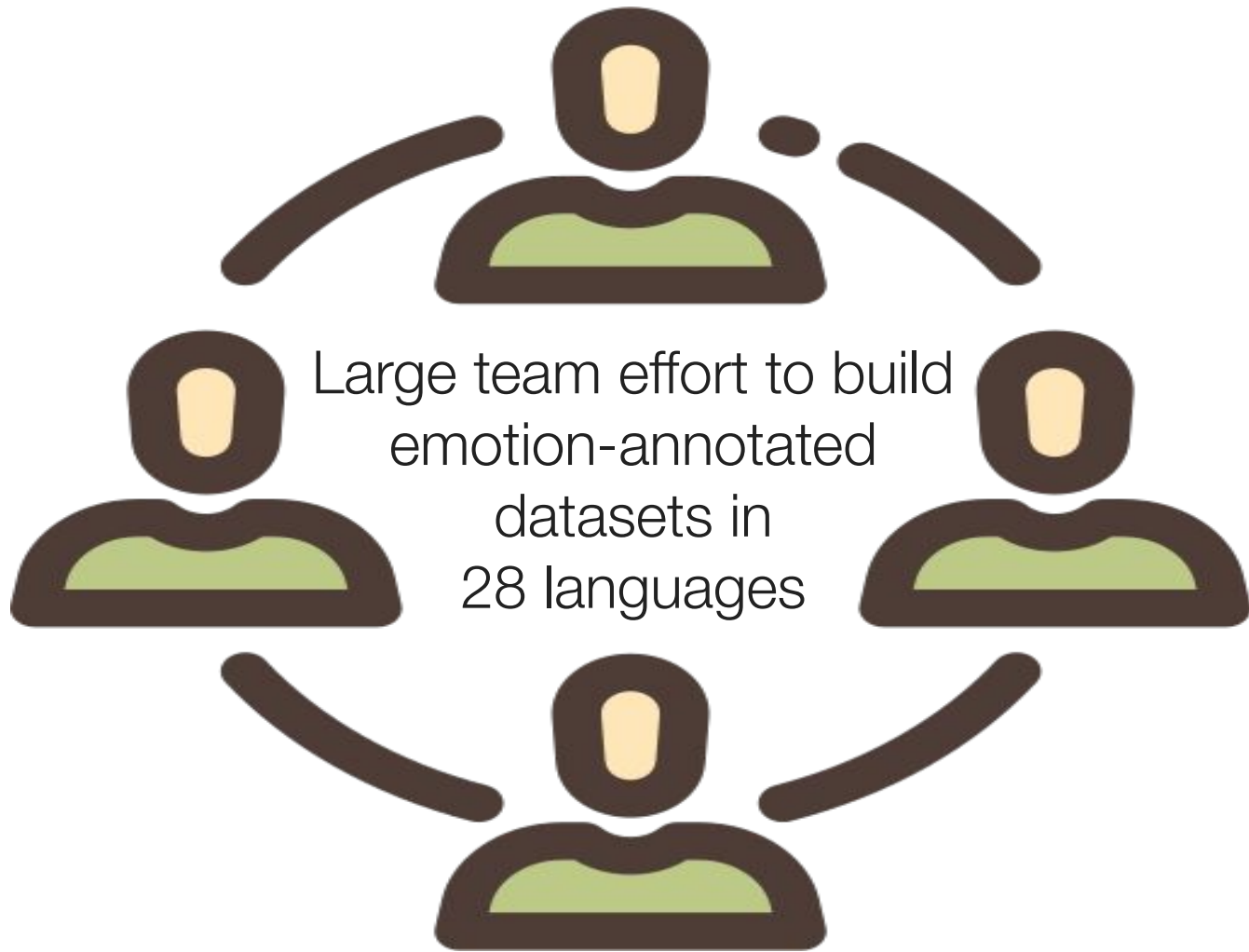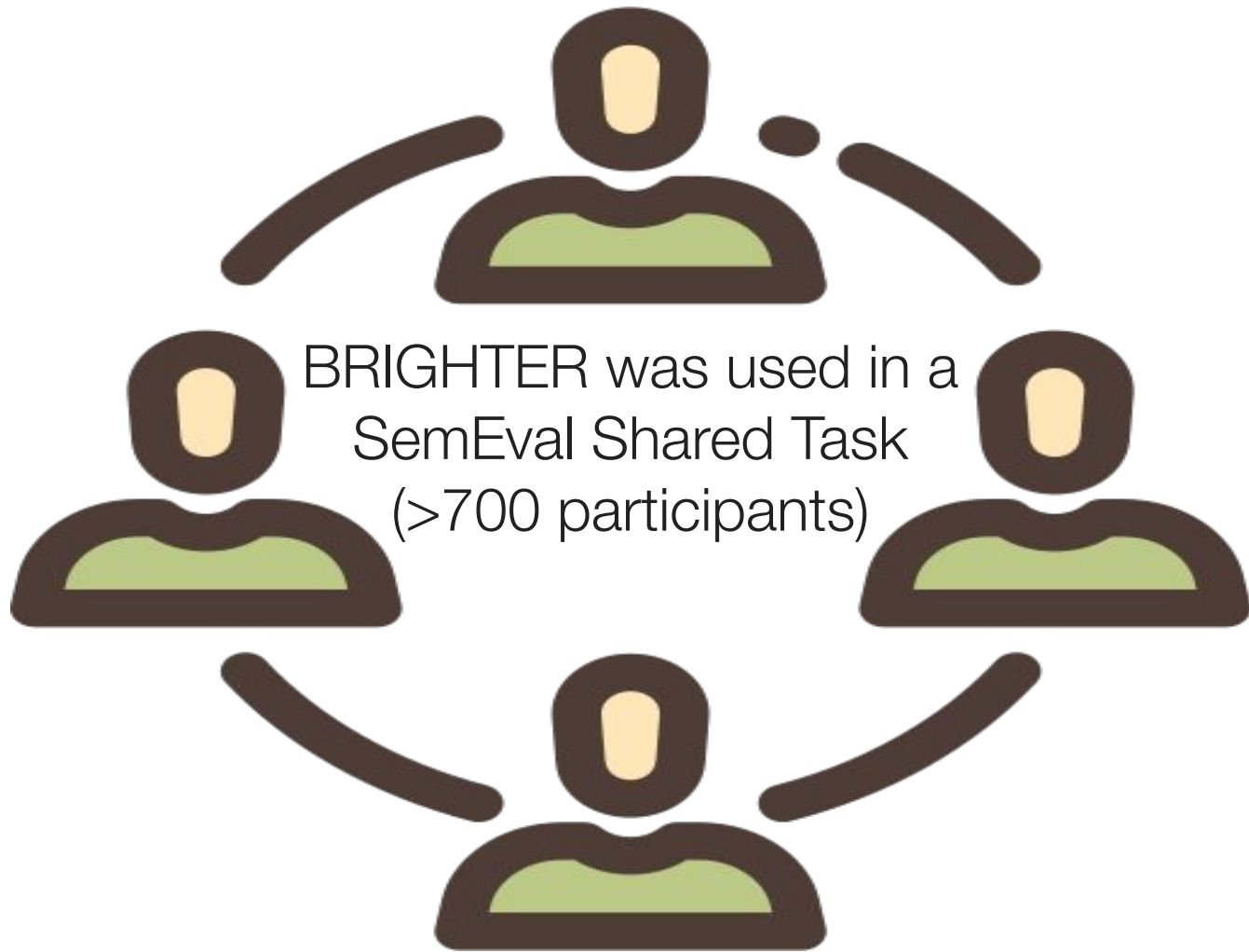# ✨ BRIGHTER ✨

ACL 2025
VIENNA
JULY 27 - AUGUST 1

# BRIdging the Gap in Human-Annotated Textual Emotion Recognition Datasets for 28 Languages

S. H. Muhammad*, **N. Ousidhoum***, I. Abdulmumin, J. P. Wahle, T. Ruas, M. Beloucif, C. de Kock, N. Surange, D. Teodorescu, I. S. Ahmad, D. I. Adelani, A. F. Aji, F. D. M. A. Ali, I. Alimova, V. Araujo, N. Babakov, N. Baes, A.-M. Bucur, A. Bukula, G. Cao, R. Tufiño, R. Chevi, C. I. Chukwuneke, A. Ciobotaru, D. Dementieva, M. S. Gadanya, R. Geislinger, B. Gipp, O. Hourrane, O. Ignat, F. I. Lawan, R. Mabuya, R. Mahendra, V. Marivate, A. Piper, A. Panchenko, C. H. Porto Ferreira, V. Protasov, S. Rutunda, M. Shrivastava, A. C. Udrea, L. D. A. Wanzare, S. Wu, F. V. Wunderlich, H. M. Zhafran, T. Zhang, Y. Zhou, S. M. Mohammad.

https://brighter-dataset.github.io

Large team effort to build
emotion-annotated
datasets in
28 languages

BRIGHTER was used in a SemEval Shared Task (>700 participants)

# BRIGHTER: Coverage

BRIGHTER primarily covers low-resource languages from Africa, Asia, Eastern Europe, Latin America

# BRIGHTER: Coverage of 28 languages

# BRIGHTER: Emotion Recognition Datasets

- BRIGHTER focuses on **perceived emotions**

  - I.e., emotion(s) most people think the speaker might have felt given a text snippet uttered by them

- The datasets are multi-labeled

# Dataset Construction
## Data Collection

- We targeted emotionally rich text (e.g., personal narratives)

- Eventually, we used various sources depending on the availability text data

  - Social media (e.g., Reddit in English, German, Romanian, others)

  - Speeches (e.g., in Afrikaans)

  - A translated novel (e.g., in Algerian Arabic)

  - News data and combined sources when text sources are scarce

# Dataset Construction

## Data Annotation

Given a text snippet, we asked the annotators to select all the emotions that apply on a 4-level intensity scale



If no emotion is selected then the text is considered neutral

# BRIGHTER: Multi-labeled Datasets



**Anger**

Isso é um tipo de violência contra a MULHER MINISTRA DAMARES.

This is a kind of violence against the WOMAN MINISTER DAMARES. Action must be taken.

Intensity:

**Sadness**

الدراري لي دايرين بيه يعيطووا و يحللوا "إدّا حليم" ،"!عطيتنيش ما وانا" ،"!وانا وانا انا ماشي لالا" ،"!كلش!".

The kids who were surrounding him were yelling and crying: "And me! and me!" "You did not give anything!" "Halim took everything!" "No not me!"

Intensity:

**Disgust** **Surprise**

Me eh convertido en lo que siempre jure destruir...

I have become what I always swore to destroy...

Intensity:

**BRIGHTER ✨✨✨**

**Fear**

Osareya maasi onnopopiha o suuli wa Mosampiikhi

Floods now threaten southern Mozambique

Intensity:

**Disgust**

शी! तुझ्या टॉवेलला वास येतो आहे!!

Eww! Your towel stinks!!

Intensity:

**Joy**

Але, божечко, як добре вдома.

Oh, God, how great to be at home.

Intensity:

# Dataset Construction
## Quality Control

**Intensity scores**
Intensity scores are kept for datasets >=5 annotators per instance (i.e., 10 languages)

**Pre-processing**
Text is processed by native speakers.

**Label determination**
Final labels are chosen based on agreement and intensity score threshold

**Annotation**
Annotators are native speakers, >=3 annotators per instance

**Reliability scores of the final datasets**
Reliability Scores >62%

# BRIGHTER: Final Datasets



Grouped bar charts showing counts and percentages of emotion labels (Anger, Disgust, Fear, Joy, Neutral, Sadness, Surprise) across languages (zul, yor, xho, vmw, ukr, tat, swe, swa, sun, rus, ron, ptmz, ptbr, pcm, mar, kin, jav, ind, ibo, hin, hau, esp, eng, deu, chn, ary, arq, afr).

# BRIGHTER: Final Datasets

Annotated a total of >100k instances

# BRIGHTER: Final Datasets

**Anger** | **Disgust** | **Fear** | **Joy** | **Neutral** | **Sadness** | **Surprise**

Dataset sizes vary from >1,6k instances to >7,6k

Languages (y-axis): zul, yor, xho, vmw, ukr, tat, swe, swa, sun, rus, ron, ptmz, ptbr, pcm, mar, kin, jav, ind, ibo, hin, hau, esp, eng, deu, chn, ary, arq, afr

**Anger**
- zul: 240 (1.4%)
- yor: 325 (1.9%)
- xho: 93 (0.5%)
- vmw: 151 (0.9%)
- rus: 816 (4.8%)
- ron: 440 (2.6%)
- ptmz: 151 (0.9%)
- ptbr: 1475 (8.7%)
- pcm: 561 (3.3%)
- mar: 528 (3.1%)
- kin: 743 (4.4%)
- jav: 205 (1.2%)
- ind: 209 (1.2%)
- ibo: 965 (5.7%)
- hin: 599 (3.5%)
- hau: 684 (4.0%)
- esp: 929 (5.5%)
- eng: 671 (4.0%)
- deu: 1605 (9.5%)
- chn: 2432 (14.3%)
- ary: 636 (3.8%)
- arq: 620 (3.7%)
- afr: 132 (0.8%)

**Disgust**
- zul: 84 (0.6%)
- yor: 135 (1.0%)
- xho: 11 (0.1%)
- vmw: 107 (0.8%)
- rus: 421 (3.1%)
- ron: 600 (4.5%)
- ptmz: 107 (0.8%)
- ptbr: 151 (1.1%)
- pcm: 3144 (23.4%)
- mar: 408 (3.0%)
- kin: 208 (1.5%)
- jav: 101 (0.8%)
- ind: 196 (1.5%)
- ibo: 898 (6.7%)
- hin: 386 (2.9%)
- hau: 552 (4.1%)
- esp: 1275 (9.5%)
- eng: 0 (0.0%)
- deu: 1718 (12.8%)
- chn: 852 (6.3%)
- ary: 116 (0.9%)
- arq: 436 (3.2%)
- afr: 88 (0.7%)

**Fear**
- zul: 94 (0.8%)
- yor: 130 (1.1%)
- xho: 44 (0.4%)
- vmw: 183 (1.6%)
- rus: 457 (4.0%)
- ron: 834 (7.3%)
- ptmz: 182 (1.6%)
- ptbr: 256 (2.2%)
- pcm: 676 (5.9%)
- mar: 544 (4.8%)
- kin: 236 (2.1%)
- jav: 113 (1.0%)
- ind: 79 (0.7%)
- ibo: 366 (3.2%)
- hin: 540 (4.7%)
- hau: 549 (4.8%)
- esp: 549 (4.8%)
- eng: 3218 (28.2%)
- deu: 506 (4.4%)
- chn: 150 (1.3%)
- ary: 202 (1.8%)
- arq: 465 (4.1%)
- afr: 252 (2.2%)

**Joy**
- zul: 172 (0.8%)
- yor: 454 (2.1%)
- xho: 707 (3.3%)
- vmw: 332 (1.6%)
- rus: 782 (3.7%)
- ron: 904 (4.3%)
- ptmz: 332 (1.6%)
- ptbr: 1198 (5.7%)
- pcm: 951 (4.5%)
- mar: 655 (3.1%)
- kin: 704 (3.3%)
- jav: 356 (1.7%)
- ind: 495 (2.3%)
- ibo: 778 (3.7%)
- hin: 644 (3.0%)
- hau: 535 (2.5%)
- esp: 1237 (5.8%)
- eng: 1375 (6.5%)
- deu: 1139 (5.4%)
- chn: 1103 (5.2%)
- ary: 504 (2.4%)
- arq: 326 (1.5%)
- afr: 958 (4.5%)

**Neutral**
- zul: 1688 (6.8%)
- yor: 2364 (9.6%)
- xho: 332 (1.3%)
- vmw: 1195 (4.8%)
- ukr: 2461 (10.0%)
- rus: 203 (0.8%)
- ron: 1192 (4.8%)
- ptmz: 1294 (5.2%)
- ptbr: 306 (1.2%)
- pcm: 625 (2.5%)
- mar: 1180 (4.8%)
- kin: 190 (0.8%)
- jav: 65 (0.3%)
- ibo: 1030 (4.2%)
- hin: 781 (3.2%)
- hau: 461 (1.9%)
- esp: 0 (0.0%) / 545 (2.2%)
- eng: 1297 (5.3%)
- deu: 1224 (5.0%)
- chn: 757 (3.1%)
- ary: 196 (0.8%)
- arq: 849 (4.1%)
- afr: 892 (3.6%)

**Sadness**
- zul: 601 (2.9%)
- yor: 1393 (6.7%)
- xho: 939 (4.5%)
- vmw: 463 (2.2%)
- ukr: 665 (3.2%)
- rus: 755 (3.6%)
- ron: 460 (2.2%)
- ptmz: 693 (3.3%)
- ptbr: 1648 (7.9%)
- pcm: 655 (3.2%)
- mar: 1053 (5.1%)
- kin: 317 (1.5%)
- jav: 296 (1.4%)
- ind: 822 (4.0%)
- ibo: 635 (3.1%)
- hin: 1084 (5.2%)
- hau: 586 (2.8%)
- esp: 1794 (8.6%)
- eng: 1083 (5.2%)
- deu: 762 (3.7%)
- chn: 415 (2.0%)

**Surprise**
- zul: 263 (1.9%)
- yor: 424 (3.1%)
- xho: 363 (2.6%)
- vmw: 167 (1.2%)
- ukr: 390 (2.8%)
- (6.5%)
- rus: 908 (6.6%)
- ron: 167 (1.2%)
- ptmz: 331 (2.4%)
- ptbr: 1643 (12.0%)
- pcm: 449 (3.3%)
- mar: 175 (1.3%)
- kin: 363 (2.6%)
- jav: 316 (2.3%)
- ind: 128 (0.9%)
- ibo: 434 (3.2%)
- hin: 583 (4.3%)
- hau: 855 (6.2%)
- esp: 1669 (12.2%)
- eng: 336 (2.5%)
- deu: 388 (2.8%)
- chn: 414 (3.0%)
- ary: 650 (4.7%)
- arq: 0 (0.0%)

x-axis: Count

# BRIGHTER: Final Datasets



| | Anger | Disgust | Fear | Joy | Neutral | Sadness | Surprise |
|---|---|---|---|---|---|---|---|
| zul | 240 (1.4%) | 84 (0.6%) | 94 (0.8%) | 172 (0.8%) | 1688 (6.8%) | 601 (2.9%) | 263 (1.9%) |
| yor | 325 (1.9%) | 135 (1.0%) | 130 (1.1%) | 454 (2.1%) | 2364 (9.6%) | 1393 (6.7%) | 424 (3.1%) |
| xho | 93 (0.5%) | 11 (0.1%) | 44 (0.4%) | 707 (3.3%) | 332 (1.3%) | 939 (4.5%) | 363 (2.6%) |
| vmw | 151 (0.9%) | 107 (0.8%) | 183 (1.6%) | 332 (1.6%) | 1195 (4.8%) | 463 (2.2%) | 167 (1.2%) |

**Variation** in **class distribution** across datasets given the different **data sources**

| | | | | | | | |
|---|---|---|---|---|---|---|---|
| ptbr | | | | | | | 1643 (12.0%) |
| pcm | 561 (3.3%) | 3144 (23.4%) | 676 (5.9%) | 951 (4.5%) | | 655 (3.2%) | 449 (3.3%) |
| mar | 528 (3.1%) | 408 (3.0%) | 544 (4.8%) | 655 (3.1%) | 625 (2.5%) | 1053 (5.1%) | 175 (1.3%) |
| kin | 743 (4.4%) | 208 (1.5%) | 236 (2.1%) | 704 (3.3%) | 1180 (4.8%) | 317 (1.5%) | 363 (2.6%) |
| jav | 205 (1.2%) | 101 (0.8%) | 113 (1.0%) | 356 (1.7%) | 190 (0.8%) | 296 (1.4%) | 316 (2.3%) |
| ind | 209 (1.2%) | 196 (1.5%) | 79 (0.7%) | 495 (2.3%) | 65 (0.3%) | 822 (4.0%) | 128 (0.9%) |
| ibo | 965 (5.7%) | 898 (6.7%) | 366 (3.2%) | 778 (3.7%) | 1030 (4.2%) | 635 (3.1%) | 434 (3.2%) |
| hin | 599 (3.5%) | 386 (2.9%) | 540 (4.7%) | 644 (3.0%) | 781 (3.2%) | 1084 (5.2%) | 583 (4.3%) |
| hau | 684 (4.0%) | 552 (4.1%) | 549 (4.8%) | 535 (2.5%) | 461 (1.9%) | 586 (2.8%) | 855 (6.2%) |
| esp | 929 (5.5%) | 1275 (9.5%) | 549 (4.8%) | 1237 (5.8%) | 0 (0.0%) | 1794 (8.6%) | 1669 (12.2%) |
| eng | 671 (4.0%) | 0 (0.0%) | 3218 (28.2%) | 1375 (6.5%) | 545 (2.2%) | 1083 (5.2%) | 336 (2.5%) |
| deu | 1605 (9.5%) | 1718 (12.8%) | 506 (4.4%) | 1139 (5.4%) | 1297 (5.3%) | 762 (3.7%) | 388 (2.8%) |
| chn | 2432 (14.3%) | 852 (6.3%) | 150 (1.3%) | 1103 (5.2%) | 1224 (5.0%) | 415 (2.0%) | 414 (3.0%) |
| ary | 636 (3.8%) | 116 (0.9%) | 202 (1.8%) | 504 (2.4%) | 757 (3.1%) | 849 (4.1%) | 650 (4.7%) |
| arq | 620 (3.7%) | 436 (3.2%) | 465 (4.1%) | 326 (1.5%) | 196 (0.8%) | 415 (2.0%) | 0 (0.0%) |
| afr | 132 (0.8%) | 88 (0.7%) | 252 (2.2%) | 958 (4.5%) | 892 (3.6%) | 365 (1.8%) | |

# Experiments
## Multi-label Emotion Classification

The results were highly language-dependent
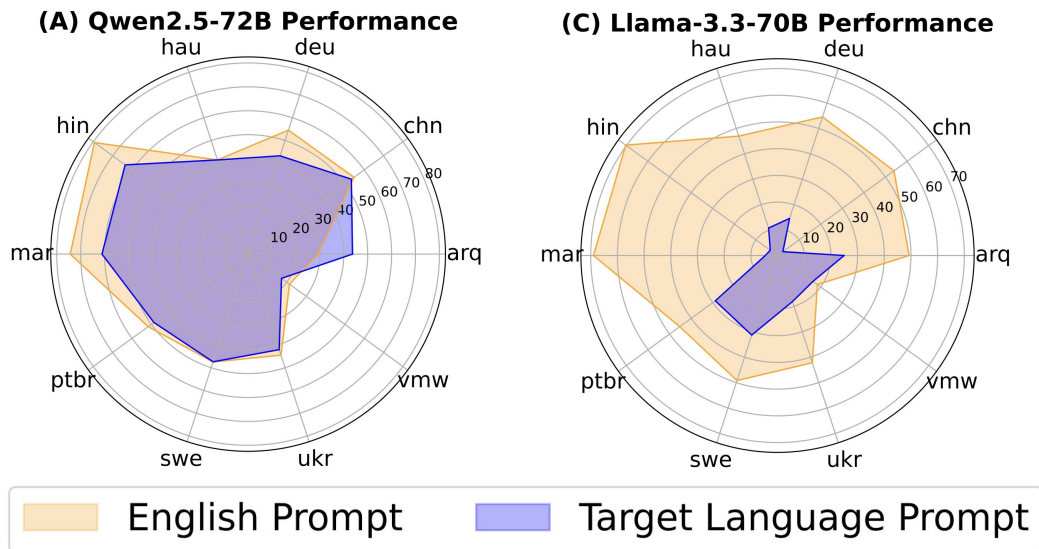
| | afr | arq | ary | chn | deu | eng | esp | hau | hin | ibo | ind | jav | kin | mar | pcm | ptbr | ptmz | ron | rus | sun | swa | swe |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| **Qwen** | 60.18 | 37.78 | 52.76 | 55.23 | 59.17 | 55.72 | 72.33 | 43.79 | 79.73 | 37.4 | 57.29 | 50.47 | 31.96 | 74.58 | 38.66 | 51.6 | 40.44 | 68.18 | 73.08 | 42.67 | 27.36 | 48.89 |
| **Dolly** | 23.58 | 38.59 | 24.27 | 27.52 | 26.86 | 42.6 | 36.41 | 29.43 | 27.59 | 24.31 | 36.61 | 36.18 | 19.73 | 25.69 | 34.41 | 25.9 | 16.7 | 43.58 | 29.72 | 32.2 | 17.63 | 21.79 |
| **Llama** | 61.28 | 55.75 | 44.96 | 53.36 | 56.99 | 65.58 | 61.27 | 50.91 | 60.59 | 33.18 | 39.2 | 41.88 | 34.36 | 67.4 | 48.67 | 45.03 | 34.06 | 71.28 | 62.61 | 46.33 | 29.47 | 50.26 |
| **Mixtral** | 53.69 | 45.29 | 35.07 | 44.91 | 51.2 | 58.12 | 65.72 | 40.4 | 62.19 | 31.9 | 54.37 | 48.37 | 26.35 | 50.36 | 45.61 | 41.64 | 36.52 | 68.51 | 61.72 | 42.1 | 26.51 | 48.61 |
| **Deep Seek** | 43.66 | 50.87 | 47.21 | 53.45 | 54.26 | 56.99 | 73.29 | 51.91 | 76.91 | 32.85 | 49.51 | 43.05 | 32.52 | 76.68 | 45 | 51.49 | 39.58 | 65.02 | 76.97 | 44.61 | 33.27 | 44.6 |

# Experiments
## Multi-label Emotion Classification

Qwen2.5-72B performed the best on average

| | afr | arq | ary | chn | deu | eng | esp | hau | hin | ibo | ind | jav | kin | mar | pcm | ptbr | ptmz | ron | rus | sun | swa | swe |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| **Qwen** | 60.18 | 37.78 | 52.76 | 55.23 | 59.17 | 55.72 | 72.33 | 43.79 | 79.73 | 37.4 | 57.29 | 50.47 | 31.96 | 74.58 | 38.66 | 51.6 | 40.44 | 68.18 | 73.08 | 42.67 | 27.36 | 48.89 |
| **Dolly** | 23.58 | 38.59 | 24.27 | 27.52 | 26.86 | 42.6 | 36.41 | 29.43 | 27.59 | 24.31 | 36.61 | 36.18 | 19.73 | 25.69 | 34.41 | 25.9 | 16.7 | 43.58 | 29.72 | 32.2 | 17.63 | 21.79 |
| **Llama** | 61.28 | 55.75 | 44.96 | 53.36 | 56.99 | 65.58 | 61.27 | 50.91 | 60.59 | 33.18 | 39.2 | 41.88 | 34.36 | 67.4 | 48.67 | 45.03 | 34.06 | 71.28 | 62.61 | 46.33 | 29.47 | 50.26 |
| **Mixtral** | 53.69 | 45.29 | 35.07 | 44.91 | 51.2 | 58.12 | 65.72 | 40.4 | 62.19 | 31.9 | 54.37 | 48.37 | 26.35 | 50.36 | 45.61 | 41.64 | 36.52 | 68.51 | 61.72 | 42.1 | 26.51 | 48.61 |
| **Deep Seek** | 43.66 | 50.87 | 47.21 | 53.45 | 54.26 | 56.99 | 73.29 | 51.91 | 76.91 | 32.85 | 49.51 | 43.05 | 32.52 | 76.68 | 45 | 51.49 | 39.58 | 65.02 | 76.97 | 44.61 | 33.27 | 44.6 |

# Experiments
## Sensitivity to the Language of the Prompt

LLMs generally perform better when prompted in English



(A) Qwen2.5-72B Performance

(C) Llama-3.3-70B Performance

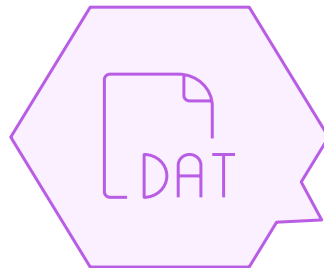English Prompt     Target Language Prompt

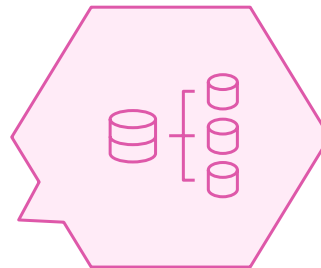# Takeaways from Additional Experiments

- LLMs still struggle with emotion recognition

- We observe large performance gaps across languages

- Performance still depends on prompt wording, number of shots and language

# BRIGHTER Public Release

**Annotations guidelines**

**Datasets**

**Individual labels**

https://brighter-dataset.github.io

# Thank you!

# Any questions?