

1) Cultural Bias in Knowledge Benchmarks. 2) DLAMA-v1: Culturally Diverse Benchmark. 3) Contrastive Sets of Facts Enrich Analysis.

DLAMA: A Framework for Curating Culturally Diverse Facts for Probing the Knowledge of Pretrained Language Models

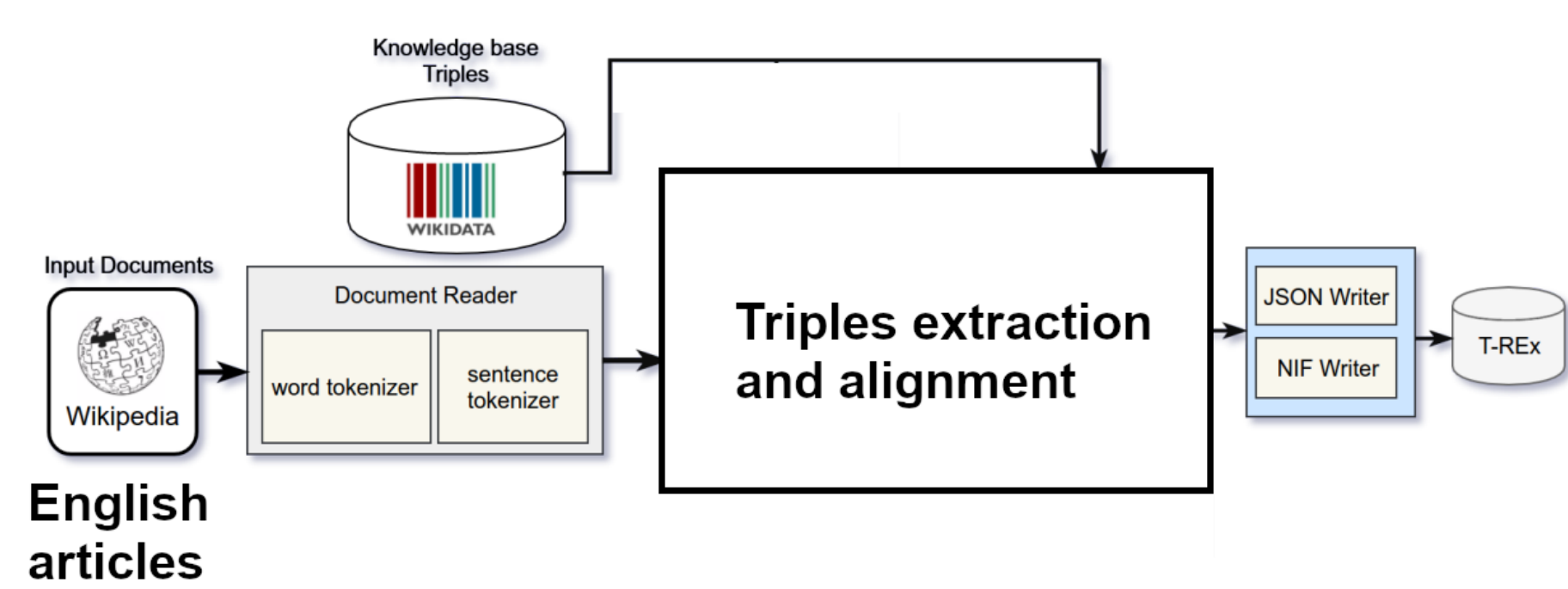
Amr Keleg, Walid Magdy

Factual Knowledge of mBERT

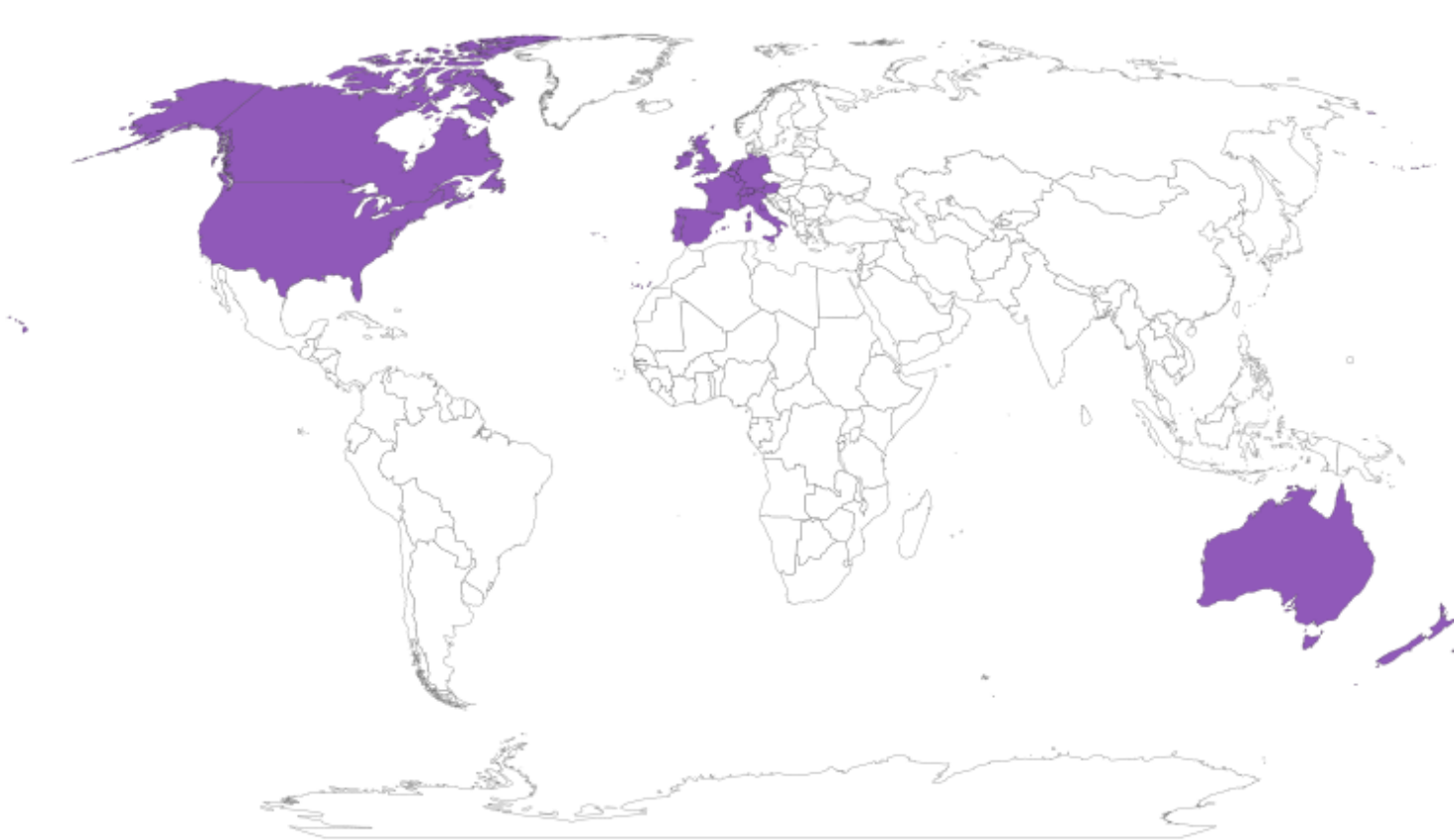
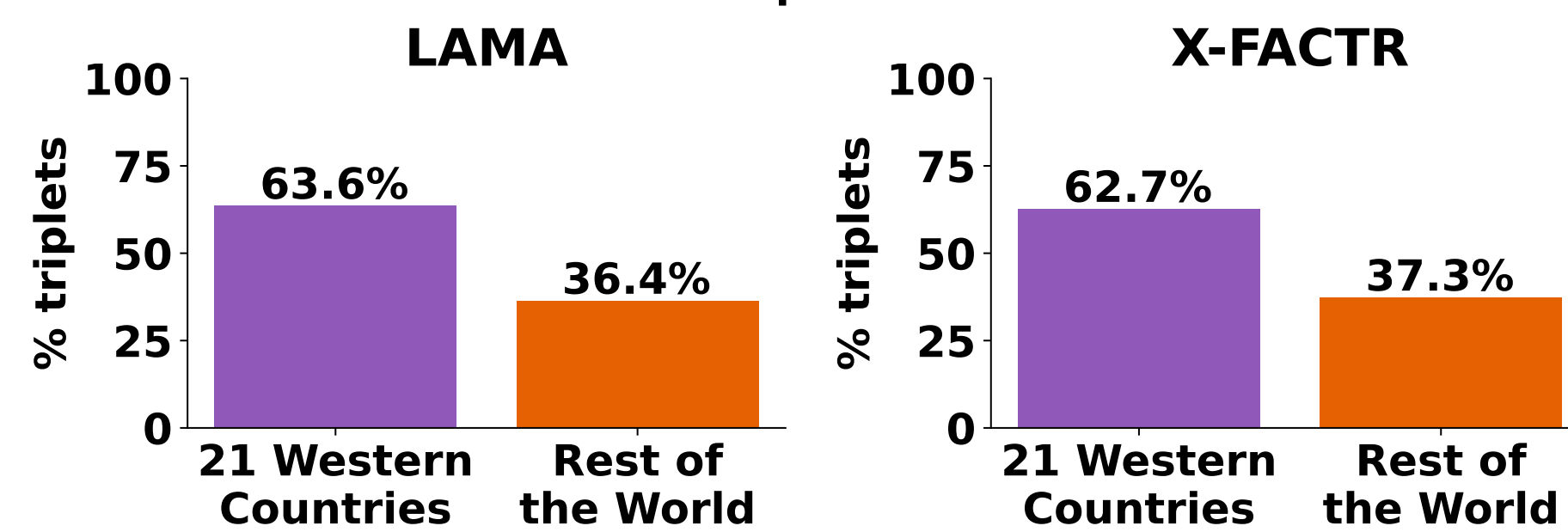
- Using prompts from mLAMA benchmark
- $ACC_{Arabic} \text{ prompts} < \frac{1}{2} ACC_{English} \text{ prompts}$
- Is mLAMA culturally representative?

Bias in Knowledge Benchmarks

- TRE-x creation steps (problems?):



- 2 benchmarks sampled from TRE-x:

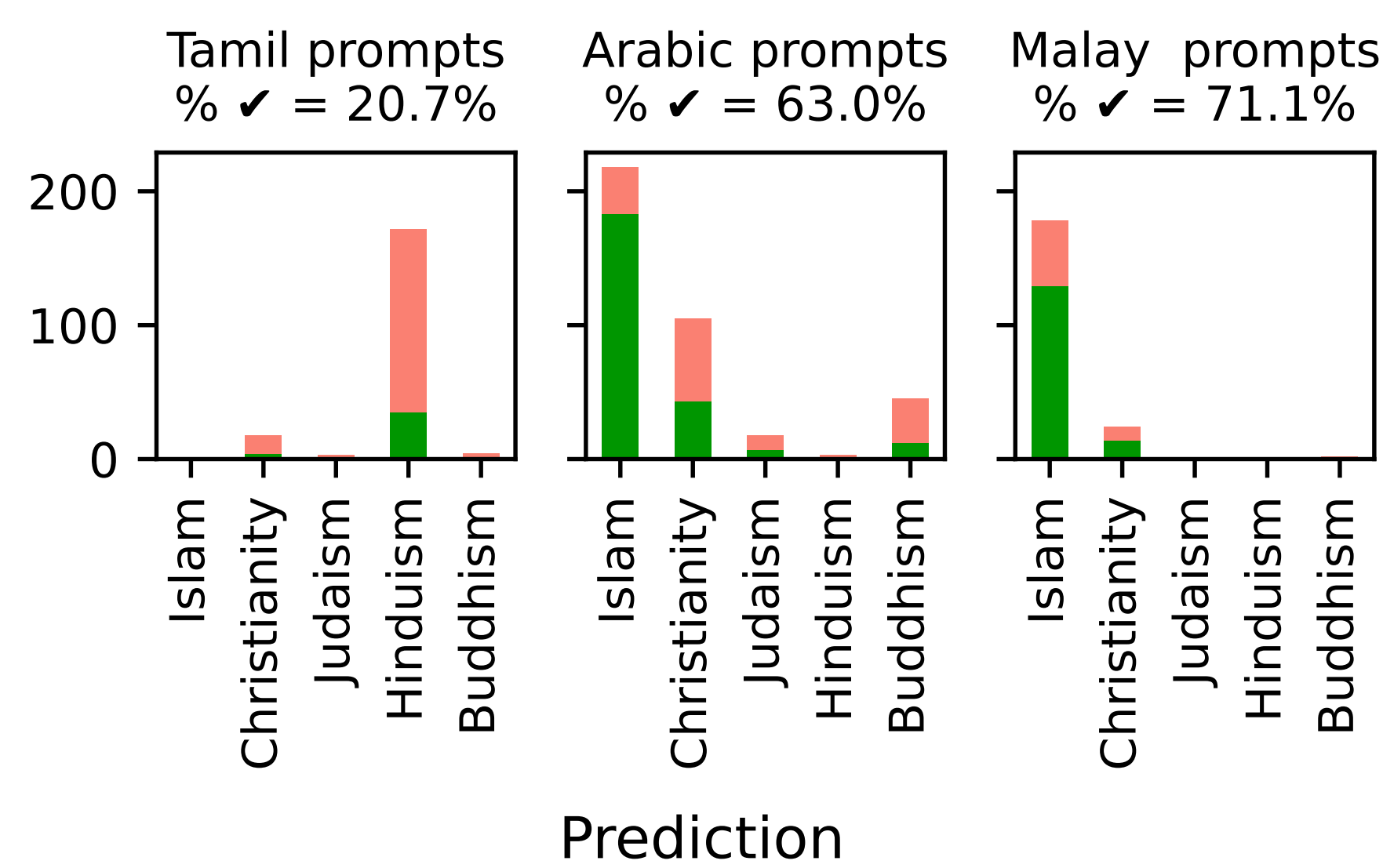


Probing Setup

- Starting from Wikidata triplets:
 - Subject: Edward I of England [X]
 - Predicate: religion or worldview [P140]
 - Object: Christianity [Y]
- Manual templates to form prompts:
 - Edward I of England is affiliated with ...
 - (a) Islam (b) Christianity (c) Judaism (d) Hinduism (e) Buddhism ...
- Choices: objects [Y] for the predicate.

Impact of Bias on Analysis

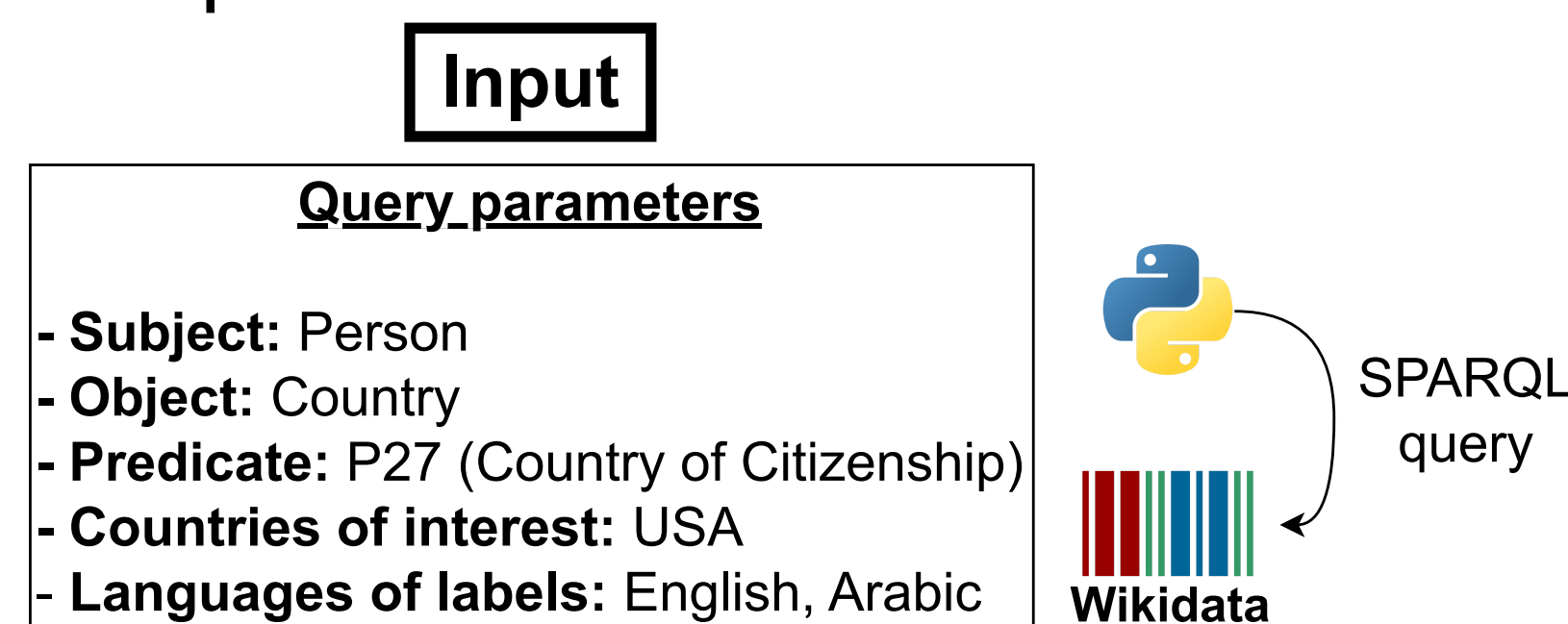
- Benchmark: mLAMA - P140
- Tamil, Arabic, Malay prompts
- Islam is the most common object.



- ⚠️ # prompts differs across languages.
- Malay prompts > Arabic prompts?

The DLAMA Framework

1. Set parameters



2. Query Wikidata triplets

subj uri	obj uri	subj wikipedia url	article size
Q81324	Q30	.../Bret_Hart	201353
Q81328	Q30	.../Harrison_Ford	81422
Q65645	Q30	.../Matthias_Pintscher	6923
...

3. Sort by article size

subj uri	obj uri	subj wikipedia url	article size
Q22686	Q30	.../Donald_Trump	417652
Q313381	Q30	.../Tom_Brady	408608
...

4. Query multilingual labels

English		Arabic	
subj label	obj label	subj label	obj label
Donald Trump	United States of America	دونالد ترامب	الولايات المتحدة
Tom Brady	United States of America	توم برايدي	الولايات المتحدة
...

The DLAMA-v1 Benchmark

- > 78K triplets from 20 predicates

- 3 bilingual contrasting sets of facts:
 - Western(en) ✖ Arab(ar)
 - Western(en) ✖ South American(es)
 - Western(en) ✖ East-Asian(ko)

Results

DLAMA-v1 (Arab-West)

Prompt Lang.	Model name	$P_{@1}$	
		Arab $N=10946$	West $N=13589$
Arabic	mBERT-base arBERT	13.7 33.6*	15.1* 23.0
English	mBERT-base BERT-base	21.2 27.5	37.7* 31.3*

- arBERT (Arab) on par with BERT-base (West)

DLAMA-v1 (South America-West)

Prompt Lang.	Model name	$P_{@1}$	
		S. America $N=13071$	West $N=13586$
Spanish	mBERT-base BETO	25.4 16.0	33.8* 26.5*
English	mBERT-base BERT-base	27.0 26.9	37.6* 31.3*

- Models not capturing S. American facts?

DLAMA-v1 (Asia-West)

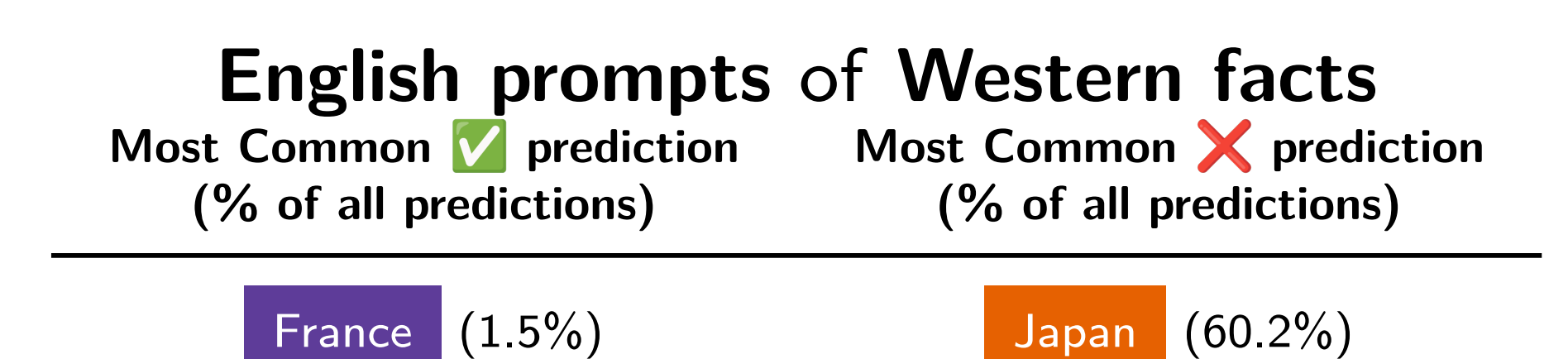
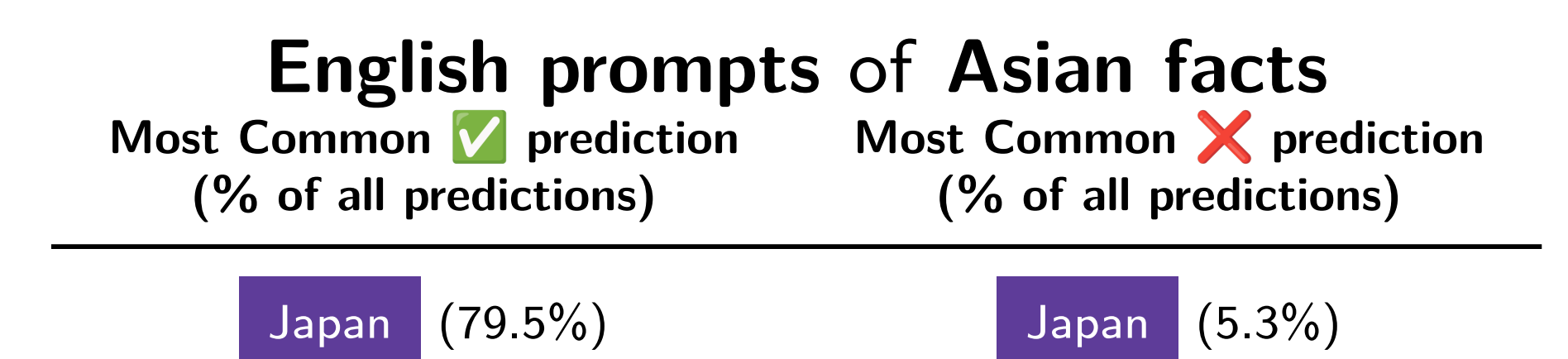
Prompt Lang.	Model name	$P_{@1}$	
		Asia $N=13479$	West $N=13588$
Korean	mBERT-base KyKim	16.4 22.1*	28.5* 19.5
English	mBERT-base BERT-base	33.0 38.3*	39.9* 31.9

- BERT-base (Asia) > BERT-base (West) ↓

Benefits of Contrastive Benchmarks

Model: English BERT-base

Prompt [P495]: [X] was created in [Y].



- Prompt bias toward Japan?



THE UNIVERSITY OF EDINBURGH
informatics



Institute for Language,
Cognition and Computation



UKRI CENTRE
FOR DOCTORAL
TRAINING