingual tasks. In  
 292 *Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing and*  
*293 the 9th International Joint Conference on Natural Language Processing (EMNLP-IJCNLP)*, pages  
 294 2485–2494, Hong Kong, China, November 2019. Association for Computational Linguistics. doi:  
 295 10.18653/v1/D19-1252. URL <https://aclanthology.org/D19-1252>.
- 296 Jeannie Y Irwin, Henk Harkema, Lee M Christensen, Titus Schleyer, Peter J Haug, and Wendy W  
 297 Chapman. Methodology to develop and evaluate a semantic representation for nlp. In *AMIA Annual*  
*298 Symposium Proceedings*, volume 2009, page 271. American Medical Informatics Association,  
 299 2009.
- 300 Fabio Kepler, Jonay Trénous, Marcos Treviso, Miguel Vera, and André F. T. Martins. OpenKiwi:  
 301 An open source framework for quality estimation. In *Proceedings of the 57th Annual Meeting of*  
*302 the Association for Computational Linguistics: System Demonstrations*, pages 117–122, Florence,  
 303 Italy, July 2019. Association for Computational Linguistics. doi: 10.18653/v1/P19-3020. URL  
 304 <https://aclanthology.org/P19-3020>.
- 305 Phillip Keung, Yichao Lu, György Szarvas, and Noah A. Smith. The multilingual Amazon reviews  
 306 corpus. In *Proceedings of the 2020 Conference on Empirical Methods in Natural Language*  
*307 Processing (EMNLP)*, pages 4563–4568, Online, November 2020. Association for Computational  
 308 Linguistics. doi: 10.18653/v1/2020.emnlp-main.369. URL <https://aclanthology.org/2020.emnlp-main.369>.
- 310 Diederik P. Kingma and Jimmy Ba. Adam: A method for stochastic optimization. In *3rd International*  
*311 Conference on Learning Representations, San Diego, CA*, 2015. URL <http://arxiv.org/abs/1412.6980>.
- 313 Taku Kudo and John Richardson. SentencePiece: A simple and language independent subword  
 314 tokenizer and detokenizer for neural text processing. In *Proceedings of the 2018 Conference*  
*315 on Empirical Methods in Natural Language Processing: System Demonstrations*, pages 66–71,  
 316 Brussels, Belgium, November 2018. Association for Computational Linguistics. doi: 10.18653/v1/  
 317 D18-2012. URL <https://aclanthology.org/D18-2012>.
- 318 S. Kullback and R. A. Leibler. On Information and Sufficiency. *The Annals of Mathematical Statistics*,  
 319 22(1):79 – 86, 1951. doi: 10.1214/aoms/1177729694. URL <https://doi.org/10.1214/aoms/1177729694>.
- 321 Fuli Luo, Wei Wang, Jiahao Liu, Yijia Liu, Bin Bi, Songfang Huang, Fei Huang, and Luo Si.  
 322 VECO: Variable and flexible cross-lingual pre-training for language understanding and generation.  
 323 In *Proceedings of the 59th Annual Meeting of the Association for Computational Linguistics*

- 324 and the 11th International Joint Conference on Natural Language Processing (Volume 1: Long  
325 Papers), pages 3980–3994, Online, August 2021. Association for Computational Linguistics. doi:  
326 10.18653/v1/2021.acl-long.308. URL <https://aclanthology.org/2021.acl-long.308>.
- 327 Jiaqi Mu and Pramod Viswanath. All-but-the-top: Simple and effective postprocessing for word  
328 representations. In *International Conference on Learning Representations*, 2018. URL <https://openreview.net/forum?id=HkuGJ3kCb>.
- 330 Long Ouyang, Jeffrey Wu, Xu Jiang, Diogo Almeida, Carroll Wainwright, Pamela Mishkin, Chong  
331 Zhang, Sandhini Agarwal, Katarina Slama, Alex Gray, John Schulman, Jacob Hilton, Fraser Kelton,  
332 Luke Miller, Maddie Simens, Amanda Askell, Peter Welinder, Paul Christiano, Jan Leike, and  
333 Ryan Lowe. Training language models to follow instructions with human feedback. In Alice H. Oh,  
334 Alekh Agarwal, Danielle Belgrave, and Kyunghyun Cho, editors, *Advances in Neural Information  
335 Processing Systems*, 2022. URL <https://openreview.net/forum?id=TG8KACxEON>.
- 336 Serge Sharoff Pierre Zweigenbaum and Reinhard Rapp. Overview of the third bucc shared task:  
337 Spotting parallel sentences in comparable corpora. In Reinhard Rapp, Pierre Zweigenbaum,  
338 and Serge Sharoff, editors, *Proceedings of the Eleventh International Conference on Language  
339 Resources and Evaluation (LREC 2018)*, Paris, France, may 2018. European Language Resources  
340 Association (ELRA). ISBN 979-10-95546-07-8.
- 341 Edoardo Maria Ponti, Goran Glavaš, Olga Majewska, Qianchu Liu, Ivan Vulić, and Anna Korhonen.  
342 XCOPA: A multilingual dataset for causal commonsense reasoning. In *Proceedings of the 2020  
343 Conference on Empirical Methods in Natural Language Processing (EMNLP)*, pages 2362–2376,  
344 Online, November 2020. Association for Computational Linguistics. doi: 10.18653/v1/2020.  
345 emnlp-main.185. URL <https://aclanthology.org/2020.emnlp-main.185>.
- 346 Pranav Rajpurkar, Jian Zhang, Konstantin Lopyrev, and Percy Liang. SQuAD: 100,000+ questions  
347 for machine comprehension of text. In *Proceedings of the 2016 Conference on Empirical Methods  
348 in Natural Language Processing*, pages 2383–2392, Austin, Texas, November 2016. Association  
349 for Computational Linguistics. doi: 10.18653/v1/D16-1264. URL <https://aclanthology.org/D16-1264>.
- 351 Nils Reimers and Iryna Gurevych. Making monolingual sentence embeddings multilingual us-  
352 ing knowledge distillation. In *Proceedings of the 2020 Conference on Empirical Methods in  
353 Natural Language Processing (EMNLP)*, pages 4512–4525, Online, November 2020. Associa-  
354 tion for Computational Linguistics. doi: 10.18653/v1/2020.emnlp-main.365. URL <https://aclanthology.org/2020.emnlp-main.365>.
- 356 Uma Roy, Noah Constant, Rami Al-Rfou, Aditya Barua, Aaron Phillips, and Yinfei Yang. LAReQA:  
357 Language-agnostic answer retrieval from a multilingual pool. In *Proceedings of the 2020 Confer-  
358 ence on Empirical Methods in Natural Language Processing (EMNLP)*, pages 5919–5930, Online,  
359 November 2020. Association for Computational Linguistics. doi: 10.18653/v1/2020.emnlp-main.  
360 477. URL <https://aclanthology.org/2020.emnlp-main.477>.
- 361 Sebastian Ruder, Noah Constant, Jan Botha, Aditya Siddhant, Orhan Firat, Jinlan Fu, Pengfei Liu,  
362 Junjie Hu, Dan Garrette, Graham Neubig, and Melvin Johnson. XTREME-R: Towards more  
363 challenging and nuanced multilingual evaluation. In *Proceedings of the 2021 Conference on  
364 Empirical Methods in Natural Language Processing*, pages 10215–10245, Online and Punta Cana,  
365 Dominican Republic, November 2021. Association for Computational Linguistics. doi: 10.18653/  
366 v1/2021.emnlp-main.802. URL <https://aclanthology.org/2021.emnlp-main.802>.
- 367 Maarten Sap, Hannah Rashkin, Derek Chen, Ronan Le Bras, and Yejin Choi. Social IQa: Com-  
368 monsense reasoning about social interactions. In *Proceedings of the 2019 Conference on Em-  
369 pirical Methods in Natural Language Processing and the 9th International Joint Conference  
370 on Natural Language Processing (EMNLP-IJCNLP)*, pages 4463–4473, Hong Kong, China,  
371 November 2019. Association for Computational Linguistics. doi: 10.18653/v1/D19-1454. URL  
372 <https://aclanthology.org/D19-1454>.
- 373 Holger Schwenk, Guillaume Wenzek, Sergey Edunov, Edouard Grave, Armand Joulin, and Angela  
374 Fan. CCMATRIX: Mining billions of high-quality parallel sentences on the web. In *Proceedings of the  
375 59th Annual Meeting of the Association for Computational Linguistics and the 11th International*

- 376 *Joint Conference on Natural Language Processing (Volume 1: Long Papers)*, pages 6490–6500,  
377 Online, August 2021. Association for Computational Linguistics. doi: 10.18653/v1/2021.acl-long.  
378 507. URL <https://aclanthology.org/2021.acl-long.507>.
- 379 Lucia Specia, Frédéric Blain, Marina Fomicheva, Chrysoula Zerva, Zhenhao Li, Vishrav Chaudhary,  
380 and André F. T. Martins. Findings of the WMT 2021 shared task on quality estimation. In  
381 *Proceedings of the Sixth Conference on Machine Translation*, pages 684–725, Online, November  
382 2021. Association for Computational Linguistics. URL <https://aclanthology.org/2021.wmt-1.71>.
- 384 Virginia Teller. Speech and Language Processing: An Introduction to Natural Language Processing,  
385 Computational Linguistics, and Speech Recognition. *Computational Linguistics*, 26(4):638–641,  
386 12 2000. ISSN 0891-2017. doi: 10.1162/089120100750105975. URL <https://doi.org/10.1162/089120100750105975>.
- 388 Nattapong Tiyajamorn, Tomoyuki Kajiwara, Yuki Arase, and Makoto Onizuka. Language-  
389 agnostic representation from multilingual sentence encoders for cross-lingual similarity esti-  
390 mation. In *Proceedings of the 2021 Conference on Empirical Methods in Natural Language  
391 Processing*, pages 7764–7774, Online and Punta Cana, Dominican Republic, November 2021.  
392 Association for Computational Linguistics. doi: 10.18653/v1/2021.emnlp-main.612. URL  
393 <https://aclanthology.org/2021.emnlp-main.612>.
- 394 Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Lukasz  
395 Kaiser, and Illia Polosukhin. Attention is all you need. In *Advances in Neural Information  
396 Processing Systems 30, NIPS 2017 4-9 December 2017, Long Beach, CA, USA*, pages 5998–6008,  
397 2017. URL <http://papers.nips.cc/paper/7181-attention-is-all-you-need>.
- 398 Feng Wang and Huaping Liu. Understanding the behaviour of contrastive loss. In *Proceedings of the  
399 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 2495–2504,  
400 June 2021.
- 401 Xiangpeng Wei, Rongxiang Weng, Yue Hu, Luxi Xing, Heng Yu, and Weihua Luo. On learning uni-  
402 versal representations across languages. In *International Conference on Learning Representations*,  
403 2021. URL <https://openreview.net/forum?id=Uu1Nw-eeTxJ>.
- 404 Guillaume Wenzek, Marie-Anne Lachaux, Alexis Conneau, Vishrav Chaudhary, Francisco Guzmán,  
405 Armand Joulin, and Edouard Grave. CCNet: Extracting high quality monolingual datasets from  
406 web crawl data. In *Proceedings of the Twelfth Language Resources and Evaluation Conference*,  
407 pages 4003–4012, Marseille, France, May 2020. European Language Resources Association. ISBN  
408 979-10-95546-34-4. URL <https://aclanthology.org/2020.lrec-1.494>.
- 409 Adina Williams, Nikita Nangia, and Samuel Bowman. A broad-coverage challenge corpus  
410 for sentence understanding through inference. In *Proceedings of the 2018 Conference of  
411 the North American Chapter of the Association for Computational Linguistics: Human Lan-  
412 guage Technologies, Volume 1 (Long Papers)*, pages 1112–1122, New Orleans, Louisiana,  
413 June 2018. Association for Computational Linguistics. doi: 10.18653/v1/N18-1101. URL  
414 <https://aclanthology.org/N18-1101>.
- 415 Yinfei Yang, Gustavo Hernandez Abrego, Steve Yuan, Mandy Guo, Qinlan Shen, Daniel Cer, Yun-  
416 hsuan Sung, Brian Strope, and Ray Kurzweil. Improving multilingual sentence embedding using  
417 bi-directional dual encoder with additive margin softmax. In *Proceedings of the Twenty-Eighth  
418 International Joint Conference on Artificial Intelligence, IJCAI-19*, pages 5370–5378. International  
419 Joint Conferences on Artificial Intelligence Organization, 7 2019a. doi: 10.24963/ijcai.2019/746.  
420 URL <https://doi.org/10.24963/ijcai.2019/746>.
- 421 Yinfei Yang, Yuan Zhang, Chris Tar, and Jason Baldridge. PAWS-X: A cross-lingual adversar-  
422 ial dataset for paraphrase identification. In *Proceedings of the 2019 Conference on Empirical  
423 Methods in Natural Language Processing and the 9th International Joint Conference on Natu-  
424 ral Language Processing (EMNLP-IJCNLP)*, pages 3687–3692, Hong Kong, China, Novem-  
425 ber 2019b. Association for Computational Linguistics. doi: 10.18653/v1/D19-1382. URL  
426 <https://aclanthology.org/D19-1382>.

- 427 Yuan Zhang, Jason Baldridge, and Luheng He. PAWS: Paraphrase adversaries from word scrambling.  
428 In *Proceedings of the 2019 Conference of the North American Chapter of the Association for*  
429 *Computational Linguistics: Human Language Technologies, Volume 1 (Long and Short Papers)*,  
430 pages 1298–1308, Minneapolis, Minnesota, June 2019. Association for Computational Linguistics.  
431 doi: 10.18653/v1/N19-1131. URL <https://aclanthology.org/N19-1131>.
- 432 Michał Ziemska, Marcin Junczys-Dowmunt, and Bruno Pouliquen. The United Nations parallel  
433 corpus v1.0. In *Proceedings of the Tenth International Conference on Language Resources and*  
434 *Evaluation (LREC'16)*, pages 3530–3534, Portorož, Slovenia, May 2016. European Language  
435 Resources Association (ELRA). URL <https://aclanthology.org/L16-1561>.
- 436 Pierre Zweigenbaum, Serge Sharoff, and Reinhart Rapp. Overview of the second BUCC shared  
437 task: Spotting parallel sentences in comparable corpora. In *Proceedings of the 10th Workshop*  
438 *on Building and Using Comparable Corpora*, pages 60–67, Vancouver, Canada, August 2017.  
439 Association for Computational Linguistics. doi: 10.18653/v1/W17-2512. URL <https://www.aclweb.org/anthology/W17-2512>.  
440

---

#### Basic Prompt for XNLI/ANLI

---

**Task Description:** Read the following and determine the relationship between Hypothesis and Premise. Choose relation from “contradiction”, “neutral”, or “entailment”.

**Hypothesis:** Yo... no puedo pensar por qué deberías hablarme así, dijo ella, con menos de lo que le había asegurado antes.

**Premise:** Ella era una buena amiga de él, por esto le dolía que le hablara así.

---

#### Basic Prompt for MultiSTS

---

**Task Description:** Read the following sentences and measure the real-valued meaning similarity between these two sentences. You can choose the meaning similarity score, ranging from 0 for no meaning overlap to 5 for meaning equivalence.

**Sentence1:** A person is on a baseball team.

**Sentence2:** Eine Person spielt in einem Team Basketball.

---

#### Basic Prompt for QE

---

**Task Description:** Read the Source sentence and its Translation, and estimate the quality of the Translation. You can rate the translation from 0-1 according to the perceived translation quality.

**Source:** În Franță a început stagnarea demografică de lungă durată, refacerea durând o generație.

**Translation:** In France, long-term demographic stagnation has started, restoring a generation.

---

#### Basic Prompt for XCOPA

---

**Task Description:** Read the Premise and determine which choice is the effect(or cause) of the Premise . Choose from “Choice1” or “Choice2”.

**Premise:** Kuki kurukuna wasiman haykurqanku.

**Choice1:** Kuki kurukunaqa wasimanta chinkarqanku.

**Choice2:** Kuki kuruqa wasip kurkunta mikhurqanku.

---

#### Basic Prompt for MultiEURLEX

---

**Task Description:** Read the following sentences and determine the legal topic of the given sentence. Legal topic should choose from ‘international organisations’, ‘social questions’, ‘production’, ‘technology and research’, ‘environment’, ‘energy’, ‘transport’, ‘law’, ‘finance’, ‘education and communications’, ‘trade’, ‘agriculture’, ‘forestry and fisheries’, ‘economics’, ‘agri-foodstuffs’, ‘EUROPEAN UNION’, ‘science’, ‘politics’, ‘international relations’, ‘industry’, ‘geography’, ‘business and competition’, ‘employment and working conditions’.

**Sentence:** NEUVOSTON ASETUS (EU) N:o 1390/2013, annetti 16 päivänä joulukuuta 2013, Euroopan unionin ja Komorien liiton kesken näiden välisessä kalastuskumppanuussopimuksessa määritetyjen kalastusmahdollisuuksien ja taloudellisen korvauksen vahvistamisesta hyväksytyn pöytäkirjan mukaisten kalastusmahdollisuuksien jakamisesta ...

---

#### Basic Prompt for MultiARC

---

**Task Description:** Read the following review and predict a 5-star scale rating (1 means the poorest experience and 5 represents excellent or outstanding performance) that can best match the review.

**Review:** no me llego el articulo me lo mando por correos normal sin seguimiento y nunca me llego tota un desastre

---

#### Basic Prompt for PAWS-X

---

**Task Description:** Read the following sentences and determine whether two sentences are paraphrases. Return yes or no.

**Sentence1:** La excepción fue entre fines de 2005 y 2009 cuando jugó en Suecia con Carlstad United BK, Serbia con FK Borac Čačak y el FC Terek Grozny de Rusia.

**Sentence2:** La excepción se dio entre fines del 2005 y 2009, cuando jugó con Suecia en el Carlstad United BK, Serbia con el FK Borac Čačak y el FC Terek Grozny de Rusia.

---

Table 7: Prompts of ChatGPT on each task.

Model	en-de	en-zh	et-en	ne-en	ro-en	ru-en	si-en	en-cs	en-ja	km-en	ps-en	Avg.
XLM-R	0.412	0.566	0.797	0.812	0.891	0.774	0.578	0.547	0.335	0.612	0.635	<b>0.632</b>
INFOXLM	0.517	0.534	0.775	0.834	0.890	0.788	0.581	0.564	0.325	0.635	0.616	<b>0.641</b>
HICTL	0.495	0.579	0.792	0.835	<b>0.904</b>	0.787	0.575	0.556	0.342	0.625	0.648	<b>0.649</b>
EMMA-X	<b>0.580</b>	<b>0.589</b>	<b>0.809</b>	<b>0.854</b>	0.897	<b>0.829</b>	<b>0.593</b>	<b>0.577</b>	<b>0.370</b>	<b>0.641</b>	<b>0.651</b>	<b>0.672</b>

Table 8: WMT21-QE-Task1 results (Pearson) across different input languages.

Model	ar	de	el	en	es	hi	ru	th	tr	vi	zh	Avg.
XLM-R	34.1	42.4	39.3	44.8	44.0	37.3	41.7	38.6	40.9	40.4	39.5	<b>40.3</b>
INFOXLM	39.7	52.6	39.2	55.1	53.4	36.8	51.0	28.5	41.1	48.9	47.3	<b>44.9</b>
HICTL	40.3	53.2	41.7	56.3	54.3	39.6	51.7	30.1	42.8	48.9	48.5	<b>46.1</b>
EMMA-X	<b>45.1</b>	<b>58.4</b>	<b>45.4</b>	<b>60.6</b>	<b>59.8</b>	<b>41.4</b>	<b>56.3</b>	<b>34.7</b>	<b>47.1</b>	<b>54.6</b>	<b>53.4</b>	<b>50.6</b>

Table 9: LAReQA results (mean average precision@20, mAP@20) across different input languages.

Model	ar	de	en	es	fa	ja	pl	ro	ta	tr	uk	Avg.
XLM-R	34.6	66.0	62.6	64.8	27.1	47.8	64.8	33.7	17.8	62.3	53.2	<b>48.6</b>
INFOXLM	40.8	71.6	66.3	68.7	48.7	61.0	66.7	39.2	42.0	64.6	58.1	<b>57.1</b>
HICTL	41.7	68.5	64.2	65.6	45.6	51.9	67.6	40.4	32.8	65.5	58.9	<b>54.8</b>
EMMA-X	<b>50.2</b>	<b>78.7</b>	<b>69.1</b>	<b>63.7</b>	47.9	<b>59.6</b>	<b>70.0</b>	<b>50.2</b>	<b>43.5</b>	<b>68.0</b>	<b>60.9</b>	<b>59.6</b>

Table 10: Mewsl-X results (mean average precision@20, mAP@20) across different input languages.

Model	et	ht	id	it	qu	sw	ta	th	tr	vi	zh	Avg.
XLM-R	73.8	67.4	77.8	72.2	52.3	70.9	72.1	<b>74.6</b>	73.4	73.2	75.7	<b>71.2</b>
INFOXLM	75.1	73.4	<b>78.3</b>	80.7	65.6	69.1	72.7	73.9	76.9	77.8	77.5	<b>74.6</b>
HICTL	75.9	73.1	77.8	<b>81.2</b>	65.5	73.8	72.6	73.2	76.1	75.4	78.0	<b>74.8</b>
ChatGPT	80.6	64.1	85.6	89.2	47.4	75.9	56.4	67.3	82.2	81.5	85.8	<b>74.2</b>
EMMA-X	<b>76.8</b>	<b>74.0</b>	77.6	79.8	<b>76.2</b>	<b>74.4</b>	<b>77.8</b>	74.2	<b>77.6</b>	<b>82.6</b>	<b>89.6</b>	<b>78.2</b>

Table 11: XCOPA results (accuracy) across different input languages.

Model	de	fr	ru	zh	Avg.
XLM-R	76.1	72.3	62.3	60.8	<b>67.9</b>
INFOXLM	81.3	78.2	76.0	74.2	<b>77.4</b>
HICTL	80.5	79.2	76.0	74.8	<b>77.6</b>
EMMA-X	<b>85.1</b>	<b>82.8</b>	<b>81.3</b>	<b>78.3</b>	<b>81.9</b>

Table 12: BUCC results (F1) across different languages.

Model	en	de	es	fr	ja	ko	zh	Avg.
XLM-R	95.7	92.2	92.7	92.5	84.7	85.9	87.1	<b>90.1</b>
INFOXLM	<b>97.7</b>	94.6	<b>95.2</b>	95.1	88.9	89.0	90.2	<b>93.0</b>
HICTL	97.4	94.2	95.0	94.2	89.1	89.5	90.2	<b>92.8</b>
ChatGPT	71.9	67.8	67.9	67.0	58.3	54.7	61.4	<b>64.2</b>
EMMA-X	97.3	<b>95.6</b>	94.7	<b>96.0</b>	<b>92.9</b>	<b>89.8</b>	93.0	<b>94.2</b>

Table 13: PAWS-X results (accuracy) for each language.