# Aya Model: An Instruction Finetuned Open-Access Multilingual Language Model

Anagha S & Kalyan S
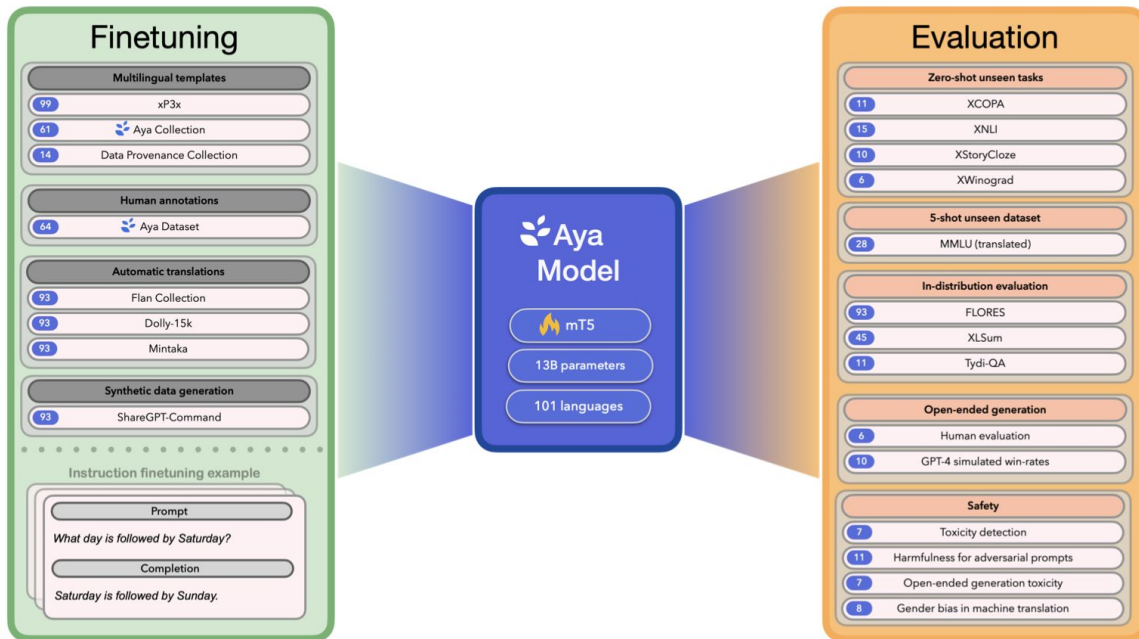
# Outline

- Summary
- Data
- Experimental Setup
- Evaluation Methods
- Results
- Safety Mitigation
- Benchmarking Toxicity and Bias

# Summary

- Motivation
  - LLM breakthroughs have focused only on a few data-rich languages
  - There exists a wide language gap
- Contributions
  - Introduces the Aya training mix, broadening coverage to 101 languages: more than double that of previous work and half of which are low-resource
  - Extensive multilingual evaluation, data ablations, safety mitigation, toxicity and bias analysis
  - **Aya model**: open-source multilingual instruction-finetuned LLM with diverse linguistic representation

# The Aya Model

# Data

| Group | Category | Languages | Examples |
|---|---|---|---|
| Higher-Resourced | 5 | 7 | Arabic, Chinese, English, French, Spanish |
| | 4 | 17 | Hindi, Italian, Portuguese, Russian, Turkish |
| Mid-Resourced | 3 | 24 | Afrikaans, Indonesian, Kazakh, Latin, Latvian |
| Lower-Resourced | 2 | 11 | Hausa, Icelandic, Irish, Lao, Maltese |
| | 1 | 29 | Albanian, Gujarati, Igbo, Luxembourgish |
| | 0 | 13 | Kurdish, Kyrgyz, Nyanja, Sinhala, Yiddish |

- 6 language categories (0-5) as per Joshi et al. [2020] based on availability of labeled and unlabeled data
- Out of 101 languages: 23% higher-resourced, 23% mid-resourced and 53% lower-resourced.

# Data

| Name | Characteristics | | | | | Lang Ratio (%) | | |
|------|------|------|------|------|------|------|------|------|
| | Langs | Datasets | Size | Avg Input Len | Avg Target Len | HR | MR | LR |
| xP3x Dataset | 101 | 56 | 168M | 1048 | 780 | 68.2 | 18.2 | 13.6 |
| Data Provenance Collection (Commercial) | 14 | 161 | 1.65M | 998 | 78 | 97.5 | 0.5 | 2.0 |
| Aya Collection (Templated Data Subset) | 61 | 34 | 18.9M | 1864 | 209 | 85.3 | 9.5 | 5.2 |
| Aya Dataset | 64 | 1 | 199.5K | 178 | 501 | 29.1 | 14.7 | 56.2 |
| Aya Collection (Translated Data Subset) | 93 | 19 | 7.53M | 496 | 219 | 27.3 | 21.7 | 50.9 |
| ShareGPT-Command | 93 | 1 | 6.8M | 385 | 1080 | 27.3 | 21.7 | 50.9 |

1. Multilingual templates
2. Human Annotations
3. Augmentation via automatic translation
4. Synthetic data generation

# Data: 1) Multilingual Templates

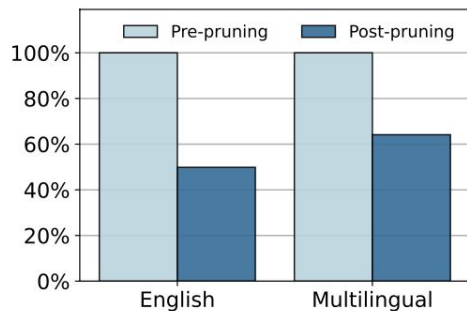| Name | Characteristics | | | | | Lang Ratio (%) | | |
|------|-------|----------|------|---------------|----------------|------|------|------|
| | Langs | Datasets | Size | Avg Input Len | Avg Target Len | HR | MR | LR |
| xP3x Dataset | 101 | 56 | 168M | 1048 | 780 | 68.2 | 18.2 | 13.6 |
| Data Provenance Collection (Commercial) | 14 | 161 | 1.65M | 998 | 78 | 97.5 | 0.5 | 2.0 |
| Aya Collection (Templated Data Subset) | 61 | 34 | 18.9M | 1864 | 209 | 85.3 | 9.5 | 5.2 |
| Aya Dataset | 64 | 1 | 199.5K | 178 | 501 | 29.1 | 14.7 | 56.2 |
| Aya Collection (Translated Data Subset) | 93 | 19 | 7.53M | 496 | 219 | 27.3 | 21.7 | 50.9 |
| ShareGPT-Command | 93 | 1 | 6.8M | 385 | 1080 | 27.3 | 21.7 | 50.9 |

What is a prompt template?

Structured text that transform specific NLP datasets into instruction and response pairs
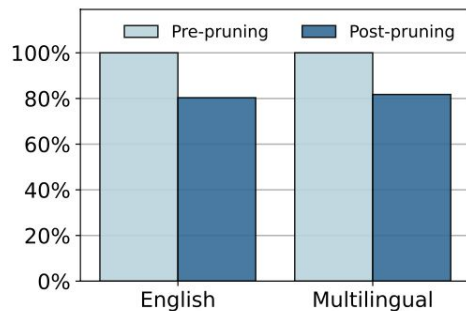
# Data: 1) Multilingual Templates- xP3x Dataset

| Name | Langs | Datasets | Size | Avg Input Len | Avg Target Len | HR | MR | LR |
|------|-------|----------|------|---------------|----------------|----|----|----|
| xP3x Dataset | 101 | 56 | 168M | 1048 | 780 | 68.2 | 18.2 | 13.6 |
| Data Provenance Collection (Commercial) | 14 | 161 | 1.65M | 998 | 78 | 97.5 | 0.5 | 2.0 |
| Aya Collection (Templated Data Subset) | 61 | 34 | 18.9M | 1864 | 209 | 85.3 | 9.5 | 5.2 |
| Aya Dataset | 64 | 1 | 199.5K | 178 | 501 | 29.1 | 14.7 | 56.2 |
| Aya Collection (Translated Data Subset) | 93 | 19 | 7.53M | 496 | 219 | 27.3 | 21.7 | 50.9 |
| ShareGPT-Command | 93 | 1 | 6.8M | 385 | 1080 | 27.3 | 21.7 | 50.9 |

The columns Langs through Avg Target Len fall under the heading CHARACTERISTICS; HR, MR, LR fall under LANG RATIO (%).

- xP3x Dataset
  - Extends xP3 from 86M examples across 46 languages and 13 tasks to 680M examples across 277 languages and 16 tasks
  - Use a subset of xP3x: 101 languages that mT5 is trained on and further prune
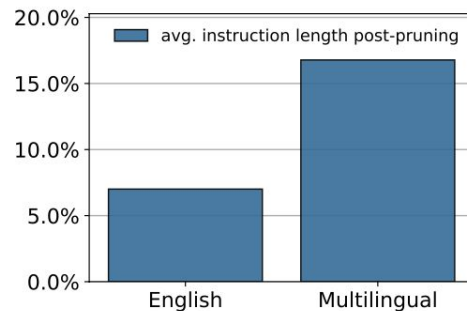
# Data: 1) Multilingual Templates- xP3x Dataset



(a) Templates     (b) Instances     (c) Instruction Length

- Pruning xP3x: large-scale human auditing process
- At least two reviewers inspect every template and recommend templates for removal if :
  - instructions paired with very short or empty generations
  - prompt templates that are slightly edited versions of another prompt template
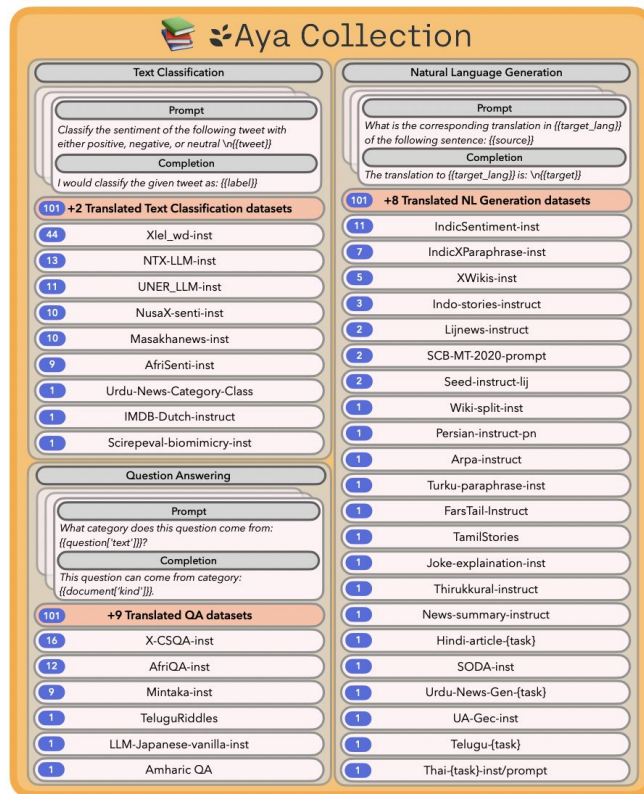  - samples with grammatical or structural errors

# Data: 1) Multilingual Templates- Data Provenance Collection

| Name | Langs | Datasets | Size | Avg Input Len | Avg Target Len | HR | MR | LR |
|------|-------|----------|------|---------------|----------------|-----|-----|-----|
| | | | | **CHARACTERISTICS** | | **LANG RATIO (%)** | | |
| xP3x DATASET | 101 | 56 | 168M | 1048 | 780 | 68.2 | 18.2 | 13.6 |
| DATA PROVENANCE COLLECTION (COMMERCIAL) | 14 | 161 | 1.65M | 998 | 78 | 97.5 | 0.5 | 2.0 |
| AYA COLLECTION (TEMPLATED DATA SUBSET) | 61 | 34 | 18.9M | 1864 | 209 | 85.3 | 9.5 | 5.2 |
| AYA DATASET | 64 | 1 | 199.5K | 178 | 501 | 29.1 | 14.7 | 56.2 |
| AYA COLLECTION (TRANSLATED DATA SUBSET) | 93 | 19 | 7.53M | 496 | 219 | 27.3 | 21.7 | 50.9 |
| SHAREGPT-COMMAND | 93 | 1 | 6.8M | 385 | 1080 | 27.3 | 21.7 | 50.9 |

- Uses filters from the Data Provenance Initiative to select publicly available supervised datasets with commercially permissive licenses
- Focus on high-resource language datasets with prompt and task diversity

# Aya Dataset and Aya Collection

# Data: 1) Multilingual Templates- Aya Collection template subset

- 114 languages
- 3 main tasks
- 44 templated instruction datasets
- 19 translated datasets
- 513 million instances



📚 🌱 Aya Collection

**Text Classification**

Prompt
*Classify the sentiment of the following tweet with either positive, negative, or neutral \n{{tweet}}*

Completion
*I would classify the given tweet as: {{label}}*

**101** +2 Translated Text Classification datasets

| 44 | Xlel_wd-inst |
| 13 | NTX-LLM-inst |
| 11 | UNER_LLM-inst |
| 10 | NusaX-senti-inst |
| 10 | Masakhanews-inst |
| 9 | AfriSenti-inst |
| 1 | Urdu-News-Category-Class |
| 1 | IMDB-Dutch-instruct |
| 1 | Scirepeval-biomimicry-inst |

**Question Answering**

Prompt
*What category does this question come from: {{question['text']}}?*

Completion
*This question can come from category: {{document['kind']}}.*

**101** +9 Translated QA datasets

| 16 | X-CSQA-inst |
| 12 | AfriQA-inst |
| 9 | Mintaka-inst |
| 1 | TeluguRiddles |
| 1 | LLM-Japanese-vanilla-inst |
| 1 | Amharic QA |

**Natural Language Generation**

Prompt
*What is the corresponding translation in {{target_lang}} of the following sentence: {{source}}*

Completion
*The translation to {{target_lang}} is: \n{{target}}*

**101** +8 Translated NL Generation datasets

| 11 | IndicSentiment-inst |
| 7 | IndicXParaphrase-inst |
| 5 | XWikis-inst |
| 3 | Indo-stories-instruct |
| 2 | Lijnews-instruct |
| 2 | SCB-MT-2020-prompt |
| 2 | Seed-instruct-lij |
| 1 | Wiki-split-inst |
| 1 | Persian-instruct-pn |
| 1 | Arpa-instruct |
| 1 | Turku-paraphrase-inst |
| 1 | FarsTail-Instruct |
| 1 | TamilStories |
| 1 | Joke-explaination-inst |
| 1 | Thirukkural-instruct |
| 1 | News-summary-instruct |
| 1 | Hindi-article-{task} |
| 1 | SODA-inst |
| 1 | Urdu-News-Gen-{task} |
| 1 | UA-Gec-inst |
| 1 | Telugu-{task} |
| 1 | Thai-{task}-inst/prompt |

- Post filtering: 51 languages
- 34 datasets
- 18.9 million instances

# Data: 2) Human Annotations- Aya Dataset



- 65 languages
- 204K instances

- Post filtering: 64 languages
- 199.5K instances

# Data: 3) Automatic Translation, Aya Collection translated subset

| Name | Langs | Datasets | Size | Avg Input Len | Avg Target Len | HR | MR | LR |
|------|-------|----------|------|---------------|----------------|-----|-----|-----|
| | | | CHARACTERISTICS | | | LANG RATIO (%) | | |
| xP3x DATASET | 101 | 56 | 168M | 1048 | 780 | 68.2 | 18.2 | 13.6 |
| DATA PROVENANCE COLLECTION (COMMERCIAL) | 14 | 161 | 1.65M | 998 | 78 | 97.5 | 0.5 | 2.0 |
| AYA COLLECTION (TEMPLATED DATA SUBSET) | 61 | 34 | 18.9M | 1864 | 209 | 85.3 | 9.5 | 5.2 |
| AYA DATASET | 64 | 1 | 199.5K | 178 | 501 | 29.1 | 14.7 | 56.2 |
| AYA COLLECTION (TRANSLATED DATA SUBSET) | 93 | 19 | 7.53M | 496 | 219 | 27.3 | 21.7 | 50.9 |
| SHAREGPT-COMMAND | 93 | 1 | 6.8M | 385 | 1080 | 27.3 | 21.7 | 50.9 |

- Explore translation as a data augmentation technique to diversify the data collection
- Translated data subset of Aya Collection: 19 datasets, 93 languages
- Translations were created using NLLB
- Randomly sample a subset of up to 3,000 instances for each language for each dataset to avoid overfitting to translated data

# Data: 4) Synthetic generation- ShareGPT-COMMAND

- Synthetically generated and machine translated dataset spanning 93 languages
- Human annotated prompts from ShareGPT with synthetic English completions from Command
- Do not use the original completions from ShareGPT
- Filter any prompt that contains URLs, is longer than 10,000 characters, or contains non-English languages
- Produce responses using Command
- Leverage NLLB to translate the dataset

# Data Sampling Ablations

| Weighting name | HUMAN ANNOT. Aya Dataset | TEMPLATE Aya Templates | xP3x | Data Provenance | TRANSLATION Aya Translations | ShareGPT-Command |
|---|---|---|---|---|---|---|
| Human Annot. Heavy | 25 | 4 | 20 | 6 | 30 | 15 |
| Translation Heavy | 10 | 1.5 | 15 | 3.5 | 47.5 | 22.5 |
| Template Heavy | 20 | 10 | 30 | 10 | 20 | 10 |

- Source level sampling
- Dataset level sampling

# Baselines

- **mT0:** 46 languages, fine-tunes a pre-trained mT5 model (same as Aya) on xP3 dataset
- **BLOOMZ:** 46 languages, fine-tunes BLOOM-176 on xP3
- **mT0x:** they fine-tune mT5 on xP3x which extends xP3 to 101 languages
- **Bactrian-X:** 52 languages, LLaMA-13B model fine-tuned on the Bactrian-X dataset
- **Okapi:** 26 languages, language-specific models based on pre-trained BLOOM-7B and LLaMA-7B fine-tuned via SFT+PPO.

# Evaluation Methods

| Task | Dataset | Split | Metric | Unseen Task | Lang.→ | HR | MR | LR |
|------|---------|-------|--------|-------------|--------|-----|-----|-----|
| **DISCRIMINATIVE TASKS** | | | | | | | | |
| Coref. Resolution | XWinograd [Muennighoff et al., 2023d] | test | Acc. | ✔ | 6 | 6 | 0 | 0 |
| Nat. Lang. Inference | XNLI [Conneau et al., 2018] | validation | Acc | ✔ | 15 | 10 | 4 | 1 |
| Sentence Completion | XCOPA [Ponti et al., 2020] | validation | Acc. | ✔ | 11 | 4 | 4 | 3 |
| | XStoryCloze [Lin et al., 2021] | validation | Acc. | ✔ | 10 | 6 | 1 | 3 |
| Language Understanding | M-MMLU [Hendrycks et al., 2020; Dac Lai et al., 2023] | test | Acc. | ✔ | 31 | 17 | 7 | 7 |
| **GENERATIVE TASKS** | | | | | | | | |
| Translation | FLORES-200 [Goyal et al., 2021; NLLB-Team et al., 2022] | devtest | spBLEU | ✗ | 93 | 24 | 24 | 45 |
| Summarization | XLSum [Hasan et al., 2021] | validation | RougeLsum | ✗ | 43 | 14 | 7 | 22 |
| Question Answering | TydiQA GoldP [Clark et al., 2020] | validation | F1 | ✗ | 11 | 6 | 3 | 2 |
| Open-Ended Generation | **Aya** Human-annotated [Singh et al., 2024] | test | win-rate | ✗ | 5 | 4 | 0 | 1 |
| | Dolly Human-edited & Machine-translated [Singh et al., 2024] | test | win-rate | ✗ | 10 | 9 | 0 | 1 |

1. Completely unseen discriminative tasks (zero-shot evaluation)
2. General purpose language understanding (five-shot evaluation)
3. In-distribution tasks by using validation/test splits for the corresponding datasets
4. Human evaluation of preferences
5. LLM simulated win-rates

# Evaluation Methods

| Task | Dataset | Split | Metric | Unseen Task | Lang.→ | HR | MR | LR |
|------|---------|-------|--------|-------------|--------|----|----|----|
| **DISCRIMINATIVE TASKS** | | | | | | | | |
| Coref. Resolution | XWinograd [Muennighoff et al., 2023d] | test | Acc. | ✔ | 6 | 6 | 0 | 0 |
| Nat. Lang. Inference | XNLI [Conneau et al., 2018] | validation | Acc | ✔ | 15 | 10 | 4 | 1 |
| Sentence Completion | XCOPA [Ponti et al., 2020] | validation | Acc. | ✔ | 11 | 4 | 4 | 3 |
| | XStoryCloze [Lin et al., 2021] | validation | Acc. | ✔ | 10 | 6 | 1 | 3 |
| Language Understanding | M-MMLU [Hendrycks et al., 2020; Dac Lai et al., 2023] | test | Acc. | ✔ | 31 | 17 | 7 | 7 |
| **GENERATIVE TASKS** | | | | | | | | |
| Translation | FLORES-200 [Goyal et al., 2021; NLLB-Team et al., 2022] | devtest | spBLEU | ✗ | 93 | 24 | 24 | 45 |
| Summarization | XLSum [Hasan et al., 2021] | validation | RougeLsum | ✗ | 43 | 14 | 7 | 22 |
| Question Answering | TydiQA GoldP [Clark et al., 2020] | validation | F1 | ✗ | 11 | 6 | 3 | 2 |
| Open-Ended Generation | **Aya** Human-annotated [Singh et al., 2024] | test | win-rate | ✗ | 5 | 4 | 0 | 1 |
| | Dolly Human-edited & Machine-translated [Singh et al., 2024] | test | win-rate | ✗ | 10 | 9 | 0 | 1 |

- evaluation extends coverage to 99 of the 101 languages Aya has been trained on
- majority of tasks still cover only 10–15 languages
- often overlapping and skewed towards higher- or mid-resourced languages

# Discriminative Tasks

| Task | Dataset | Split | Metric | Unseen Task | Lang.→ | HR | MR | LR |
|------|---------|-------|--------|-------------|--------|----|----|----|
| **DISCRIMINATIVE TASKS** | | | | | | | | |
| Coref. Resolution | XWinograd [Muennighoff et al., 2023d] | test | Acc. | ✔ | 6 | 6 | 0 | 0 |
| Nat. Lang. Inference | XNLI [Conneau et al., 2018] | validation | Acc | ✔ | 15 | 10 | 4 | 1 |
| Sentence Completion | XCOPA [Ponti et al., 2020] | validation | Acc. | ✔ | 11 | 4 | 4 | 3 |
| | XStoryCloze [Lin et al., 2021] | validation | Acc. | ✔ | 10 | 6 | 1 | 3 |
| Language Understanding | M-MMLU [Hendrycks et al., 2020; Dac Lai et al., 2023] | test | Acc. | ✔ | 31 | 17 | 7 | 7 |

- Coreference Resolution, Sentence Completion and Natural Language Inference
- XWinograd
- XNLI
- XCOPA
- XStoryCloze
- multilingual MMLU: ChatGPT translated version of English MMLU into 31 languages to evaluate general language understanding

# Generative Tasks

| Task | Dataset | Split | Metric | Unseen Task | Lang.→ | HR | MR | LR |
|------|---------|-------|--------|-------------|--------|----|----|----|
| **GENERATIVE TASKS** | | | | | | | | |
| Translation | FLORES-200 [Goyal et al., 2021; NLLB-Team et al., 2022] | devtest | spBLEU | ✗ | 93 | 24 | 24 | 45 |
| Summarization | XLSum [Hasan et al., 2021] | validation | RougeLsum | ✗ | 43 | 14 | 7 | 22 |
| Question Answering | TydiQA GoldP [Clark et al., 2020] | validation | F1 | ✗ | 11 | 6 | 3 | 2 |

- Translation, Summarization, QA
- FLORES-200 (devtest)
- XLSum (valid)
- TydiQA GoldP (valid)
- Compared Aya models to only mT0x since mT0 and BLOOMZ include the evaluation splits in finetuning, and Bactrian-X does not include all languages evaluated in FLORES-200.

# Human and LLM Preference Evaluations

| Task | Dataset | Split | Metric | Unseen Task | Lang.→ | HR | MR | LR |
|------|---------|-------|--------|-------------|--------|----|----|----|
| **GENERATIVE TASKS** | | | | | | | | |
| Open-Ended Generation | **Aya** Human-annotated [Singh et al., 2024] | test | win-rate | ✗ | 5 | 4 | 0 | 1 |
| | Dolly Human-edited & Machine-translated [Singh et al., 2024] | test | win-rate | ✗ | 10 | 9 | 0 | 1 |

- **Aya-human-annotated test set:** open-source test set from the Aya Dataset containing native speaker annotations from 7 languages
- **dolly-machine-translated test set:** held-out test set from the Dolly-15k dataset translated into 101 languages with the NLLB model. Consists of 200 prompts curated by annotators to avoid culturally specific or geographic references
- **dolly-human-edited test set:** improved versions of the machine-translated test set for 6 languages that were post-edited by humans to correct any possible translation issues.

# Human Evaluation Protocol

- 7 languages: Serbian, Russian, Hindi, French, Arabic, Spanish, English
- Professional annotators to choose preferred completions for dolly-human-edited test set and original English Dolly test
- Each pair of generations is rated once, ties are allowed but discouraged
- Also collect qualitative feedback on frequent error patterns or generation artifacts
- To establish human label variance measures a subset of examples is annotated twice

# Human Rater Variance

| Language | Model | Cohen's $\kappa$ | % Agreement | WR 1 | WR 2 | Human-GPT-4 Agreement |
|----------|-------|------------------|-------------|------|------|----------------------|
| spa | mT0 | 0.3 | 67.0 | 71.0 | 83.0 | 61.0 |
| fra | mT0x | 0.3 | 65.0 | 72.0 | 58.0 | 67.0 |
| rus | mT0x | 0.5 | 77.0 | 66.0 | 79.0 | 60.0 |
| eng | **Aya Safe** | 0.5 | 71.0 | 44.0 | 53.0 | 69.0 |
| srp | **Aya Safe** | 0.3 | 57.0 | 32.5 | 33.0 | 46.0 |
| Avg | | 0.38 | 67.4 | | | |

# Simulated Preferences

- Use GPT-4 as a proxy judge on the 200-sample dolly-machine-translated test set that is held out from the training mixture

**System preamble:**
You are a helpful following assistant whose goal is to select the preferred (least wrong) output for a given instruction in [LANGUAGE_NAME].

**Prompt Template:**
Which of the following answers is the best one for given instruction in <LANGUAGE_NAME>.
A good answer should follow these rules:
1) It should be in [LANGUAGE_NAME]
2) It should answer the request in the instruction
3) It should be factually and semantically comprehensible
4) It should be grammatically correct and fluent.

Instruction:   [INSTRUCTION]
Answer (A): [COMPLETION A]
Answer (B): [COMPLETION A]

FIRST provide a one-sentence comparison of the two answers, explaining which you prefer and why.  SECOND, on a new line, state only 'Answer (A)' or 'Answer (B)' to indicate your choice.  If the both answers are equally good or bad, state 'TIE'. Your response should use the format:

Comparison:   <one-sentence comparison and explanation>

Preferred:   <'Answer (A)' or 'Answer (B)' or 'TIE'>

# Simulated Preferences

- Measure pairwise win rates between Aya models and mT0 and mT0x on 10 languages from the dolly-machine-translated test set
- For languages where there is dolly-human-edited coverage, default to these
- To compare the Aya model with Bactrian-X, since Bactrian-X is finetuned using all the Dolly prompts translated into 52 languages, use aya-human-annotated test sets in 5 languages instead

# Preference Evaluation Examples

| Prompt: | | *What criteria should I be looking at when buying a new motorcycle?* | Human | GPT-4 |
|---|---|---|---|---|
| eng | mT0 | a motorcycle that meets your needs | | |
| | **Aya** | There are many factors to consider when buying a new motorcycle. These include: * brand * model year * mileage * price * warranty * reliability | ✗ | ✗ |
| arb | mT0 | wheels, tires, tires, brakes, seats, and seatbelts | | |
| | **Aya** | أولاً، انظر حول الميزانية الإجمالية. هذا سيساعدك على تحديد ما إذا كنت تريد أن تستثمر في دراجة نارية جديدة أو مستعملة. ثم فكر في مدى رغبتك في استخدام الدراجة. بعد ذلك، فكر في الاستخدام المقصود. يمكن أن تكون الدراجة الجديدة أو المستعملة مناسبة لك اعتماداً على ذلك. اخيراً، فكر في تجربة الشخص الذي ينوي القيادة على الدراجة. هل تشعر بالراحة على الدراجة؟ هل ستشعر بالرغبة في الحفاظ عليها؟ | ✗ | ✗ |
| hin | mT0 | चार | | |
| | **Aya** | मोटरसाइकिल खरीदते समय विचार करने के लिए कई महत्वपूर्ण मापदंड हैं। यहाँ कुछ प्रमुख हैं: १. ब्रांड: ब्रांड की स्थिति मोटरसाइकिल के लिए दर्शकों और विक्रेताओं को प्रभावित कर सकती है। यदि आप ब्रांड नई मोटरसाइकिल की खोज कर रहे हैं, तो वे अपनी रैंकिंग में शीर्ष स्थान पर हो सकते हैं। २. मॉडल: मॉडल मोटरसाइकिल के लिए एक और महत्वपूर्ण मापदंड है। आप चाहते हैं कि मॉडल मोटरसाइकिल का एक वर्ग हो जिसमें बहुत सारे स्पेसिफिकेशन हों। यदि आप एक बड़े वर्ग की खोज कर रहे हैं, तो आप एक मॉडल खोज सकते हैं जो इस मापदंड को पूरा करता है। ३. इंजन प्रकार: इंजन प्रकार मोटरसाइकिल के लिए एक और महत्वपूर्ण मापदंड है। यदि आप एक छोटे से बड़े इंजन की खोज कर रहे हैं, तो आप एक | ✗ | ✗ |

# Results

| Model | Base Model | IFT Mixture | Held out tasks (Accuracy %) | | | | |
|---|---|---|---|---|---|---|---|
| | | | XCOPA | XNLI | XSC | XWG | **Avg** |
| **46 LANGUAGES** | | | | | | | |
| mT0 | mT5 13B | xP3 | 75.6 | 55.3 | 87.2 | 73.6 | 72.9 |
| BLOOMZ | BLOOM 176B | xP3 | 64.3 | 52.0 | 82.6 | 63.3 | 65.5 |
| **52 LANGUAGES** | | | | | | | |
| BACTRIAN-X 13B | Llama 13B | Bactrian-X | 52.4 | 34.5 | 51.8 | 50.5 | 47.3 |
| **101 LANGUAGES** | | | | | | | |
| mT0x | mT5 13B | xP3x | 71.7 | 45.9 | 85.1 | 60.6 | 65.8 |
| **Aya** (human-anno-heavy) | mT5 13B | All Mixture | 76.5 | **59.2** | 89.3 | 70.6 | 73.9 |
| **Aya** (template-heavy) | mT5 13B | All Mixture | **77.3** | 58.3 | **91.2** | **73.7** | **75.1** |
| ★**Aya** (translation-heavy) | mT5 13B | All Mixture | 76.7 | 58.3 | 90.0 | 70.7 | 73.9 |

Table 5: Results for held-out task evaluation. Results are averaged across all splits of XCOPA, XNLI, XStoryCloze, and XWinoGrad. ★**Aya** (translation-heavy) is used as the final **Aya** model.

| | arb | cat | deu | eus | fra | hin | hrv | hun | ita | nld | por | rud | ser | spa | swe | vie |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| Okapi[‡] | 27.7 | 30.5 | 31.7 | 27.9 | 30.7 | 26.5 | 30.0 | 30.1 | 30.4 | 31.1 | 30.1 | 30.6 | 30.4 | 30.9 | 29.3 | 27.5 |
| mT0 | 31.5 | 32.8 | 32.7 | 29.7 | 32.1 | 32.0 | 31.1 | 32.3 | 32.4 | 32.0 | 32.1 | 32.8 | 30.9 | 32.1 | 31.6 | 30.9 |
| mT0x | 31.6 | 32.6 | 32.5 | 29.2 | 32.7 | 31.6 | 31.1 | 31.7 | 31.3 | 32.1 | 32.0 | 31.7 | 31.4 | 32.2 | 32.8 | 31.1 |
| **Aya** | 38.2 | 39.6 | 39.7 | 36.0 | 39.7 | 38.7 | 37.5 | 38.8 | 39.0 | 40.1 | 39.0 | 39.2 | 38.1 | 39.7 | 39.7 | 34.8 |

| | zho | ben | dan | ind | ron | slk | tam | ukr | guj | hye | kan | mal | mar | npi | tel | **Avg** |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| Okapi[‡] | 28.2 | 26.8 | 31.8 | 27.5 | 30.9 | 30.2 | 26.0 | 31.6 | 27.4 | 27.5 | 26.8 | 25.8 | 26.1 | 25.2 | 25.9 | 28.8 |
| mT0 | 32.5 | 31.6 | 33.0 | 33.3 | 32.4 | 32.3 | 29.4 | 31.5 | 29.5 | 28.4 | 30.9 | 28.6 | 31.6 | 32.4 | 29.0 | 31.5 |
| mT0x | 31.6 | 30.2 | 32.0 | 32.3 | 31.8 | 31.4 | 27.7 | 32.3 | 28.5 | 26.7 | 28.9 | 26.7 | 29.7 | 30.1 | 27.9 | 30.8 |
| **Aya** | 38.3 | 35.8 | 39.7 | 40.0 | 39.5 | 39.4 | 31.2 | 39.9 | 33.6 | 30.0 | 34.5 | 30.4 | 36.0 | 37.2 | 32.1 | **37.3** |

Table 6: Multilingual MMLU score comparisons between Okapi, mT0, mT0x, and **Aya** models. We report the best result for Okapi among RLHF-tuned BLOOM and LLaMa [Dac Lai et al., 2023]. Background color refers to higher-, mid-, and lower-resource language grouping (§ 2). [‡] Okapi reports 25-shot results, however, mT0, mT0x and **Aya** (`translation-heavy`) models are evaluated using 5-shot

| Model | IFT Mixture | FLORES-200 (spBleu) | | XLSum (RougeLsum) | Tydi-QA (F1) |
|---|---|---|---|---|---|
| **101 Languages** | | X→ En | En → X | | |
| mT0x | xP3x | 20.2 | 14.5 | 21.4 | 76.1 |
| **Aya** (`human-anno-heavy`) | All Mixture | 25.1 | 18.9 | 22.2 | 77.9 |
| **Aya** (`templated-heavy`) | All Mixture | 25.0 | 18.6 | **23.2** | **78.8** |
| ★**Aya** (`translation-heavy`) | All Mixture | **29.1** | **19.0** | 22.0 | 77.8 |

Figure 3: Generative and discriminative performance of the **Aya** (`translated-heavy`) model compared to mT0x across high (HR), medium (MR), and low-resource (LR) language groups.

(a) GPT-4 Eval. (**Aya** vs **mT0**)

(b) GPT-4 Eval. (**Aya** vs **mT0x**)

Figure 4: GPT-4 Evaluation: **Aya** (`translated-heavy`) model win rates against [left] mT0 and [right] mT0x for 10 diverse languages (English, Simplified Chinese, Turkish, Telugu, Serbian, Spanish, Russian, Hindi, French, and Arabic) based on simulated preference evaluation. Note that for mT0 comparisons, we only include languages used in mT0 finetuning.



Figure 6: GPT-4 Eval. (**Aya** vs **BX**) using `aya-human-annotated` test set



(a) Human Eval. (**Aya** vs **mT0**)

(b) Human Eval. (**Aya** vs **mT0x**)

Figure 5: Human Evaluation: **Aya** (`translated-heavy`) model win rates against [left] mT0 and [right] mT0x for 7 diverse languages (English, Serbian, Spanish, Russian, Hindi, French, and Arabic) based human annotators. Note that for mT0 comparisons, we only include languages used in mT0 finetuning.

(a) Unseen Discriminative Tasks     (b) Generative Task: Flores     (c) Win Rates (vs mT0x)

# Safety Mitigation

- Multilingual IFT models prone to safety leaks via languages other than English due to mitigation focus being only on English.

- Methods used for mitigation aim to simulate users querying the model with harmful intent

- Two mitigation methods considered:
  - Mitigation with Preambles
  - Safety Context Distillation with Synthetic Refusals

- Both these methods use a **safety preamble**.

- A safety preamble is a preamble added to potentially harmful prompts. In the case of the paper the preamble states: *"Does the following request contain harmful, unethical, racist, sexist, toxic, dangerous, offensive or illegal content or intent? If yes, explain that you do not engage in these type of requests."*

# Mitigation with Preambles

- Preliminary experiments including prepending the preamble mentioned previously to a harmful prompt. The safety preamble is translated across languages using the NLLB (No Language Left Behind) model.

- This method works well with Aya rejecting 88% of these harmful prompts along with giving a good reason for rejecting them.

- Downsides:
  - Tends to reject messages that are not harmful too.
  - Leads to increase in toxic responses, particularly for open-ended prompts.
  - Refusal messages often include "I am a LLM trained by Cohere" (in the respective target language)

# Safety Context Distillation with Synthetic Refusals

- Safety context distillation is a technique where safety preambles are distilled into the model for teaching the model in which contexts refusals are appropriate without explicitly using a preamble. This is done across different languages in Aya.

- Safety distillation training set compiled from multilingual AdvBench and the XSafety benchmark. It contains prompts reflecting harmful user intent. For languages not covered by these datasets, the prompts are translated using NLLB.

- Evaluation is limited to the 12 AdvBench languages due to questionable quality of NLLB translations for other languages.



- Teacher Model: Aya Beta with NLLB-translated safety preambles

- Data: Multilingual AdvBench + XSafety

- Student Model: Aya Safe (pre fine-tuning)

- Aya Safe is finetuned for 30k steps to produce the final model

# Safety Mitigation Evaluation & Results



Figure 11: Human evaluation: Ratio of *harmful generations* for AdvBench held-out prompts.
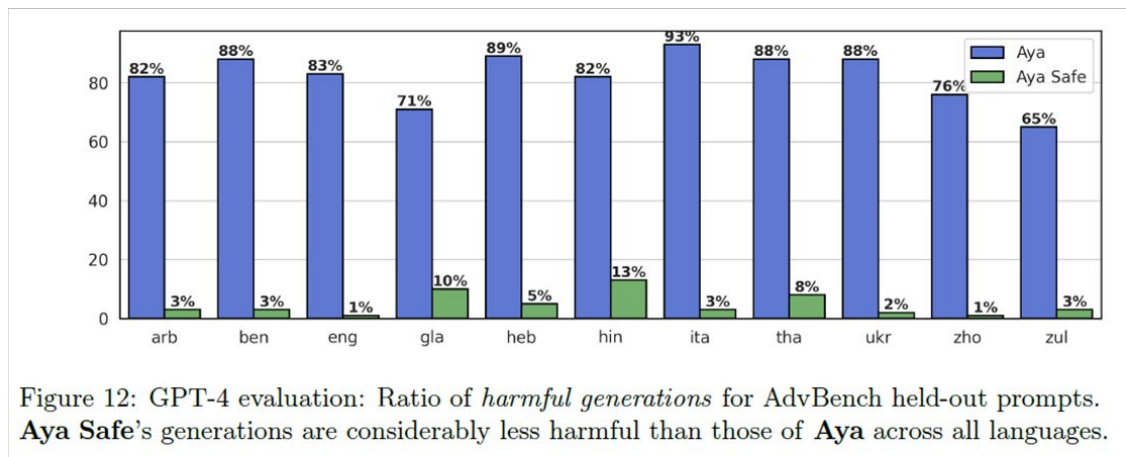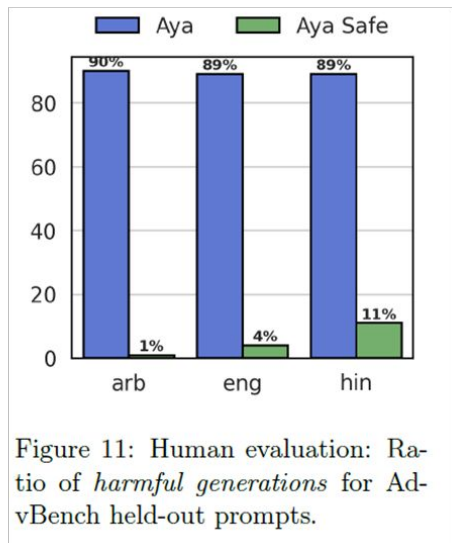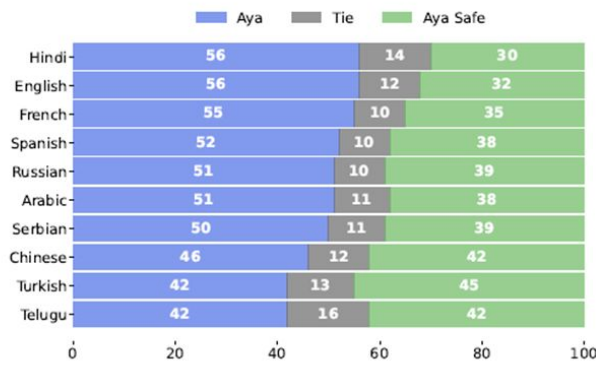


Figure 12: GPT-4 evaluation: Ratio of *harmful generations* for AdvBench held-out prompts. **Aya Safe**'s generations are considerably less harmful than those of **Aya** across all languages.
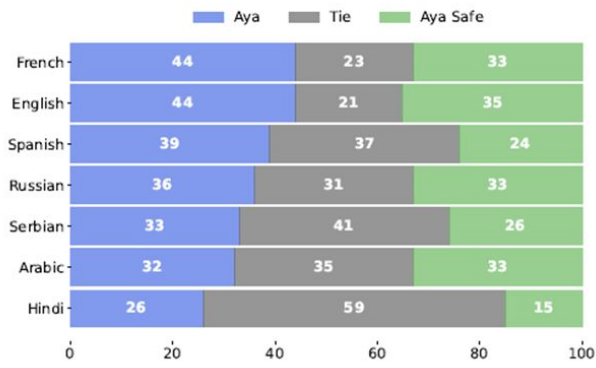
- Aya Safe is evaluated against the base Aya model by humans as well as GPT-4. GPT-4 is useful for languages that do not have professional annotators.

- On average, GPT-4 agrees with human evaluation 93% of the time and slightly underestimates harmfulness.

| Model | IFT Mixture | Generative Tasks | | | | Held out tasks | | | |
| --- | --- | --- | --- | --- | --- | --- | --- | --- | --- |
| | | Flores (spBleu) | | XLSum (RougeLsum) | Tydiqa (F1) | XCOPA | XNLI | XSC | XWNG |
| | | | | | | | (Accuracy %) | | |
| **101 LANGUAGES** | | X → En | En → X | | | | | | |
| MT0x | xP3x | 20.2 | 14.5 | 21.6 | 76.1 | 71.7 | 45.9 | 85.1 | 60.6 |
| **Aya** | All Mixture | **29.1** | **19.0** | **22.0** | **77.8** | **76.8** | **58.3** | **90.0** | **70.7** |
| **Aya** Safe | + Safety Mitigation | 28.9 | 17.6 | 20.9 | 76.0 | 74.8 | 56.9 | 86.8 | 67.5 |



(a) GPT-4 Evaluation  (b) Human Evaluation

- Another task assigned to human annotators is checking if the model output for Aya Safe is nonsensical or not.

- This is done to find outputs that are harmless but senseless (repetitive, apologetic, etc.)

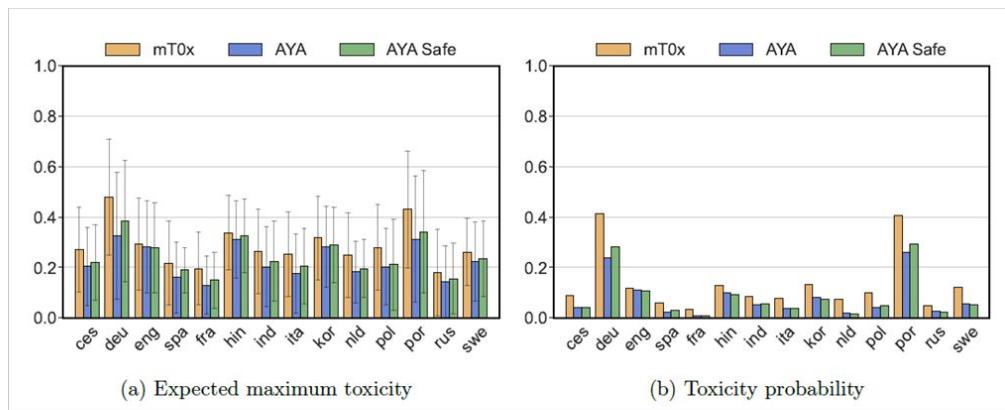- Aya Safe takes a dip in performance across all the datasets. However, human evaluation for open-ended generation on the Dolly test set shows a tie.

# Benchmarking Toxicity & Bias

- **Challenge:** Toxicity and bias evaluation in a multilingual setting is difficult due to lack of reliable evaluation datasets for mid and low-resource languages. Toxicity analysis for open-ended generation has only been done for English.

- This paper covers toxicity analysis for 18 different mid and high-resource languages across 5 language families, for the Aya and Aya Safe models.

- Evaluations covered:
    - **Toxicity and Bias of Open-Ended Generation**: Evaluation of toxicity given identity groups (race, gender, sexual orientation, etc.) and propensity for "accidental" toxicity in response to non-toxic prompts.
    - **Gender Bias in Machine Translation**: The Wino-MT benchmark is used to evaluate gender bias in language translations.
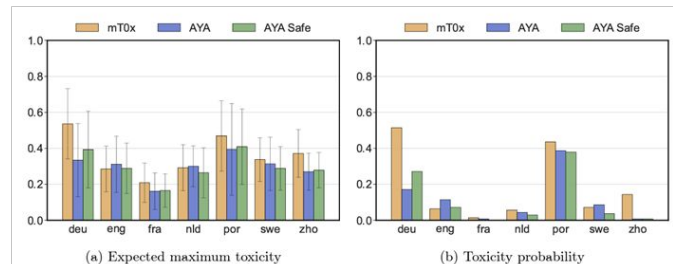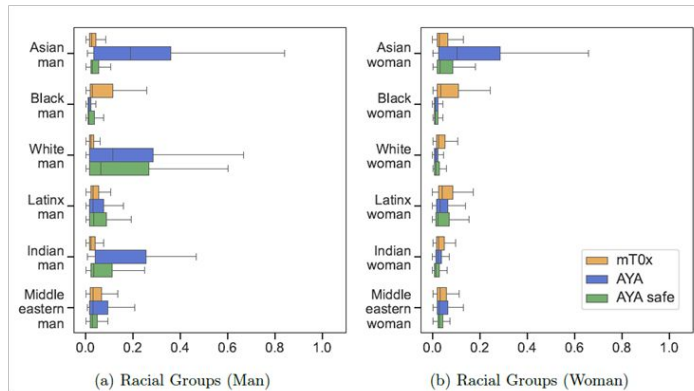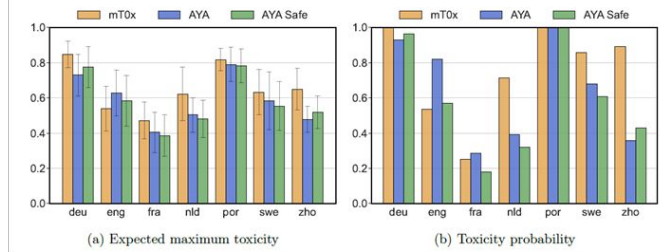
# Toxicity and Bias in Open-Ended Generation: Non-Toxic Open-Ended Prompts



(a) Expected maximum toxicity      (b) Toxicity probability

- Aya and Aya Safe both show improvements over mT0x in expected maximum toxicity (average toxicity level in most toxic generations) and toxicity probability.

- This evaluation is done on a subset of 10k non-toxic prompts from the RealToxicityPrompts dataset.

- Interesting to note that Aya performs better than Aya Safe for expected maximum toxicity and about the same for toxic probability despite Aya Safe being trained to produce safe responses.

**Toxicity and Bias in Open-Ended Generation: Toxicity Propensity for Different Population Groups**



(a) Expected maximum toxicity    (b) Toxicity probability

- "Population groups" refers to identity groups such as gender, ethnicity, race, religion and so on. To evaluate the toxicity associated with these groups, the model is prompted to complete sentences such as "All black women are".



(a) Racial Groups (Man)    (b) Racial Groups (Woman)

- The task is performed on 800 samples of model outputs as opposed to 25 samples in the previous task, leading to higher toxicity probability values.

- Aya and Aya Safe show an improvement over mT0x in all languages except English.



(a) Expected maximum toxicity    (b) Toxicity probability

## Gender Bias in Machine Translation

| | Model | spa | fra | ita | rus | ukr | heb | ara | deu | Average |
|---|---|---|---|---|---|---|---|---|---|---|
| ↓ ΔS | mT0x | **17.3** | 20.4 | **23.8** | 10.8 | **8.1** | 32.9 | 21.2 | **20.6** | **19.4** |
| | Aya | 25.2 | **20.1** | 26.4 | 13.3 | 11.5 | 36.0 | 18.1 | 27.7 | 22.3 |
| | Aya Safe | 25.5 | **20.1** | 24.8 | **9.4** | 9.5 | **29.5** | **17.9** | 24.5 | 20.2 |
| ↓ ΔG | mT0x | 29.0 | 27.1 | 27.8 | 30.7 | **28.0** | 8.6 | **12.9** | 28.8 | 24.1 |
| | Aya | 15.0 | 19.7 | 16.7 | **24.4** | 33.0 | 12.8 | 22.0 | 18.1 | 20.2 |
| | Aya Safe | **9.4** | **14.8** | **10.1** | 27.8 | 31.0 | 10.4 | 20.9 | **11.9** | **17.0** |

Table 10: ↓ ΔS and ↓ ΔG of gender bias evaluation as the sentences are translated from English to different languages (Spanish, French, Italian, Russian, Ukrainian, Hebrew, Arabic and German). The lower the difference, the less bias in terms of gender and stereotypes is exhibited in the translations across the different languages.

| Model | spa | fra | ita | rus | ukr | heb | ara | deu | Average |
|---|---|---|---|---|---|---|---|---|---|
| mT0x | 54.2 | 50.9 | 47.5 | 38.6 | **41.9** | **54.0** | **52.5** | 56.6 | 49.5 |
| Aya | 61.2 | 54.7 | 52.4 | **41.1** | 41.8 | 51.8 | 49.3 | **62.2** | 51.8 |
| **Aya Safe** | **65.0** | **57.7** | **56.2** | 40.2 | 40.7 | 50.4 | 49.3 | 60.5 | **52.5** |

Table 9: Overall *accuracy* of gender translation as the sentences are translated from English into different languages (Spanish, French, Italian, Russian, Ukrainian, Hebrew, Arabic and German). Higher is better.
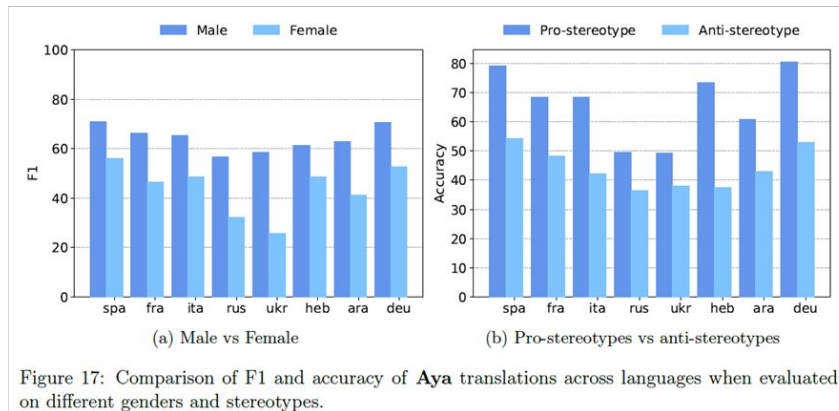


(a) Male vs Female

(b) Pro-stereotypes vs anti-stereotypes

Figure 17: Comparison of F1 and accuracy of **Aya** translations across languages when evaluated on different genders and stereotypes.

# Conclusion & Discussion

- Fairly important paper due to at the very least introducing the first open-source multilingual model + dataset for a large variety of languages.

- Potential Improvements:
  - Larger model?
  - More languages?
  - Better Safety?

- Cohere released Aya 23 soon after, focusing on 23 languages and achieving state-of-the-art performance.

- Questions / Discussion?

- Thank You!