# Making Sentence Embeddings Robust to User-Generated Content

Lydia Nishimwe

Inria, France

lydia.nishimwe@inria.fr

Benoît Sagot

Inria, France

benoit.sagot@inria.fr

Rachel Bawden

Inria, France

rachel.bawden@inria.fr

# User-Generated Content (UGC)

**Ergographic phenomena (encoding simplification)**

i **don wanna fyt witchu**

**al b** an **our l8**    **c u 2moro**

**Neologisms**

The math is not **mathing**.

**burkini**

**Transverse phenomena**

i **aint playin**    **idk**

**afaik**    **N. E. V. E. R**

**Foreign language influence**

Cette fête a l'air **fun, let's go !**

**likez** et commentez

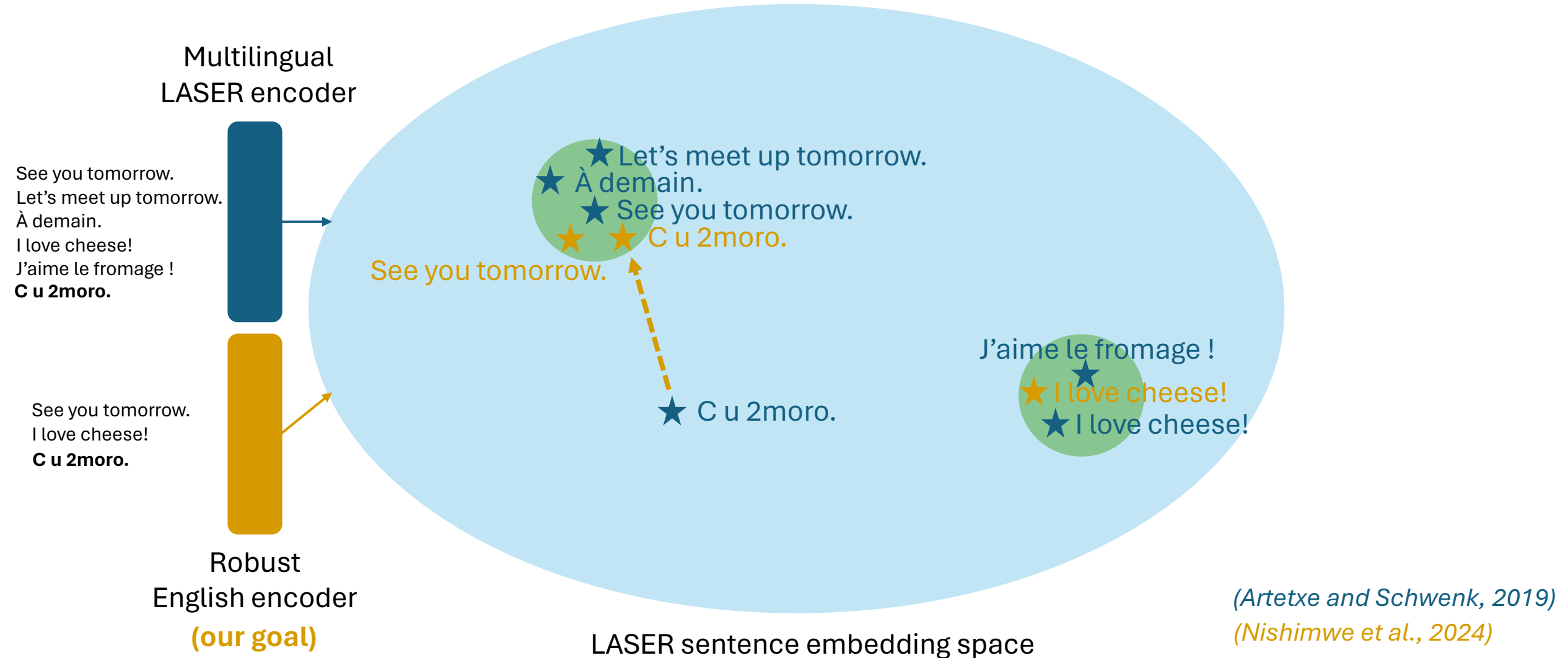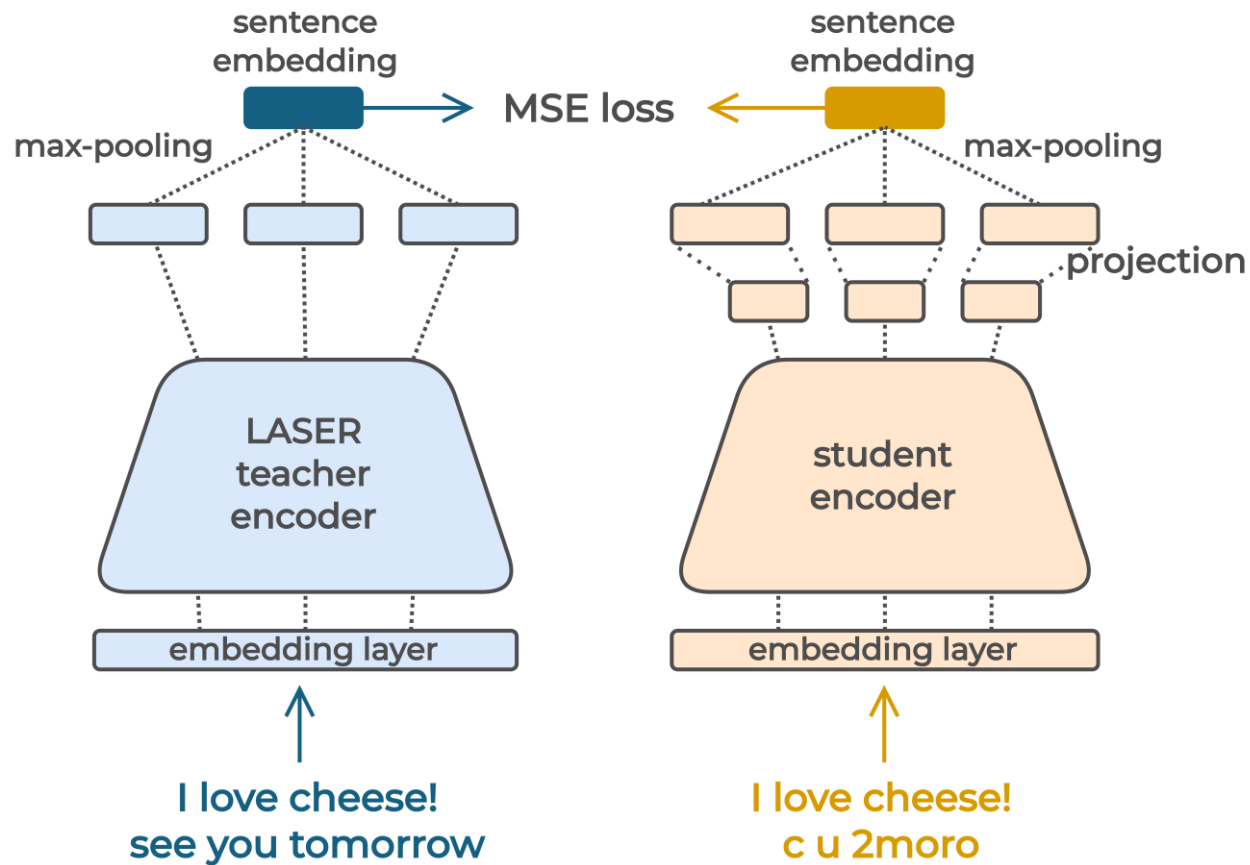**Marks of expressiveness**

supe**rrr !!!!**

<3    ☺    **!d10t**

**sh*t**

*(Seddah et al., 2012)*
*(Zalmout et al., 2019)*
*(Sanguinetti et al., 2020)*

# Multilingual sentence embeddings



Multilingual
LASER encoder

See you tomorrow.
Let's meet up tomorrow.
À demain.
I love cheese!
J'aime le fromage !
**C u 2moro.**

See you tomorrow.
I love cheese!
**C u 2moro.**

Robust
English encoder
**(our goal)**

Let's meet up tomorrow.
À demain.
See you tomorrow.
C u 2moro.
See you tomorrow.
C u 2moro.

J'aime le fromage !
I love cheese!
I love cheese!

LASER sentence embedding space

*(Artetxe and Schwenk, 2019)*
*(Nishimwe et al., 2024)*

3

# Proposed approach: Teacher-Student training



- LASER (teacher):
  - 45M parameters
  - 5-layer bi-LSTM
  - 1024 output dimension
  - fixed during training

- RoLASER [Robust LASER] (student):
  - 108M parameters
  - 12-layer Transformer
  - 768 output dimension
  - projection layer -> 1024

- c-RoLASER (student):
  - 104M parameters
  - same as RoLASER, except for
  - Character-CNN input embedding layer

*(Reimers and Gurevych, 2020; Duquenne et al., 2022; Mao and Nakagawa, 2023)*

# Generating artificial UGC (NL-Augmenter)

contractions and expansions

abbreviations, acronyms, slang

misspellings

**cont**   I am ↔ I'm

**abr1**   because → cuz

**fing**   tried → triwd

**week**   Monday ↔ Mon.

**abr2**   easy → ez

**homo**   there ↔ their

**abr3**   ASAP ↔ as soon as possible

visual and segmentation

**dysl**   lose ↔ loose

**slng**   jewellery → bling bling

**leet**   love → l0V3

**spel**   absent → apsent

**spac**   hello there → h elloth ere

*(Dhole et al., 2021)*

# Generating artificial UGC training data



p=0.1

| abr1 ✓ | abr2 ✓ | abr3 ✗ | cont ✗ | dysl ✗ | fing ✓ | homo ✗ | leet ✗ | slng ✗ | spac ✗ | spel ✗ | week ✗ |

| abr1 | abr2 | fing |

shuffle

| abr2 | fing | abr1 |

$1/4$ — fing (p=0.025)
$1/2$ — fing (p=0.05)
$1/4$ — fing (p=0.075)

$1/4$ — abr1 (p=0.05)
$1/2$ — abr1 (p=0.1)
$1/4$ — abr1 (p=0.15)

| abr2 | fing (p=0.05) | abr1 (p=0.15) |

mix_all

"Luckily **nothing** happened **to** me, but I saw a macabre scene, as **people tried to** break windows in order **to get** out."

"Luckily **nthing** happened **2** me, but I saw a macabre scene, as **ppl triwd 2** break windows in order **2 gt** out."

7

# Experimental setup

- **Training data:**
  - 2M "bilingual" standard-UGC sentences
  - 2M standard English sentences from the OSCAR dataset *(Ortiz Suárez et al., 2019)*
  - augmented with the *mix_all* transformation

- **RoLASER training:**
  - initialised with RoBERTa *(Liu et al., 2019)*
  - 98 epochs
- **c-RoLASER training:**
  - initialised with CharacterBERT *(El Boukkouri et al., 2020)*
  - 32 epochs

# Evaluation data and metrics

## Data

- MultiLexNorm *(van der Goot et al., 2021)*
  - Twitter
  - 1967 standard ↔ UGC sentences in English
- RoCS-MT *(Bawden and Sagot, 2023)*
  - Reddit
  - 1922 standard ↔ UGC sentences in English
  - translations into other 5 languages
- FLORES-200 *(NLLB Team et al., 2022)*
  - WikiNews, WikiBooks, WikiVoyage
  - parallel texts in 200 languages
  - 997 dev and 1012 test sentences
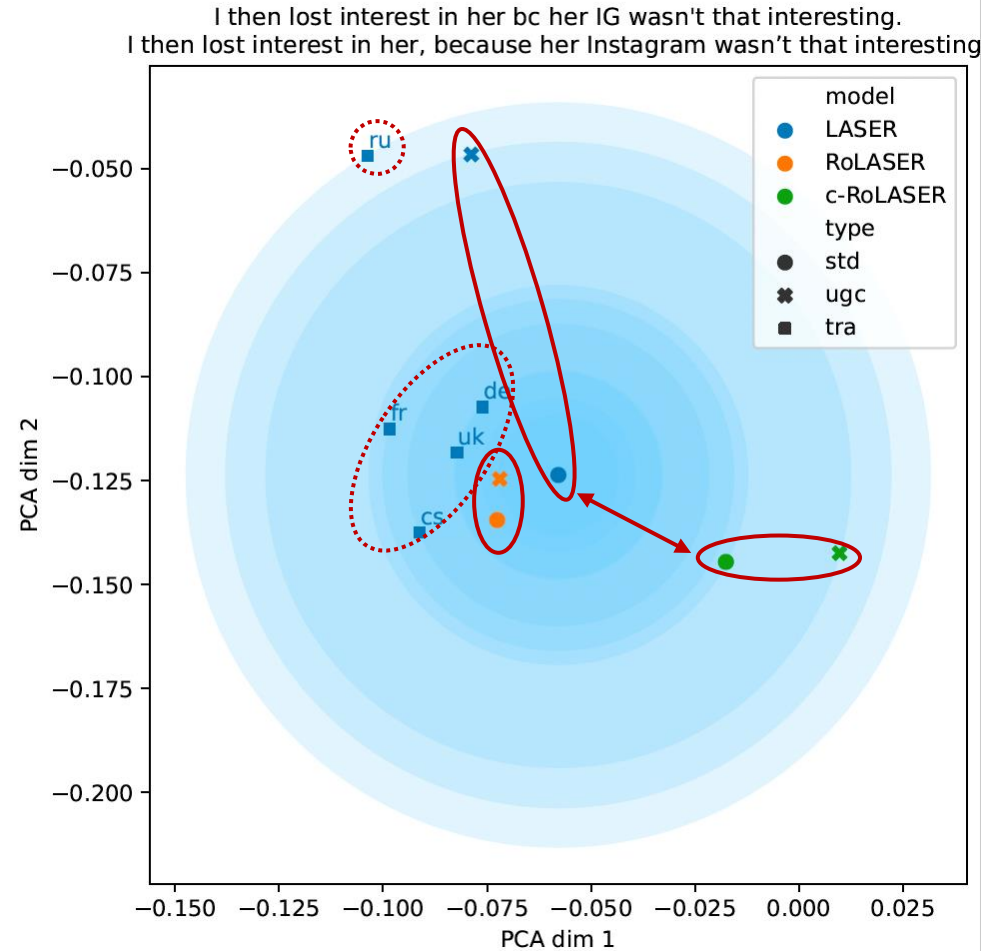  - artificially augmented

## Metrics

- Average pairwise cosine distance
- xSIM *(Artetxe and Schwenk, 2019)*
  - cross-lingual similarity search
  - proxy metric for bitext mining
  - error rate of aligning translation pairs
- xSIM++ *(Chen et al., 2023)*
  - augmenting the English set of FLORES-200
  - altering the meaning
  - minimal surface changes
  - more challenging than xSIM

# Evaluation on natural UGC



(lower is better)
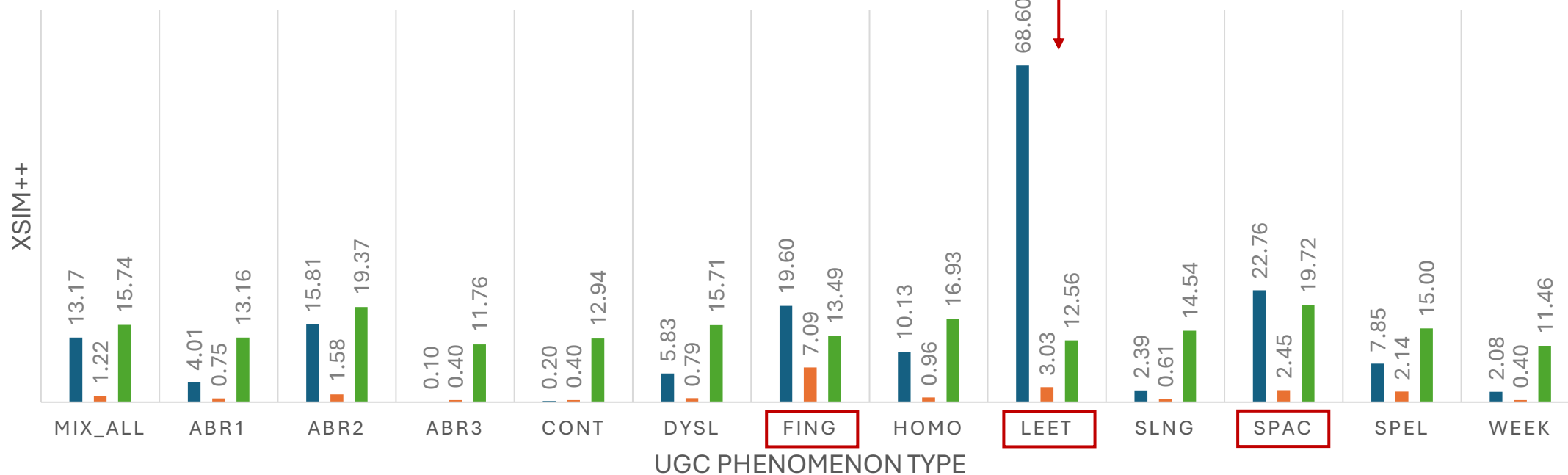
# Evaluation on artificial UGC

Hello world ➜ `__Hel` `lo` `__world`

H3ll0 w0rld ➜ `__H` `3` `ll` `0` `__w` `0` `r` `ld`

**FLORES-200**

■ LASER   ■ RoLASER   ■ c-RoLASER



XSIM++

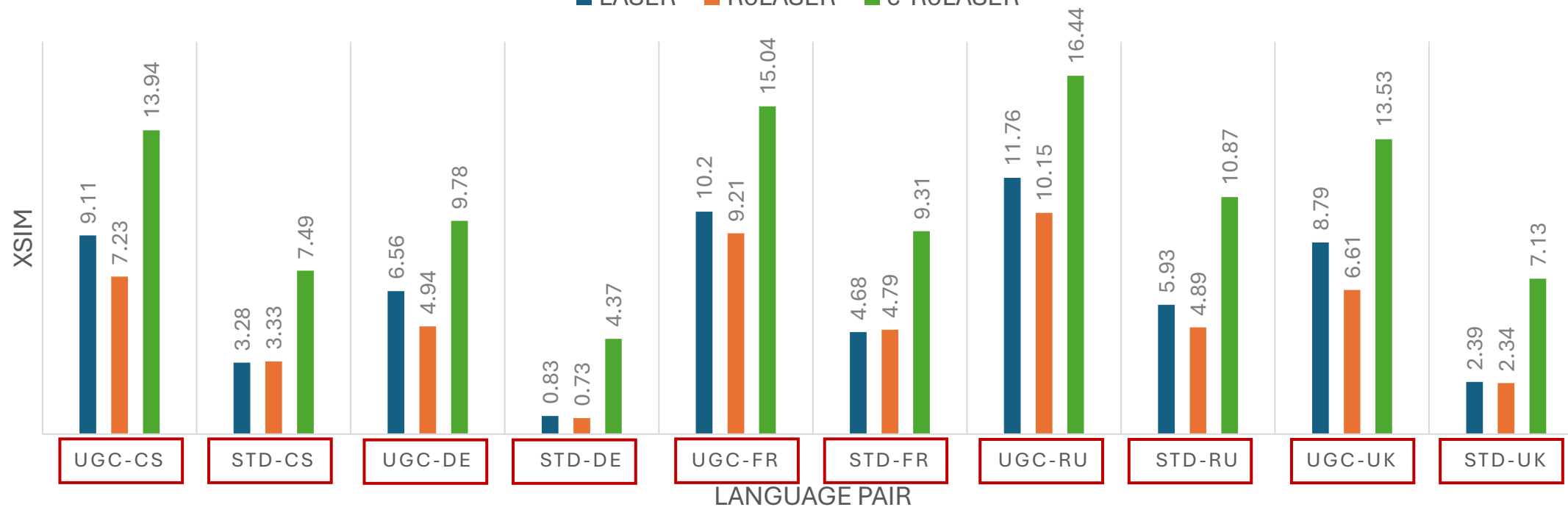| | MIX_ALL | ABR1 | ABR2 | ABR3 | CONT | DYSL | FING | HOMO | LEET | SLNG | SPAC | SPEL | WEEK |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| LASER | 13.17 | 4.01 | 15.81 | 0.10 | 0.20 | 5.83 | 19.60 | 10.13 | 68.60 | 2.39 | 22.76 | 7.85 | 2.08 |
| RoLASER | 1.22 | 0.75 | 1.58 | 0.40 | 0.40 | 0.79 | 7.09 | 0.96 | 3.03 | 0.61 | 2.45 | 2.14 | 0.40 |
| c-RoLASER | 15.74 | 13.16 | 19.37 | 11.76 | 12.94 | 15.71 | 13.49 | 16.93 | 12.56 | 14.54 | 19.72 | 15.00 | 11.46 |

UGC PHENOMENON TYPE

(lower is better)

12

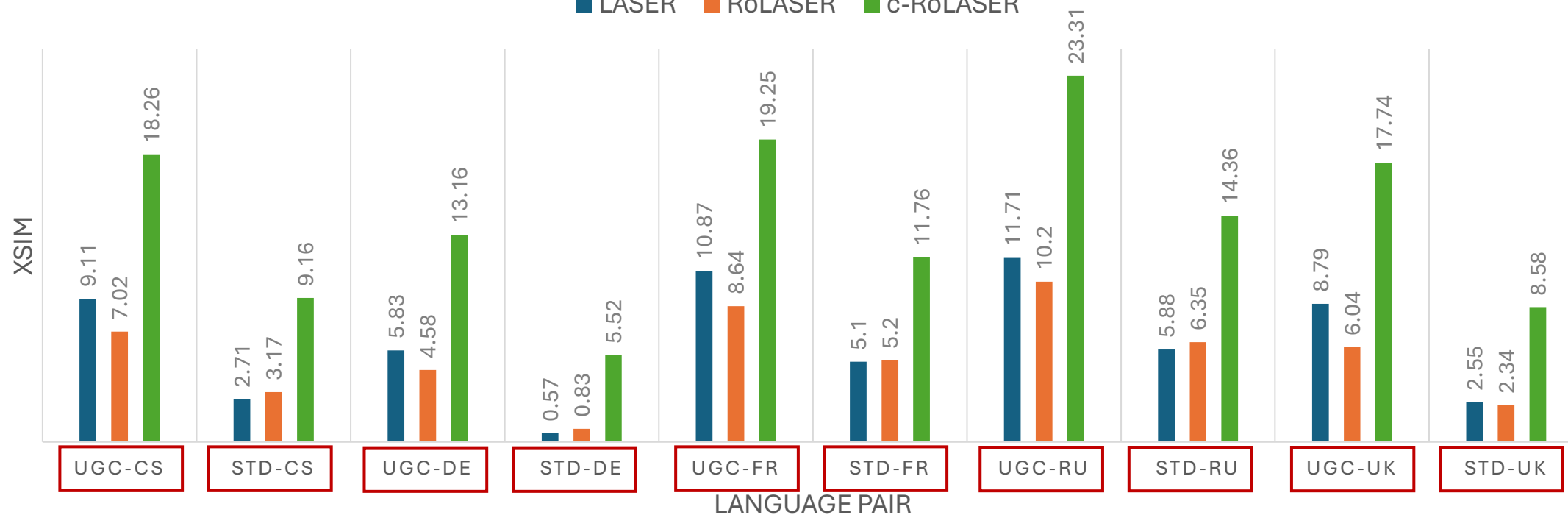# Evaluation on UGC and standard data in a multilingual setting (1)

**ROCS-MT ENGLISH→XX**



(lower is better)

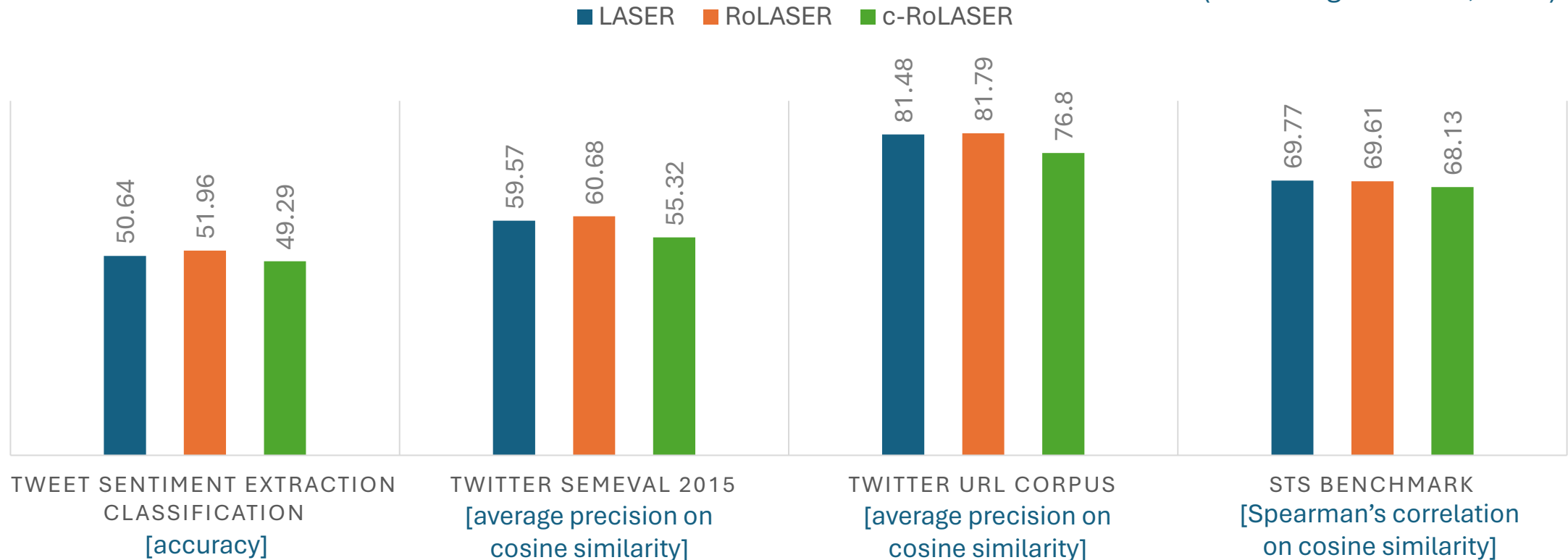# Evaluation on UGC and standard data in a multilingual setting (2)



**ROCS-MT XX→ENGLISH**

(lower is better)

# Evaluation on downstream tasks



**MTEB: MASSIVE TEXT EMBEDDING BENCHMARK**

(Muenninghoff et al., 2023)

■ LASER  ■ RoLASER  ■ c-RoLASER

| | LASER | RoLASER | c-RoLASER |
|---|---|---|---|
| TWEET SENTIMENT EXTRACTION CLASSIFICATION [accuracy] | 50.64 | 51.96 | 49.29 |
| TWITTER SEMEVAL 2015 [average precision on cosine similarity] | 59.57 | 60.68 | 55.32 |
| TWITTER URL CORPUS [average precision on cosine similarity] | 81.48 | 81.79 | 76.8 |
| STS BENCHMARK [Spearman's correlation on cosine similarity] | 69.77 | 69.61 | 68.13 |

(higher is better)

# Takeaways

**Approach:**

Making LASER more robust to UGC English

1. Teacher-Student training
2. Minimising the standard-UGC distance in the embedding space
3. Generating and training on synthetic UGC-like data

Extending RoLASER to **more languages** and their corresponding UGC phenomena...

Future work

**Results:**

RoLASER is significantly more robust than LASER
- on natural and artificial UGC
- on standard data and downstream tasks (improves/matches LASER's performance)

**Findings:**

1. c-RoLASER struggles to map its standard embeddings to LASER's
2. Most challenging UGC phenomena: character-level perturbations that shatter subword tokenisation

*Read our paper for more details and references!*

19