



Accelerating multilingual
AI through open science

The Aya Movement at a glance.

cohere.com/research/aya



The word Aya is derived from the Twi language meaning “fern” - a symbol of endurance and resourcefulness. Aya embodies our dedication to advancing multilingual AI.

3 

Models

513M 

Total Release Dataset
Size

3K 

Independent
Researchers

250+ 

Language Ambassadors

119 

Countries

204K 

Original Human
Annotations

101 

Languages

81K 

Discord Messages

The Aya models and datasets cover 101 languages with enhanced performance for 23 of them

Achinese · Afrikaans · Albanian · Amharic · **Arabic** · Arabic · Armenian · Azerbaijani · Balinese · Banjar · Basque · Belarusian · Bemba · Bengali · Bulgarian · Burmese · Catalan · Cebuano · **Chinese** · Croatian · **Czech** · Danish · **Dutch** · **English** · Esperanto · Estonian · Filipino · Finnish · Fon · **French** · Galician · Georgian · **German** · **Greek** · Gujarati · Haitian Creole · Hausa · **Hebrew** · **Hindi** · Hungarian · Icelandic · Igbo · **Indonesian** · Irish · **Italian** · **Japanese** · Javanese · Kannada · Kanuri · Kashmiri · Kazakh · Khmer · Kinyarwanda · **Korean** · Kurdish · Kyrgyz · Lao · Latvian · Ligurian · Lithuanian · Luxembourgish · Macedonian · Madurese · Malagasy · Malay · Malayalam · Maltese · Manipuri · Maori · Marathi · Minangkabau · Mongolian · Nepali · Ngaju · Northern Sotho · Norwegian · Pashto · **Persian** · **Polish** · **Portuguese** · Punjabi · **Romanian** · **Russian** · Samoan · Scottish Gaelic · Serbian · Shona · Sindhi · Sinhala · Slovak · Slovenian · Somali · Southern Sotho · **Spanish** · Sundanese · Swahili · Swedish · Tajik · Tamasheq · Tamil · Telugu · Thai · Toba Batak · **Turkish** · Twi · **Ukrainian** · Urdu · Uzbek · **Vietnamese** · Welsh · Wolof · Xhosa · Yiddish · Yoruba · Zulu

*(Languages in **bold** have better performance coverage in Aya Expanse models)

Contents

- 
- 01 The Story of Aya
 - 02 Aya Dataset & Collection
 - 03 Aya Models
 - 04 The People of Aya
 - 05 Responsibility
 - 06 The Aya Movement

01 The Story of Aya

A global initiative led by Cohere For AI to advance the state-of-art in multilingual AI and bridge gaps between people and cultures across the world.

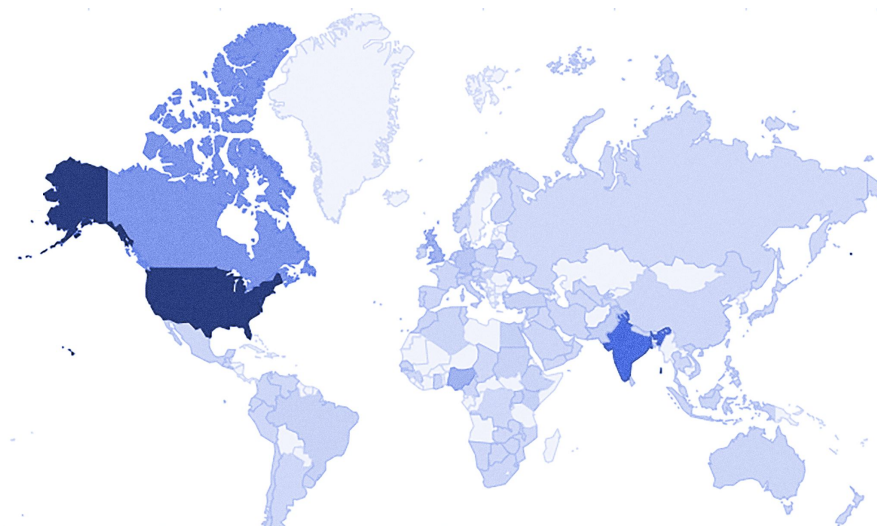
Aya is an open science project to create new models and datasets that expand the number of languages covered by AI, involving over 3,000 independent researchers across 119 countries.

But how did we get here? It all started with a vision to solve complex machine learning problems and an ambitious goal to increase access to language technology for all.

A community, ready to collaborate

The impetus for Aya came out of the [Cohere For AI](#) Open Science initiative - a community that supports independent researchers around the world connect, learn from one another, and work collaboratively to advance the field of ML research.

Starting in January, 2023, members worldwide were keen to leverage the strengths of their diversity and collaborate on something brand new - an open science project to accelerate multilingual AI, and increase access to this technology for the people of their regions.



Join our Open Science Community



Involving 3000+ researchers around the world

Aya is as much a protest against how research is done as it is a technical contribution. Most breakthroughs to-date have come from a small set of labs and countries. Aya instead started with a revolutionary premise: working with independent researchers, engineers, linguists, language enthusiasts around the world to defy expectations and build a breakthrough model.

A widening gap.

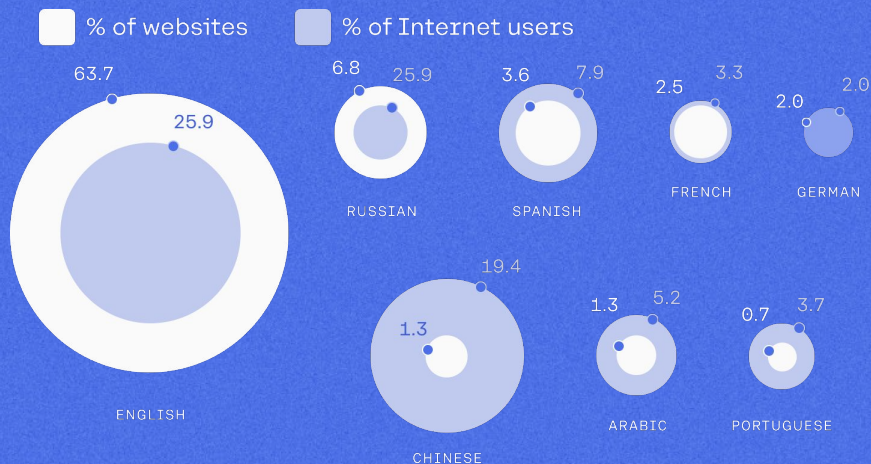
The impetus for this project stems from the stark reality that while natural language processing technologies have advanced exponentially, not all languages have been treated equally by developers and researchers. A significant drawback lies in the source of data used to train large language models, predominantly originating from the internet.

Language	# of papers per million speakers	# of speakers (in millions)
Irish	5235	0.2
Basque	2430	0.5
German	179	83
English	63	550
Chinese	11	1000
Hausa	1.5	70
Nigerian Pidgin	0.4	30

Van Esch, et al. 2022. [Writing System and Speaker Metadata for 2,800+ Language Varieties](#). In *Proceedings of the Thirteenth Language Resources and Evaluation Conference*, pages 5035–5046, Marseille, France. European Language Resources Association.

English is the internet's dominant language

Share of websites using selected languages vs. estimated share of internet users speaking those languages*



*Websites as of February 2022, internet users as of 2021.
Sources: W3Techs, Internet World Stats

This mirrors the early adoption stage of this technology, where a mere 5% of the world's population speaks English at home, yet a surprising 63.7% of internet communication is in English. This trend inadvertently widens the gap in language access to new technologies, exacerbating disproportionate representation, and perpetuating this divide further.

Richter, F. (2022, February 21). English Is the Internet's Universal Language. *Statista*. <https://www.statista.com/chart/26884/languages-on-the-internet/>

Endurance and resourcefulness

The name Aya originates from the Twi language, meaning "fern," symbolizing endurance and resourcefulness – a perfect testament to the movement's commitment to accelerating multilingual AI progress. What we didn't realize when we named the project was how much endurance and resourcefulness we would need to pull it off.



“ If you want to go fast,
go alone.
If you want to go far,
go together.”

– African Proverb

Creating together

Aya has been the largest open-science project in the field of AI. Bringing together 3,000+ collaborators from 119 countries is no small feat. In addition to all the typical challenges of working in groups, we had to take into account time differences, language barriers, various culture understandings and resource inequity.

We hope our journey will help serve as a case study for future participatory research initiatives. We share both the challenges as well as the unique advantages of working together on this mega-scale scientific initiative.

One step down a long road

The Aya models and dataset are released openly, inviting researchers and developers to build upon this progress and conduct further research and build tools to increase access for people in their communities.

By leveraging the Aya resources, you can contribute to the larger challenge of shifting the focus of technological development to encompass all communities and their unique languages.



[Visit the Aya website](https://cohere.com/research/aya)

Together, we can create the future of AI advancement that benefits all.

Let us unite, collaborate, and unleash the full potential of open science for the betterment of global communication.

02 Aya Dataset & Collection



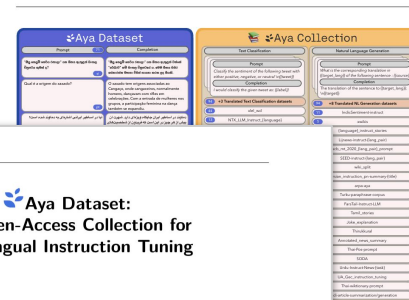


Aya Dataset

An Open-Access Collection for Multilingual Instruction Fine-Tuning

The Aya Dataset represents the most extensive compilation of multilingual instructional examples to date, and it is accessible for use under a fully permissive licensing framework.

For the full paper, read [here](#).



Aya Dataset: An Open-Access Collection for Multilingual Instruction Tuning

Abstract

Datasets are foundational to many breakthroughs in modern artificial intelligence. Many recent achievements in the space of large language models (LLMs) can be attributed to the fine-tuning of pre-trained models on a diverse set of tasks that enables an LLM to respond to instructions. Unlike pre-training, instruction fine-tuning requires the collection of specifically constructed and annotated datasets. However, the creation of such datasets has almost entirely centered on the English language. In this work, our primary goal is to bridge the language gap by building a human-curated instruction-following dataset spanning 71 languages. We worked with native speakers from around the world to collect natural instances of instructions and completions. **Aya** contributes three key resources: we develop and open source the **Aya Annotation Platform**, the **Aya Dataset**, and the **Aya Collection**. The **Aya** initiative also serves as a valuable case study in participatory research, involving 2,997 collaborators from 119 countries. We see this as an important framework for future research collaborations that aim to bridge gaps in resources.

1 Introduction

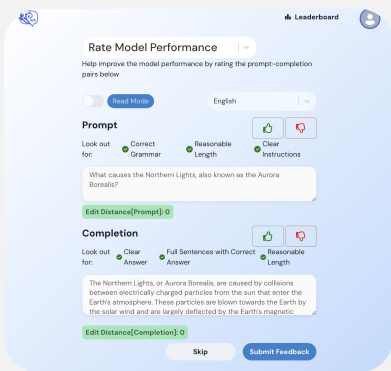
Datasets are static representations of the world, far from the rich ever-evolving environment we navigate as humans. Yet, these frozen snapshots in time are the foundation upon which progress in AI has been built. Much of the recent progress in language modelling can be attributed to fine-tuning pre-trained models on a diverse set of tasks that enable a Large Language Model (LLM) to follow instructions [McCaun et al., 2018; Sanh et al., 2022; Wei et al., 2022a; Muennighoff et al., 2023; Longpre et al., 2023a]. Instruction fine-tuning (IFT) leverages the precept that Natural Language Processing (NLP) tasks can be described via natural language instructions, such as “What were the reviews like for the Barbie movie?” or “Write a recipe from the following list of ingredients.” This process requires pairs of prompts with expected completions [Ziegler et al., 2020; Ouyang et al., 2022] aiming to capture the variety of ways an LLM can be used in downstream tasks. Yet, the very act of curating this data imparts a viewpoint about what distributions we want our model to represent and what is forgotten. So, what do these widely used datasets tell us about the assumptions underlying these breakthroughs?

More than 7,000 languages¹ are spoken around the world today, with a considerable number facing

¹<https://www.ethnologue.com/>

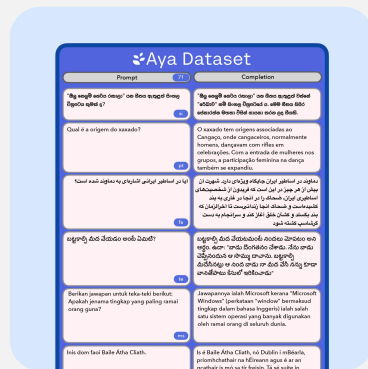


Aya contributes four key resources:



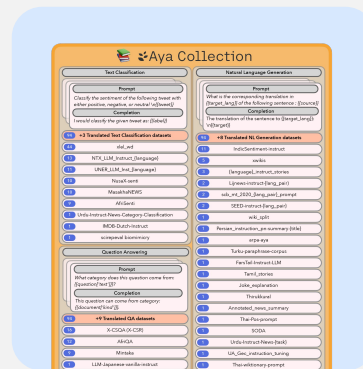
Aya Annotation Platform

A user interface for large-scale participatory research available for free. Used by **2,997 Aya contributors**



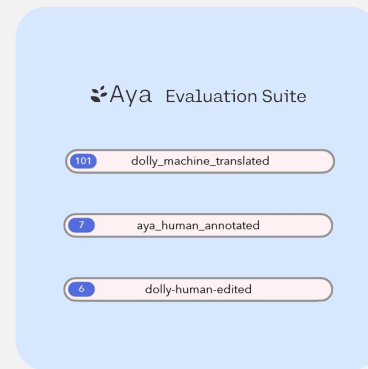
Aya Dataset

The largest human-annotated, multilingual dataset supporting **65 languages**



Aya Collection

A collection of **44 templated and 19 translated datasets**, supporting **115 languages**, to train multilingual LLMs



Aya Evaluation Suite

A high quality dataset for evaluation of LLMs. Subsets include **human-written (7 languages)**, **post-edited translations (6 languages)**, and translations of manually selected prompts (**101 languages**)

Aya Datasets at a glance

Dataset

Download

65 languages

Human-written instances
from fluent native speakers

204K instances

Collection

Download

115 languages

Templating and Translating
existing datasets

513M instances

Evaluation

Download

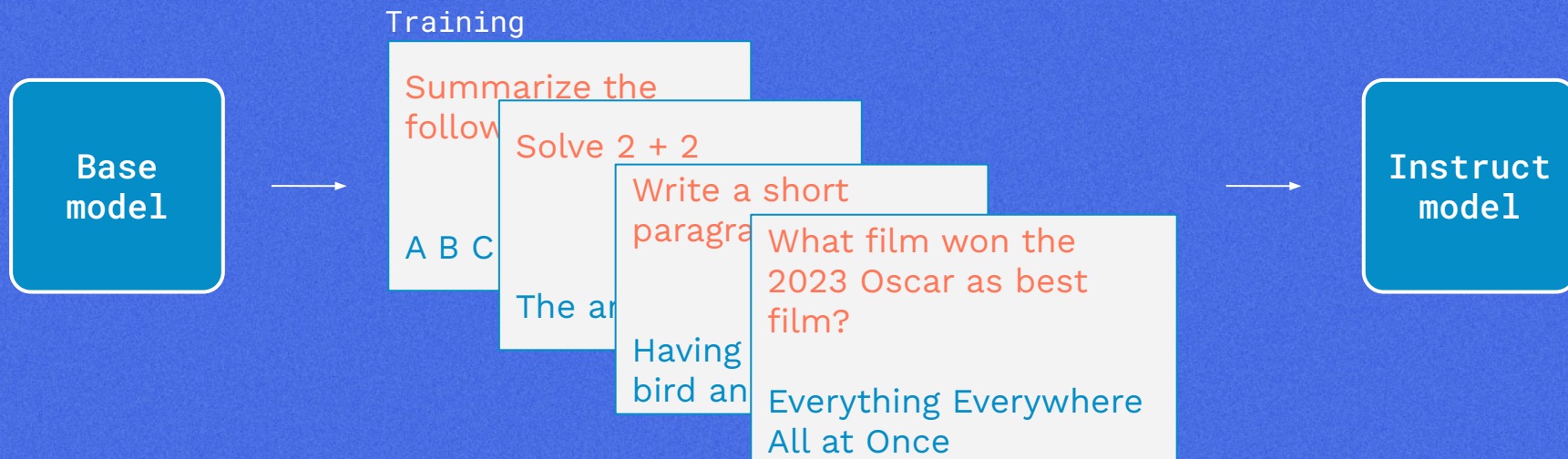
101 languages

Mixture of human-curated,
postedits, and translations

23K instances

What Is Instruction Fine-Tuning?

Instruction Fine-Tuning (IFT) is a form of model training that enables models to better understand and act upon instructions. It is based on the idea that we can use everyday language to ask a model to perform a task and in return the model generates an accurate response in natural language.



Challenges With Multilingual Data Quality and Coverage

To effectively train foundational models with multilingual instructions, we need access to large volumes of quality multilingual instructional data.

This has been plagued by three challenges:



Data scarcity



Low quality data



Lack of qualified contributors for low-resource languages

Without robust multilingual datasets to train models, we risk:



Introducing biases towards languages not included.



Marginalizing speakers of languages not included.



Creating a performance-divide for languages with limited datasets.



Introducing security flaws.

The **Aya** Dataset

The largest
human-curated
multilingual dataset
for fine-tuning LLMs
to follow instructions.

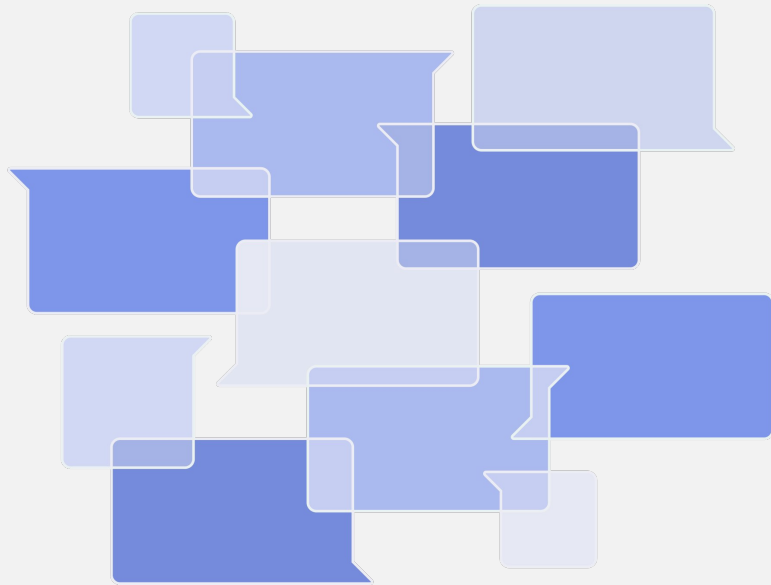
The **Aya** Collection

The **Aya** Evaluation Suite

The Largest Human-Curated Dataset from Native and Fluent Speakers

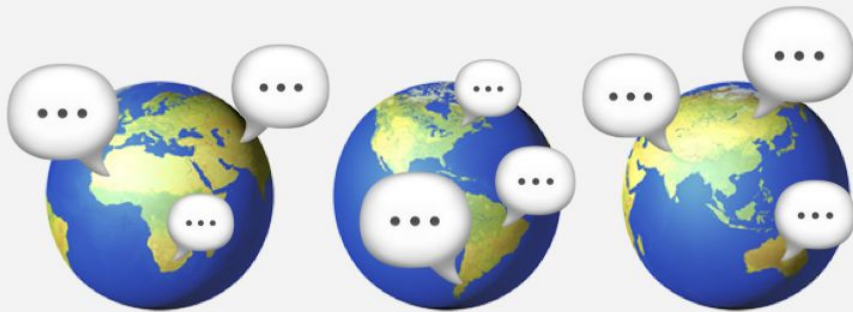
Human-curated data from native and fluent speakers can be hard to come by. It can be costly and difficult to orchestrate.

By leveraging best practices from open-source and crowdsourced science projects, we were able to create the Aya Dataset – the largest collection to date of human-curated and annotated multilingual instruction data.



Aiming for Worldwide Coverage of Languages

Behind each datapoint for each language is a person familiar with the nuances of the language. This level of expertise provides the subtle distinctions and variations in meaning that make each language unique in practice.



Criteria for Inclusion in Aya Dataset

The **Aya Dataset** includes all original annotations and a subset of all re-annotations that vary to a certain extent from the originals.

In order to ensure linguistic diversity and quality, we included languages that were varied, with at least 50 contributions, and with naturally long prompts and corresponding completions.

65 languages

33 high-resource

12 mid-resource

31 low-resource languages

The goal was to include as many languages as possible without lowering the overall quality of the dataset. The table below lists details of the **Aya Dataset**.

Aya Dataset Statistics (number of pairs of prompts and completions obtained through various annotation tasks)

			Count
Original Annotations			138,844
Re-Annotations	xP3 datasets		2,895
	Translated datasets		7,757
	Templated datasets		11,013
	Original Annotations		43,641
Aya Dataset Total			204,114

The Aya Dataset

The Aya Collection

A combination of
human-annotated,
translated, and
templated data.

The Aya Evaluation Suite

An Overview of the Aya Collection

How do we make the world's largest multilingual instruction dataset?



Human Annotated

Human-annotated data is information that has been manually reviewed, labelled, and/or annotated by human annotators, leveraging their native knowledge of a language to provide context and enhance machine learning algorithms.



Translated

Translated multilingual data is when machine translation tools convert text from one language to another, making use of an existing dataset in one language to create the set in another.



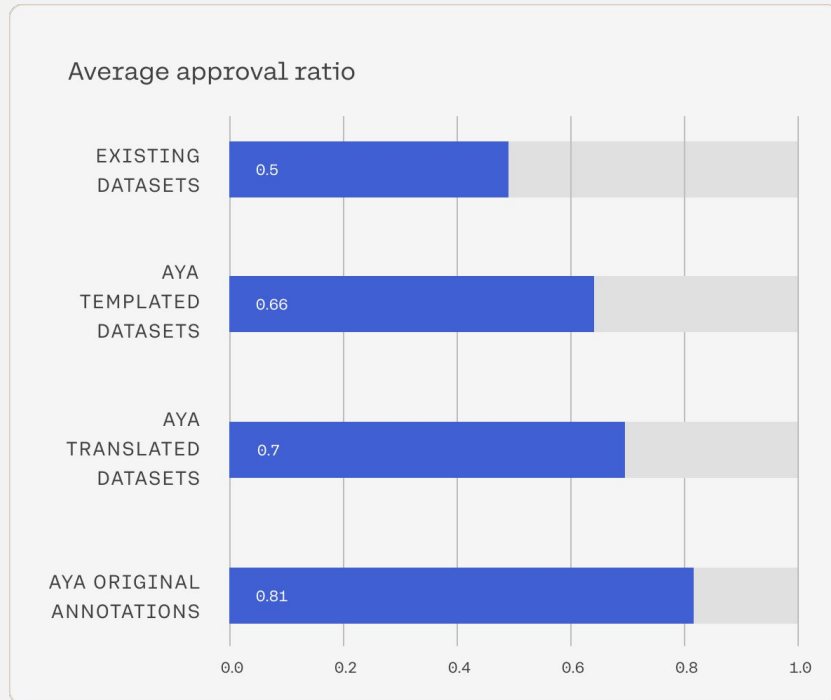
Templated

Templated is created by annotators writing templates and then applying them to datasets to reformat existing NLP datasets into instruction-style.

Aya Collection Surpasses Previous Multilingual Datasets in terms of quality

The quality of instruction data significantly influences the performance of the fine-tuned language model.

Through a global assessment, we enlisted annotators to assess the quality of various multilingual data collections. This process revealed that Aya's original annotations received the highest approval ratings from both native and fluent speakers.



Expanding Data Diversity and Task Coverage

Increasing diversity while maintaining high quality will result in more robust and powerful [1, 2]

We focused on existing datasets templated for instructions and finding tasks that require asking questions and answering based on small pieces of information.

The collection includes 3 main tasks,

- 1) Question Answering
- 2) Natural Language Generation
- 3) Text Classification

and 12 fine-grained task types.

Task Taxonomy of NLP tasks in the Aya Collection

Main Task Type	Fine-grained Task Type
Question Answering	—
Natural Language Generation	Summarization Translation Paraphrasing Dialogue Text Simplification
Text Classification	Sentiment ANALYSIS Information Extraction Named Entity Recognition Event Linking Natural Language Inference Document Representation

The Aya Dataset

The Aya Collection

The Aya Evaluation Suite

A diverse multilingual dataset to assess open-ended generation capabilities of LLMs.

Building an Evaluation Suite

We curate and release an evaluation suite tailored for multilingual models.

This set is a valuable contribution in tackling the scarcity of multilingual data, a challenge that becomes even more apparent when considering evaluation sets.

To strike a balance between language coverage and the quality that comes with human oversight, we create an evaluation suite that includes:

- (1) **human-curated** examples in a limited set of languages,
- (2) automatic **translations** of handpicked examples in an extensive number of languages, and
- (3) **human-post-edited** translations in a few languages.

Human-
curated examples

7 languages

1750 instances

Translations of
hand-picked
examples from
Dolly-15k

101 languages

20K instances

Human-post-
edited translations

6 languages

1200 instances

Limitations of the Aya Dataset

All research has limitations. Below we outline the top challenges faced by the Aya project and results.



Language and dialect coverage: 115 languages (Aya Dataset and Aya Collection) is only a tiny fraction of the world's linguistic diversity.



Uneven distribution of contributions: Relatively few contributors accounted for the most annotations.



Cultural or personal bias: limited representation can lead to a narrow selection of cultural viewpoints.



Gendered pronouns: featuring languages with gendered pronouns or lacking gender-neutral ones, requires careful response crafting to maintain gender neutrality.



Formality distinctions: released dataset contains many languages that have varying levels of standardization and differing style guidelines for formal language like honorifics.



Toxic or offensive speech: the annotation platform does not contain specific flags for toxic, harmful, or offensive speech, so it is possible that malicious users could submit unsafe data.



Accounting for mislabeled data: the annotation platform does not contain any components that enable re-labeling the assigned language of annotations.



Coverage of tasks in Aya Collection: the collection only includes 3 main tasks (Question Answering, Natural Language Generation, Text Classification) and 12 fine-grained task types.

03 Aya Models



Introducing the Aya Models

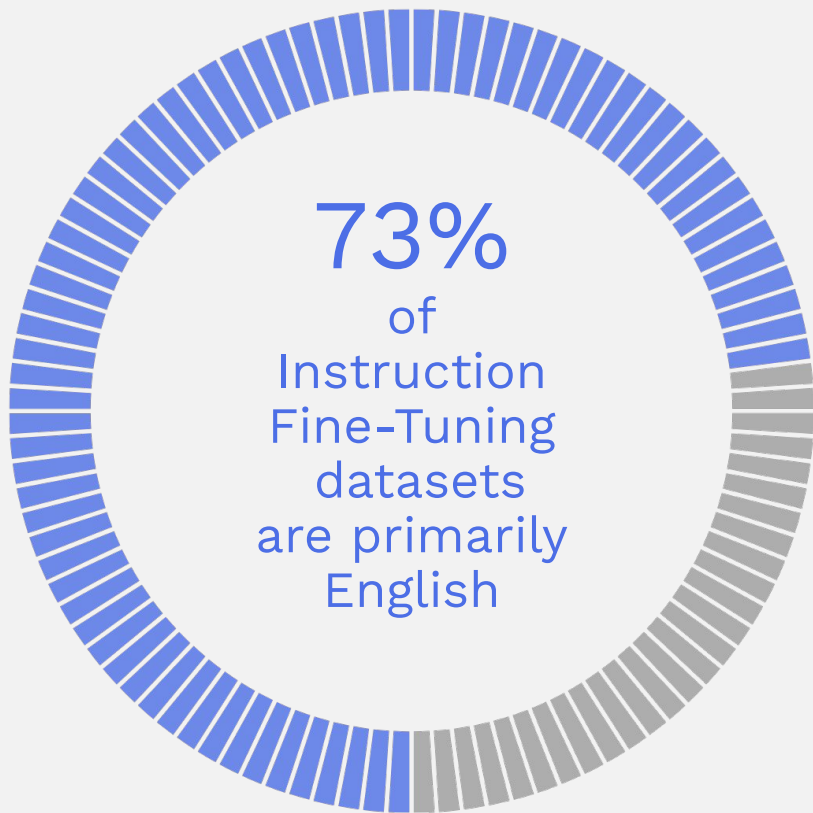
The landscape of modern machine learning has been profoundly shaped by datasets. Yet, this progress has predominantly favored a few data-rich languages due to legacy use and lack of accessible resources. The global linguistic diversity is not represented.

This skew contrasts sharply with a core machine learning principle: **training data should mirror the real-world's vast linguistic diversity.**

We face a glaring inclusivity gap.

“ The limits of my
language
means the limits of my
world. ”

– Ludwig Wittgenstein



The Aya Model aims to bridge this divide, pushing for multilingual IFT datasets that truly reflect our world's rich tapestry of languages, making machine learning not just smarter, but more equitable and representative.

Prompt:

What are some languages spoken in Mexico?

Output:

The three most spoken languages in Mexico are Spanish, Nahuatl, and Maya.

The Aya Models Explained

The Aya Models are designed to tackle linguistic inequality. They can execute tasks in response to prompts given in any supported language. This eliminates the need for multilingual speakers to default to English when writing prompts.

Our goal is to greatly expand the coverage of languages to 101, far beyond the current coverage of previous instruction fine-tuned multilingual models.

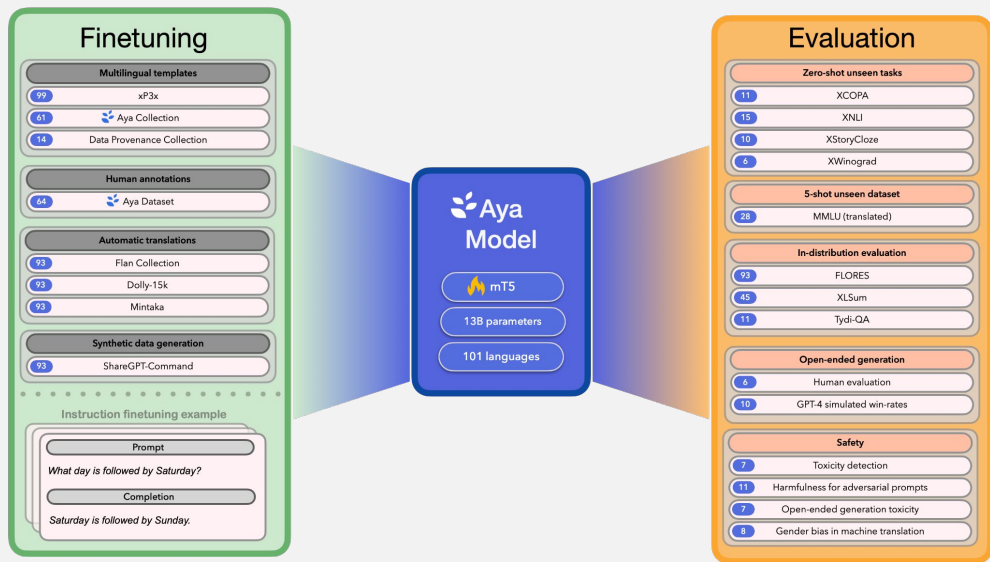


Figure 2: Aya 101 involved extensive contributions to both the breadth of IFT training dataset, optimization techniques including weighting of datasets and introducing more extensive evaluation of performance across varied tasks.



There are three models in the Aya family



STATE OF THE ART, ACCESSIBLE RESEARCH LLM

Aya Expanse - 8B



STATE OF THE ART RESEARCH LLM

Aya Expanse - 32B



MASSIVELY MULTILINGUAL RESEARCH LLM

Aya 101

Our first Aya
model was
Aya 101



Representing Linguistic Diversity through Aya 101

To create a model with diverse linguistic representation, we focused on four areas:



Expansion of Language Coverage

We more than doubled the number of languages with 2.5x the size of the starting dataset.



Broadening Multilingual Evaluation

We benchmark on 99 languages with 4 different evaluation categories using 10 datasets.



Leading Multilingual Performance

The Aya Model consistently outperforms various baselines across all multilingual benchmarks.



Safety

We evaluate our model for gender bias, social bias, harmfulness, and toxicity across languages.

Recipe for building Aya 101



We fine-tuned pretrained multilingual T5 (mT5) language model using instructions in 101 languages

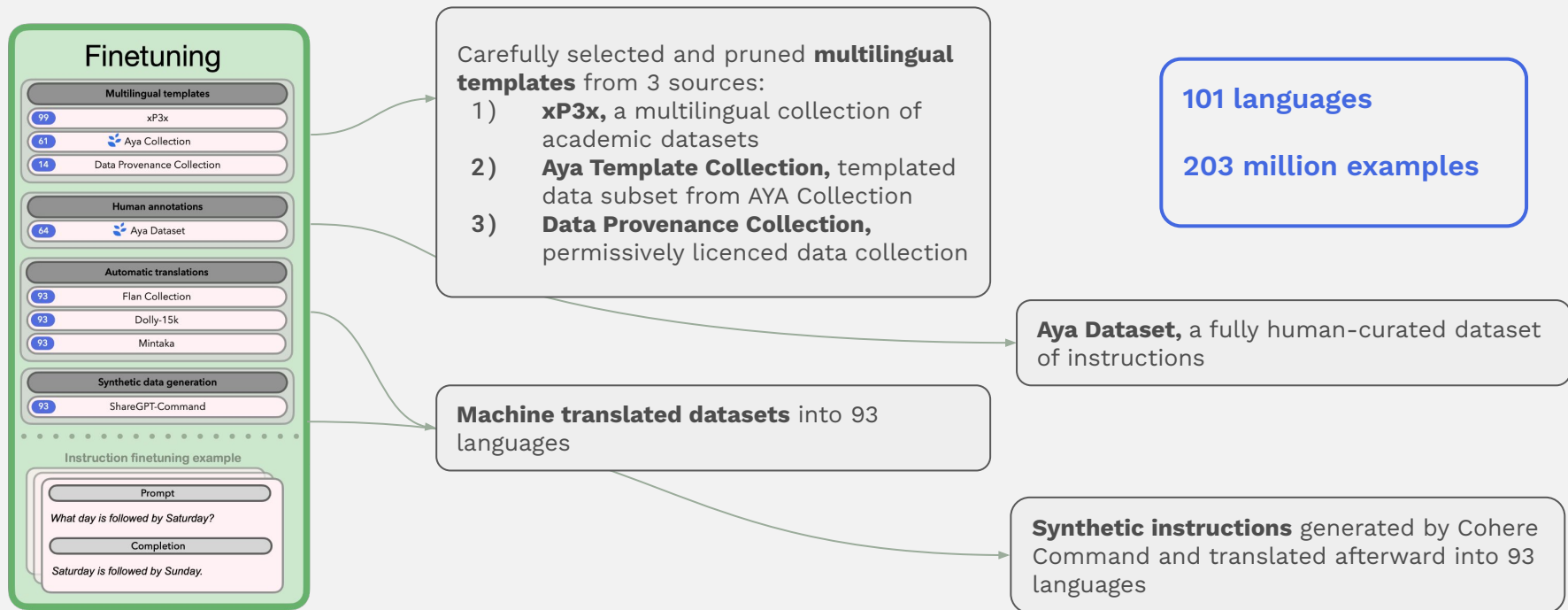


We carefully selected data sources and further prune them to have high quality and diverse set of instruction datasets

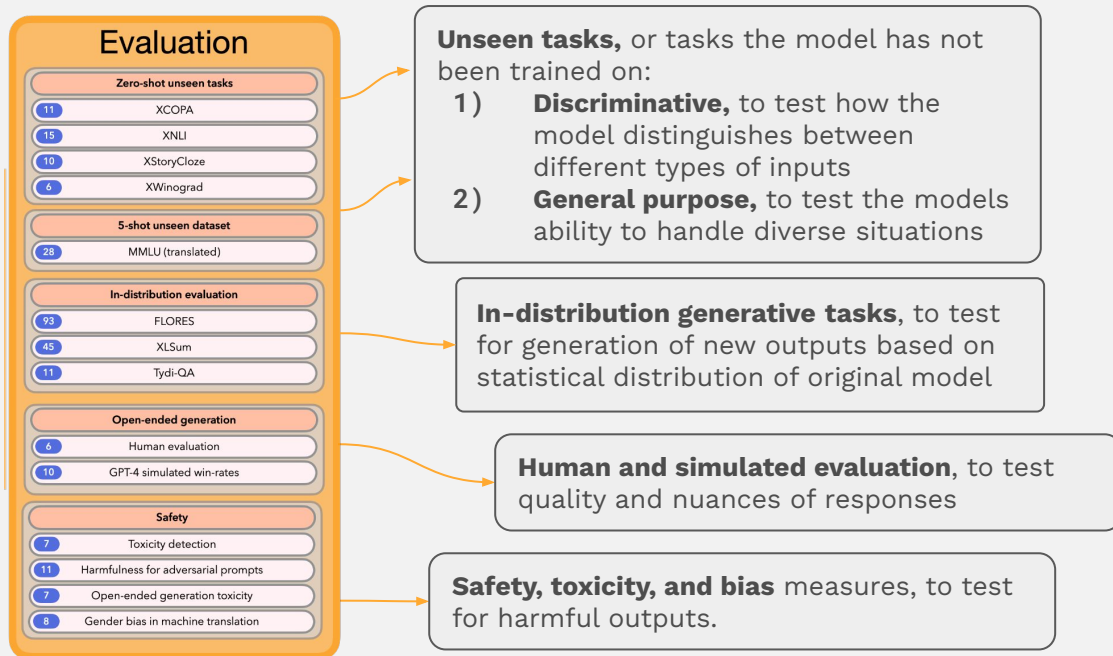


We balanced different data sources during fine-tuning, resulting in high performance across several category of tasks

Building a Massively Multilingual and Diverse Instruction Fine-tuning Mixture



Creating a Massively Multilingual Evaluation Suite



Evaluation at a glance:

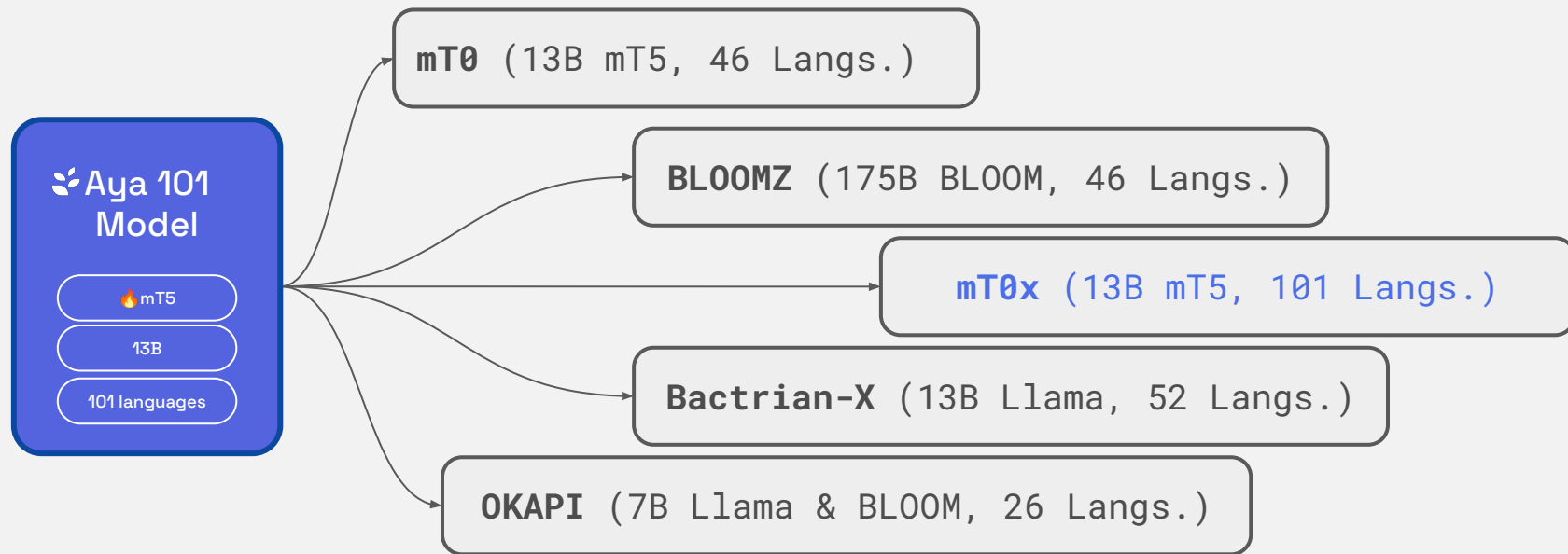
99 languages

13 datasets

6 distinct evaluation types:

- Unseen zero-shot tasks
- General purpose unseen dataset (5-shot)
- In-distribution generative tasks
- Human eval
- LLM simulated eval
- Safety eval

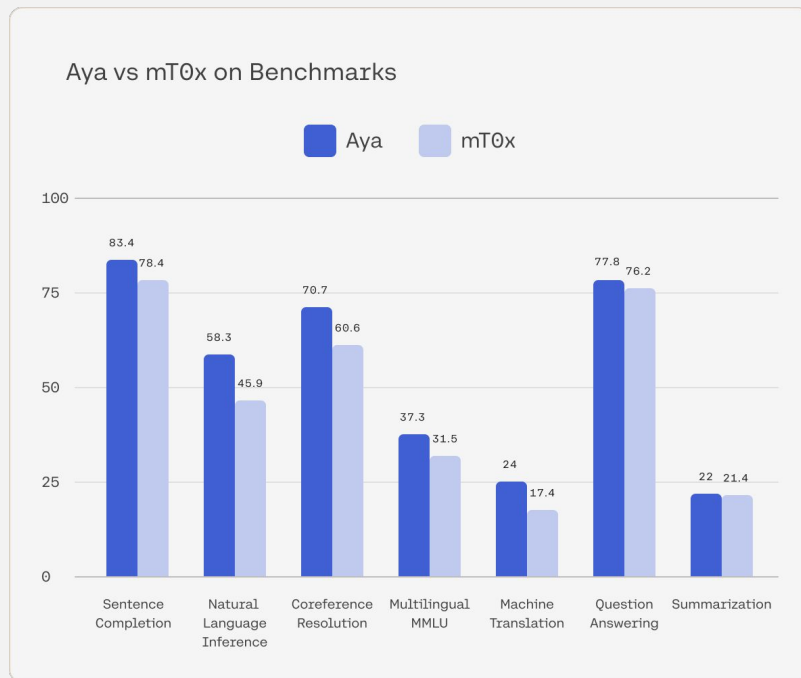
Aya 101 Compared With Multiple Baselines



Advancing Multilingual Performance

Aya 101 achieves superior performance compared to mT0x in the multilingual benchmarks.

These benchmarks include a collection of unseen tasks and in-distribution generative tasks in total covering 100 languages. The Aya model outperforms mT0x in all tasks showing its multilingual capabilities in different task types.

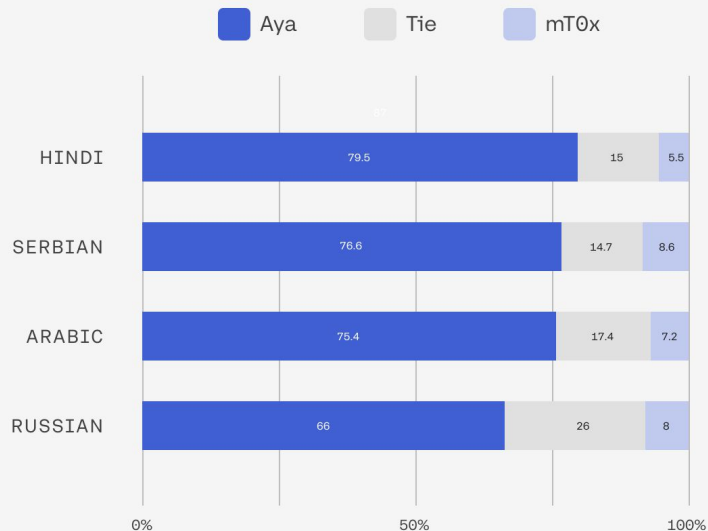


Aya 101 Win Rates

Aya 101 follows instructions and generates responses of significantly higher quality than mT0x.

According to the human evaluation where the professional annotators compared models' responses for given instructions in multiple languages, **the Aya Model is preferred by an average of 77% times.**

Aya win rates against mT0x



Advancing the state of art with Aya ExpansE



Introducing Aya Expanse

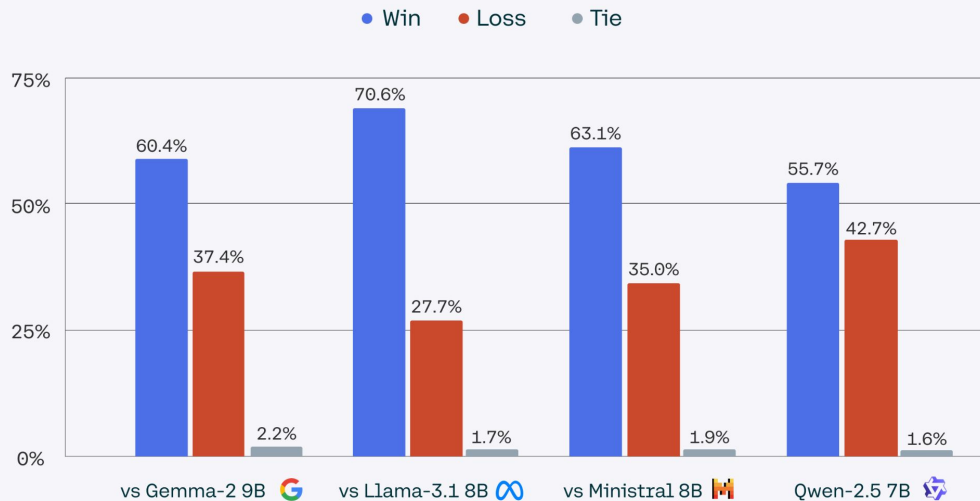
The Aya Expanse models advance the **state of the art in modeling 23 languages which cover half of the world's population:**

Arabic, Chinese, Czech, Dutch, English, French, German, Greek, Hebrew, Hebrew, Hindi, Indonesian, Italian, Japanese, Korean, Persian, Polish, Portuguese, Romanian, Russian, Spanish, Turkish, Ukrainian, and Vietnamese

Leading Multilingual Performance

Aya Expanse achieves **superior performance across 23 languages on difficult, diverse instruction following tasks** when compared to other open weights models including Gemma, Llama, Mistral, and Qwen.

Aya Expanse 8B Win Rates (m-ArenaHard)



Builds on Several Years of dedicated Multilingual Research

Achieving Aya Expanse's leading multilingual performance **required combining years of multiple, dedicated multilingual research efforts**

C4AI Multilingual AI Revolution: Breaking Language Barriers

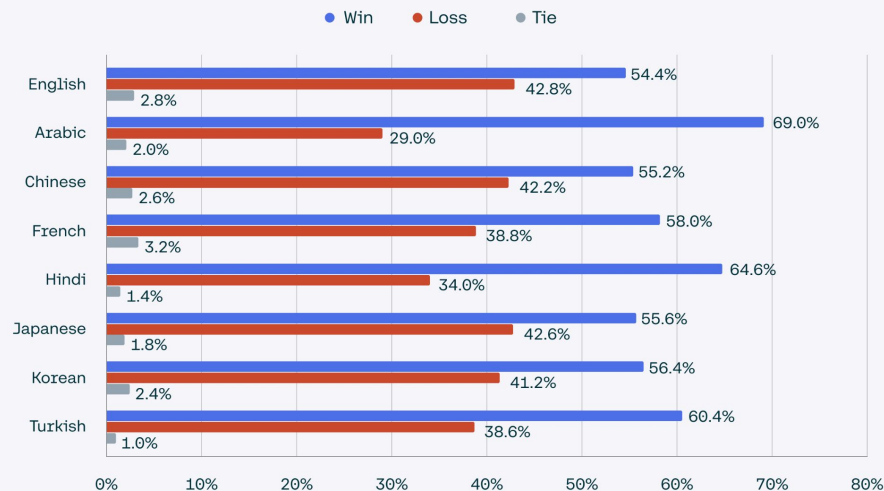


Leading Multilingual Performance

Aya Expanse 8B **outperforms Gemma-2 9B across all 23 languages including English!**

This shows that is possible to advance multilingual performance more equitably for lower resource languages without cannibalizing performance in higher resource languages like English

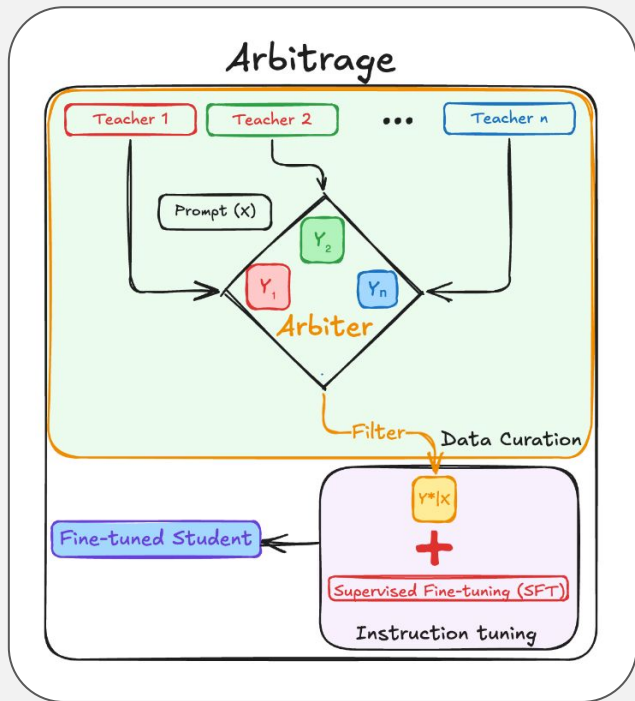
Aya Expanse 8B Language Specific Win Rates vs Gemma-2 9B (m-ArenaHard)



Training Aya Expanse: Arbitrage

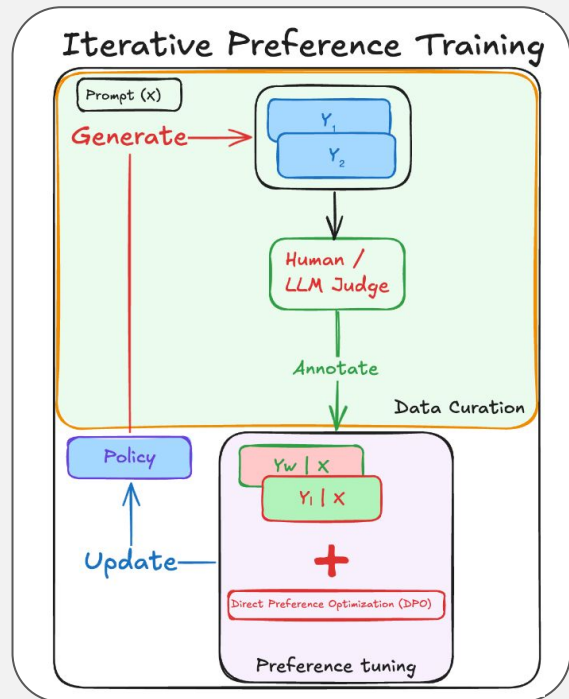
Multilingual Arbitrage: fine-tuning an LLM on the best completion (as determined by an arbiter) from a pool of teacher models

Multilingual Arbitrage enables strategic distillation from a pool of models where any individual teacher model may only be strong in small set of languages or domains



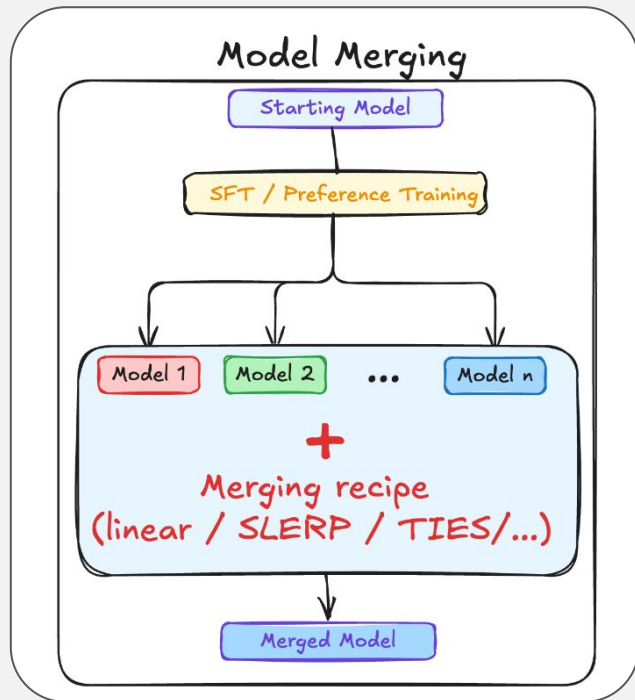
Training Aya Expanse: Preference Training

Aya Expanse is preference-trained by contrasting the best and worst completions from the arbitrage stage, **steering completions away from features of low-quality multilingual completions**



Training Aya Expanse: Model Merging

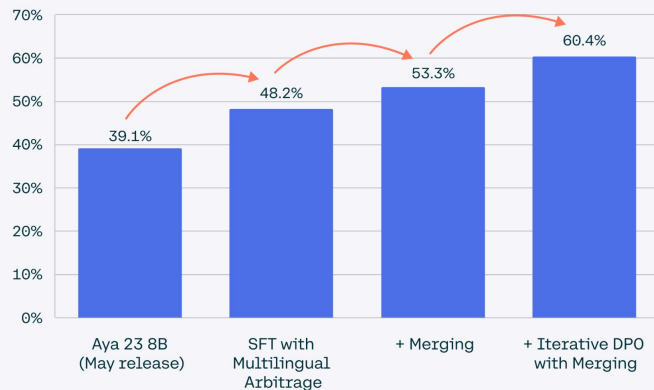
During SFT and RLHF stages of training Aya Expanse, **multiple models trained on different language subsets of the training data are merged** together to produce a single, more performant model across all languages



Training Aya Expanse: Summary

Multilingual arbitrage, multilingual preference training, and model merging were **all critical steps in achieving Aya Expanse's remarkable performance**

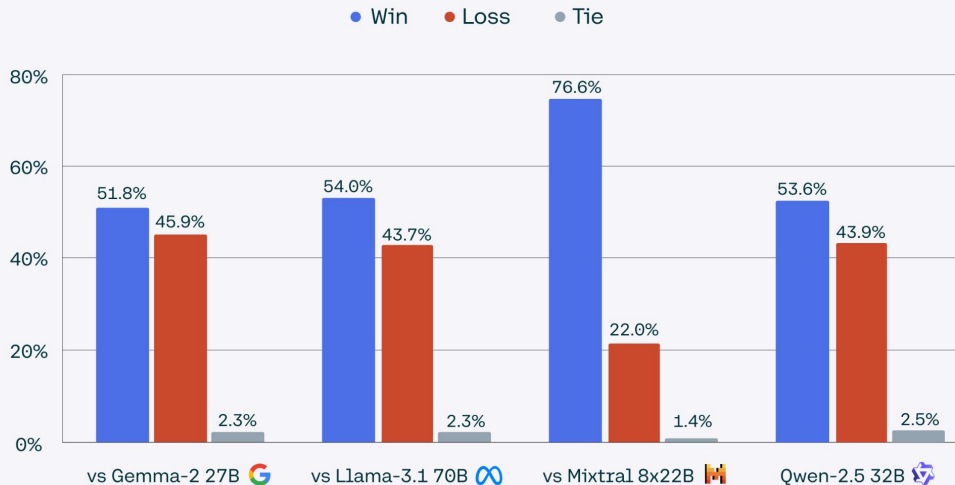
Step by Step Improvement in Win-Rates Against Gemma-2 9B



Scaling Aya Expanse

The same training recipe scales to 32B parameter scale, **outperforming competitor open weights models including LLMs with many more parameters!**

Aya Expanse 32B Win Rates (m-ArenaHard)



Aya Expanse Team

Core Aya Expanse Team

Madeline Smith, Marzieh Fadaee, Ahmet Üstün, Beyza Ermis, Sara Hooker, John Dang, Shivalika Singh, Arash Ahmadian, Daniel D'souza, Alejandro Salamanca, Aidan Peppin, Arielle Bailey, Meor Amer, Sungjin Hong, Manoj Govindassamy, Sandra Kublik

Wider Cohere For AI and Cohere Contributors

Acyr Locatelli, Adrien Morisot, Jon Ander Campos, Sara Elsharkawy, Eddie Kim, Julia Kreutzer, Nick Frosst, Aidan Gomez, Ivan Zhang





Aya Expanse Language Ambassadors













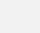
56

We create breakthroughs together. Ambassadors represent 45 countries and 23 languages. Before the launch of Aya Expanse, we invited 110 ambassadors to join us to shape how Aya worked for communities all over the world.

 Mehmet Emre Akbulut
 Samer Attrah
 MUHDIN AWOL
 Kenza Benkirane
 Mann Bhanushali
  Isabella Bicalho Frazeto
 Danylo Boiko
 Sabrina Boumaiza
 Samuel Cahyawijaya
 Samuel Cahyawijaya
 Emirhan Çelik
 Ryan Chan
 Aurélien-Morgan CLAUDON
 Urszula Czerwinska
 Joana da Matta
 Nguyễn Đạt
 Manasvi Dawane
 Akanksha Devkar
 Sharad Duwal
 Abdeljalil EL MAJJODI
 Shafagh Fadaei
 Neil Fernandes
 Silvia Fernandez
 Hamidreza Ghader
 Manuel Goulão
 Bassam Gouti
 María Grandury
 Miguel Guerrero

 Siddhesh Gunjal
 Mohammed Hamdy
 Hafedh Hichri
  Nhu Hoang Anh Quynh
 Kyle Howard
 Jiwung Hyun
 Joseph Marvin Imperial
 Burin Intachuen
 Ryan Junejo
 Juan Junqueras
  Karthik Reddy Kanjula
 Albert Kao
 Morteza Kashani
 Ahmed Khaled
 Niharika Khanna
  Dipika Khullar
 Christopher Klamm
 Nazar Kohut
  Alkis Koudounas
 Diana Kozachek
 Katrina Lawrence
 James León
 Jiazheng Li
 Nicolò Loddò
 Dante, Fu On Lok
 Iro Malta
 Harras Mansoor
 Bhavnick Minhas

 Shachar Mirkin
 Roa'a Mohammad
 Yiyang Nan
  Sree Harsha Nelaturu
 Jekaterina Novikova
 Roni Obaid
 Olympiah Otieno
 Enes Özgözler
 Yavuz Alp Sencer Ozturk
 Carlos Patiño
 Jebish Purbey
 Maria Quijano Jesurum
 Swati Rajwal
 Didi Ramsaran Chin
 Divyaraj Rana
 Aditya Retnanto
 Rodrigo Ribeiro Gomes
 Esra'a Saleh
  Roshan Santhosh
  Drishti Sharma
 Kinza Sheikh
  Aditya Shrivastava
 Vivek Silimkhan
 Marjana Skenduli
 Soham Sonar
 Gürkan Soykan
 David Styveen
 Anthony Susevski

 Adrian Szymczak
 Joanne Tan
 Quentin Tardif
  Ameer Taylor
 Yiorgos Tsalikidis
 Roman Tymtsiv
 Muhammad Saad Uddin
 Louis Ulmer
 Sundar Sripada V. S.
 Freddie Vargus
 Vlad Vasilescu
 Karan Verma
 Henry Vo
  Minh Chien Vu
 Hieu Vu
 Azmine Toushik Wasi
 Warren Williams
 Joseph Wilson
 Gusti Winata
 Ege Yakut
 Eray Yapağcı
 Taha Yassine
 Serhan YILMAZ
 Hanna Yukhymenko
 Mike Zhang

04 The People of Aya



The Frontiers of Participatory Research

Language is a deeply social phenomenon for its everyday users. It thrives on a network of social relations. However, there is no template or rulebook for working with 3000+ researchers and enthusiasts around the world. Instead, we kept in mind some guiding principles:



Fluid Ownership and Growth

A decentralized model supports fluid leadership and flexible role adoption. It empowers members to take initiative independent of hierarchical position or level of involvement.



Organizational Structure

Asynchronous communication channels facilitate rich and timely collaborations.



Inclusion and Access

Bypass academic norms that often marginalize non-English speakers and people without formal academic credentials.



Participating motivators

Not based on financial remuneration but on ideals of community, identity, and social justice.

Whenever we engage with data, we are also engaging with the connections that data has to the people who produce it, prepare it, and distribute it.

The Journey of Aya

Watch [The Journey of Aya](#), a short documentary in which our collaborators tell the story of how Aya came to be.

Aya 101 Core team 1/2

Listed in alphabetical order.

The Core Team has been responsible for various technical elements of making our Aya 101 models and dataset a reality. Their contributions varied across building an accessible user interface, establishing strong baselines, exploring data augmentation strategies, ensure responsible deployment, and coordinating regional contributions.



Aisha Alaagib
Cohere For AI
Community



Emad A. Alghamdi
King Abdulaziz U
ASAS, AI



Zaid Alyafei
King Fahd University of
Petroleum and
Minerals or KFUPM



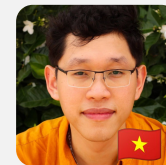
Viraat Aryabumi
Cohere For AI



Max Bartolo
Cohere



Neel Bhandari
Cohere For AI
Community



Vu Minh Chien
Cohere For AI
Community



Daniel D'souza
Cohere For AI



Irem Ergun
Cohere



Ellie Evans
Cohere For AI
Community



Marzieh Fadaee
Cohere For AI



**Hakimeh
(Shafagh) Fadaei**
Cohere For AI
Community



**Sebastian
Gehrmann**
Bloomberg LP



**Ramith
Hettiarachchi**
MIT



Sara Hooker
Cohere For AI



Sarah Jafari
Cohere For AI



Börje Karlsson
Beijing Academy of
Artificial Intelligence
(BAAI)



Amr Kayid
Cohere



Farhan Khot



Wei-Yin Ko
Cohere



Julia Kreutzer
Cohere For AI

Aya 101 Core team 2/2

Listed in alphabetical order.

The Core Team has been responsible for various technical elements of making our Aya 101 models and dataset a reality. Their contributions varied across building an accessible user interface, establishing strong baselines, exploring data augmentation strategies, ensure responsible deployment, and coordinating regional contributions.



Dominik Krzeminski
Cohere For AI Community



Shayne Longpre
MIT



Marina Machado
Cohere



Abinaya Mahendiran
Cohere For AI Community



Deividas Mataciunas
Cohere For AI Community



Oshan Mudannayake
Cohere For AI Community



Niklas Muennighoff
Cohere For AI Community

T



Laura O'Mahony
University of Limerick, Limerick, Ireland



Ifeoma Okoh
Cohere For AI Community



Gbemileke Onilude
Carnegie Mellon University



Hui-lee Ooi
Cohere For AI Community



Jay Patel
Binghamton University, NY, USA



Herumb Shandilya
Cohere For AI Community



Shivalika Singh
Cohere For AI Community



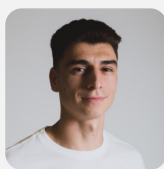
Madeline Smith
Cohere For AI



Luisa Souza Moura
Cohere



Ahmet Üstün
Cohere For AI



Freddie Vargus
Cohere For AI Community



Joseph Wilson
University of Toronto



Mike Zhang
IT University of Copenhagen



Yong Zheng Xin
Brown University
Cohere For AI Community

Aya 101 Language Ambassadors 1/3

Listed in alphabetical order.

Language Ambassadors spread the word about Aya to speakers of their language, recruit new contributors, support those contributors to understand the goals of Aya data collection efforts, and celebrate progress.



Diana Abagyan
Russian



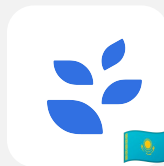
Muhammad
Abdullahi
Somali



Elyanah Aco
Filipino



Henok
Ademtew
Amharic



Adil
Kazakh



Emad A.
Alghamdi
Arabic



Zaid Alyafeai
Arabic



Ahmad Anis
Urdu



Daniel Avila
Spanish



Michael
Bayron
Cebuano



Nathanael Carraz
Rakotonirina
Malagasy



Alberto Mario
Ceballos Arroyo
Spanish



Yi Yi Chan Myae
Win Shein
Burmese



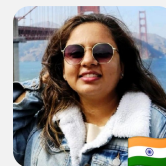
Vu Minh Chien
Vietnamese



Caroline Shamiso
Chitongo
Zulu



Ionescu
Cristian
Romanian



Rupal Darji
Gujarati



Suchandra
Datta
Bengali



Rokhaya
Diagne
Wolof



Irem Ergun
Turkish



Hakimeh
(Shafagh) Fadaei
Persian

Aya 101 Language Ambassadors 2/3

Listed in alphabetical order.

Language Ambassadors spread the word about Aya to speakers of their language, recruit new contributors, support those contributors to understand the goals of Aya data collection efforts, and celebrate progress.



Surya Krishna
Guthikonda
Telugu



Aleksandra
Hadžić
Serbian



Shamsuddeen
Hassan
Muhammad
Hausa



Ramith
Hettiarachchi
Sinhala



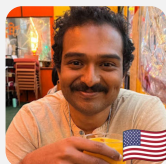
Mochamad
Wahyu Hidayat
Sundanese



Rin Intachuen
Thai



Eldho Ittan
George
Malayalam



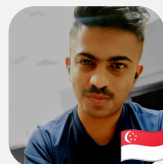
Ganesh
Jagadeesan
Hindi



Murat
Jumashev
Kyrgyz



Börje Karlsson
Portuguese and
Swedish



Abhinav
Kashyap
Kannada



JiWoo Kim
Korean



Alkis
Koudounas
Italian



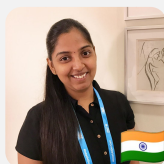
Kevin Kudakwashe
Murera
Shona



Falalu Ibrahim
Lawan
Hausa



Wen-Ding Li
Traditional Chinese



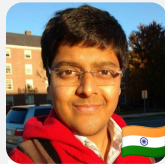
Abinaya
Mahendiran
Tamil



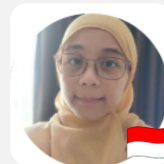
Mouhamadane
Mboup
Wolof



Oleksander
Medyuk
Ukrainian



Pratik Mehta
Hindi



Iftitahu Nimah
Javanese

Aya 101 Language Ambassadors 3/3

Listed in alphabetical order.

Language Ambassadors spread the word about Aya to speakers of their language, recruit new contributors, support those contributors to understand the goals of Aya data collection efforts, and celebrate progress.



Solam Nyangiwe
Xhosa



Laura O'Mahony
Irish



Ifeoma Okoh
Igbo



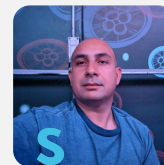
Hui-Lee Ooi
Malay



Iñigo Parra
Basque



Jay Patel
Gujarati



Hanif Rahman
Pashto

B



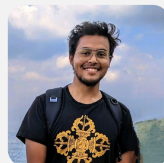
Olanrewaju Samuel
Yorùbá



Suman Sapkota
Nepali



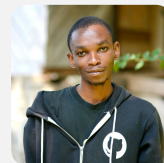
Giacomo Sarchioni
Italian



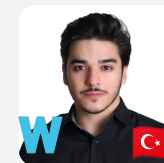
Rashik Shrestha
Nepali



Bhavdeep Singh
Punjabi



Sean Andrew Thawe
Chichewa



Alperen Ünlü
Turkish

W



Joseph Wilson
French



Emilia Wiśnios
Polish



Yang Xu
Simplified Chinese



Zheng-Xin Yong (Yong)
Malay



Mike Zhang
Dutch

K

Top 50 Quality Champions 1/2

Collaborators listed in ascending order based on Aya Quality Score.

These collaborators lead the way in ensuring the textual data contributed to Aya 101 was of high quality including being free of grammatical errors, safe and factually correct, and robust completions to support model training.


 Vu Minh Chien

 Hui-Lee Ooi

 Gamage Omega Ishendra


 Surya Krishna Guthikonda

 Hoang Anh Quynh Nhu

 Moses Oyeleye

 Amarjit Singh Sachdeva

 Mike Zhang

 Almazbekov Bekmyrza
Ruslanovich


 Ramla Abdullahi
Mohamed

 Börje F. Karlsson

 Regina Sahani Lourdes De
Silva Goonetilleke

 Zaid Alyafeai

 Yong Zheng Xin

 Yavuz Alp Sencer Öztürk

 Mohammed Hamdy

 Anitha Ranganathan

 Ramith Hettiarachchi

 Ooi Hui Yin

 Caroline Shamiso
Chitongo

 Bhavdeep Singh Sachdeva

 Valentyn Bezshapkin

Top 50 Quality Champions 2/2

Collaborators listed in ascending order based on Aya Quality Score.


These collaborators lead the way in ensuring the textual data contributed to Aya 101 was of high quality including being free of grammatical errors, safe and factually correct, and robust completions to support model training.

 Yang Xu

 Dominik Krzeminski

 Iftitahu Nimah


 Muna Mohamed Abdinur

 Nurbaeva Zhiidegul
Talaibekovna

 Younes Bensassi Nour

 Eldho Ittan George

 Caio Dallaqua

 Hakimeh (Shafagh) Fadaei

 Henok Ademtew

 Vijayalakshmi Varadharajan

 Yogesh Haribhau Kulkarni

 Laura O'Mahony

 Jay Patel

 Luísa Souza Moura

 Rama Hasiba


 Geoh Zie Ee

 Gabriela Vilela Heimer

 Pratham Prafulbhai
Savaliya

 Deividas Mataciunas

 Ifeoma Okoh

 Alberto Mario Ceballos
Arroyo

 Basiiru Silla

 Yiorgos Tsalikidis

Dataset Champions

Collaborators listed in alphabetical order.

Aya 101 Dataset Champions sourced, formatted and submitted open-source datasets in their languages to be included in the Aya collection.



Diana Abagyan



Md. Tahmid Hossain



Abinaya Mahendiran



Henok Ademtew



Eldho Ittan George



Desik Mandava



Ahmad Anis



Ganesh Jagadeesan



Iftitahu Nimah



Hakimeh (Shafagh) Fadaei



Börje F. Karlsson



Wannaphong Phatthiyaphaibun



Hamidreza Ghader



Surya Krishna Guthikonda




Mike Zhang


5000 Contribution Points

Collaborators listed in descending order of most points earned.

These contributors achieved at least 5000 Contributions Points via the Aya data collection user interface.

 Moses Oyeleye

 Vu Minh Chien

 Ramla Abdullahi
Mohamed


 Gamage Omega Ishendra

 Nitta Sitakrishna

 Surya Krishna Guthikonda

 Hui-Lee Ooi


 Hoang Anh Quynh Nhu


 Nurbaeva Zhiidegul
Talaibekovna


 Muna Mohamed Abdinur

 Amarjit Singh Sachdeva

 Yang Xu


 Almazbekov Bekmyrza
Ruslanovich

 Ahmed Mohamed Hussein
Malin

 Bhavdeep Singh Sachdeva

 Yong Zheng Xin

 Yavuz Alp Sencer Öztürk

 Regina Sahani Lourdes De
Silva Goonetilleke

 Yogesh Haribhau Kulkarni

 Zaid Alyafeai

 L N Deepak

 Caroline Shamiso
Chitongo

 Börje F. Karlsson

 Younès Bensassi Nour

1000 Contribution Points 1/3

These contributors achieved at least 1000 Contributions Points via the Aya data collection user interface.


Contributors listed in descending order from most points.

 Sudharshini AJ	 Gabriela Vilela Heimer	 Sefika Efeoglu	 Rafael Panisset Motta
 Maryam Sabo Abubakar	 Júlia Souza Moura	 Abdishakuur Mohamed Hussein	 Jay Patel
 Mr. A. Karthik	 Suchandra Datta	 Hakimeh (Shafagh) Fadaei	 Zalkarbek Tilenbaev
 Mike Zhang	 Laura O'Mahony	 Luísa Souza Moura	 Meghana Denduluri
 Caio Dallaqua	 Valentyn Bezshapkin	 Iñigo Parra	 Abdou Sall
 Rokhaya Diagne	 Makomborero Magaya	 Razafindrakotonjatovo Zo Anjatiana Henitsoa Kokoly	 Nathanaël Carraz Rakotonirina
 Anitha Ranganathan	 Taqi Haider	 Aidaiym Omurbekovna	 Dr. Maharasan.K.S
 Eldho Ittan George	 R. A. Nirmal Sankalana	 Ripal Darji	 Khaleel Jageer
 Dominik Krzeminski	 Basiiru Silla	 Mr. MARAPPAN .A	 Falalu Ibrahim Lawan
 Rama Hasiba	 Ramith Hettiarachchi	 NDIMBIARISOA Valdo Tsiaro Hasina	 Iftitahu Nimah
 Dev Haral	 Yat Kan Eden Cheung		 Armeen Kaur Luthra

1000 Contribution Points 2/3

These contributors achieved at least 1000 Contributions Points via the Aya data collection user interface.








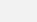


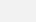
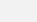


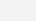
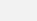
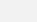

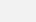


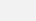
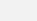

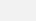
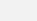
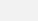
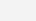
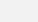
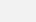
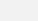
Contributors listed in descending order from most points.

 Elyanah Marie Aco	 Alberto Mario Ceballos Arroyo	 Md. Tahmid Hossain	 Ainura Nurueva
 Adeer Khan	 Geoh Zie Ee	 Henok Ademteu	 Hollie O'Shea
 Ooi Hui Mei	 Andriatsalama Fiononantsoa Jaofera	 Mohammed Nasiru	 Wannaphong Phatthiyaphaibun
 Deividas Mataciunas	 Tsaramanga Jeanny Fidelica	 Harena Finaritra Ranaivoarison	 Abubakr Labaran Salisu
 Betel Addisu	 Sean Andrew Thawe	 Mansi Kamlesh Patel	 Ooi Hui Yin
 Randriamanantena Manitra Luc	 Ratsimba Ranto Sarobidy	 Marina Fontes Alcântara Machado	 RAKOTONIRINA Tokinantenaina Mathieu Razokiny
 K.Chinnaraju	 Srinadh Vura	 Tahina Mahatoky	 Robinson Rodrigo Silva Oliveira
 Mouhamadane Mboup	 Benmeridja Ahmed Younes	 Ramarozatovomampionona Todisoa Nirina Mickael	 Hanif Rahman
 Filamatra Manampy Fanantenana Rasolofoniaina	 Elshaday Desalegn Asfaw	 Ana Carolina Correia Pierote	 Maminirina Rahenintsoa
 Amandeep Singh			

1000 Contribution Points 3/3

Contributors listed in descending order from most points.


















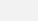
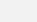


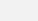
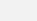


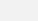
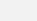

These contributors achieved at least 1000 Contributions Points via the Aya data collection user interface.

 Krishna Chhatbar	 Ifeoma Okoh	 Ijeoma Irene Okoh	 G. A. Jalina Hirushan Gunathunga
 J.Nirmala	 Sumi Shakya	 Ajayi Akinloluwa Irawomitan	 Ogba Stephen Kesandu
 Tharin Edirisinghe	 Alkis Koudounas	 Zarlykov Kelsinbek	 Tiana Kaleba Andriamanaja
 Randrianarison Diarintsoa Fandresena No HerijaonaHerijaona	 Mohamad Aboufoul	 Micol Altomare	 Andriamiadanjato Mioraniaina
 Andrianarivony Harijaona Fanirintsoa	 Emad A. Alghamdi	 Yadnyesh Chakane	
 Rakotondrainibe Nirisoa Tendry	 Jothika. S	 Rafidy Julie Tassia	
 Bekbolot Abdirasulov	 Razakahasina Fanomezana Sarobidy	 Rabin Adhikari	
 Joseph Marvin Imperial	 Valério Viégas Wittler	 Chinwendu Peace Anyanwu	
	 Anish Gasi Shrestha	 Dr. S.P. Balamurugan	
	 Joseph Wilson		

500 Contribution Points

These contributors achieved at least 500 Contributions Points via the Aya data collection user interface.

Contributors listed in descending order from most points.

 M.Neelavathi	 Easwaran K	 Santiago Pedroza Díaz	 Ruqayya Nasir Iro
 Sabita Rajbanshi	 Ahmad Mustafa Anis	 Siyu Wang	 Geetharamani R.
 Silambarasan U.	 Dr.G.Thilagar	 Randinu Jayaratne	 Sandesh Pokhrel
 Dr.A.Prasanth	 Gan Chin Chin	 Rithara Kithmanthie	 Orozbai Topchubek uulu
 Sara Salvador	 Bhanu Prakash Doppalapudi	 Bhanu Prakash Doppalapudi	 Prajapati Maitri R.
 Dr A.Jeba Christy	 Abdullahi Adan Hassan	 TSuman Sapkota	 Francisco Valente
 Mr.V.Balakrishnan	 Sara Hooker	 Charindu Abeysekara	 Gaurav Jyakhwa
 Abinaya Mahendiran	 Amjad Abdulkhaliq Alkhatabi	 Afifah binti Mohd Shamsuddin	 Mrs. G. Sangeetha
 Solam	 Muhamad Audi Bin Pasha	 Verassree Rajaratnam	 Ahmet Güneyli
 Rashik Shrestha			

Public Release and Engineering Team 1/2

Collaborators listed in alphabetical order.

The public release team is responsible for bringing Aya to the world. From building and deployment of the model, planning the launch event, creating *The Journey of Aya* documentary, hosting the model and coordinating outreach efforts.

 Viraat Aryabumi

 Saurabh Baji

 Max Bartolo

 Claude Beaupré

 Phil Blunsom

 Tomeu Cabot

 Isabelle Camp

 Jon Ander Campos

 Claire Cheng

 Linus Chui

 Jenna Cook

 Natasha Deichmann

 Roy Eldar

 Irem Ergun

 Beyza Ermis

 Marzieh Fadaee

 Ramy Farid


 Nick Frosst

 Josh Gartner

 Aidan Gomez

 Manoj Govindassamy

 Rod Hajjar

 Sara Hooker

 Monica Iyer

 Sarah Jafari

 Amr Kayid

 Julia Kedrzycki

 Wei-Yin Ko


Public Release and Engineering Team 1/2

Collaborators listed in alphabetical order.

The public release team is responsible for bringing Aya to the world. From building and deployment of the model, planning the launch event, creating *The Journey of Aya* documentary, hosting the model and coordinating outreach efforts.

 Martin Kon

 Kim Moir


 Sudip Roy

 Chris Taeyoung Kim

 Dave Kong

 Luísa Moura

 Sebastian Ruder

 Yi Chern Tan

 Julia Kreutzer

 Alyssa Pothier

 Astrid Sandoval

 Ahmet Üstün

 Kyle Lastovica

 Brittawnya Prince

 Shubham Shukla

 Jaron Waldman

 Tali Livni

 Daniel Quainoo

 Madeline Smith

 Donglu Wang

 Marina Machado

 Jess Rosenthal

 Trish Starostina

 Lauren Waters

 Abigail Mackenzie-Armes

 Kate Svetlakova

 Ivan Zhang

Safety Evaluation

Our multilingual human evaluation annotators help us understand model quality across languages. They support our evaluations of where models differ and uncover safety and quality issues.

Faraaz Ahmed

Bruno Guratti

Arishi Maisara

Alizé Qureshi

April Alcantara

Maryam Helmy

Brenda Malacara

Manuela Ramirez Naranjo

Kirill Borisov

Ricardo Joaquin Hornedo
Aldeco

Annika Maldonado

Boris Sehovac

Owen Chung

Nishi Jain

Simar Malhan

Ankit Sharma

Laura De Vuono

Milica Jez

Jullia Naag

Hana Sherafati Zanganeh

Sama Elhansi

Dina Kliuchareva

Sasha O'Marra

Ambuj Upadhyay

Sonja Gavric

Finlay Korol-O'Dwyer

Uros Popic

Susheela Willis

Marwan Genena

Rachel Lo

Naeesha Puri

Linda Yanes

Robin Gershman

Juan Lozano

Elina Qureshi

Joanna Yulo

Stuti Govil

Partner Organizations

These organizations supported Aya by hosting events, providing resources, and/or spreading awareness of the project, thereby facilitating contributions and boosting language inclusion efforts.



Universiti Malaysia Sarawak
Faculty of Computer Science and
Information Technology



Google Developer Student Clubs
Thapar Institute of Engineering and
Technology, Patiala, under the leadership
of Siya Sindhani

Linguistics Circle
Nigeria



GalsenAI



Google Developer Student Clubs
P P Savani University, Surat

Google Developer Student Club
P P Savani University, Surat, Gujarat



Rotaract Club
University of Moratuwa, Sri Lanka, led by
Nawoda Thathsarani, Jalina Hirushan and
Chamod Perera



SIMAD iLab



KG College of Arts and Science

KG College of Arts and Science
Coimbatore

































Tensorflow
User Group Surat, Gujarat
















Aya Expanse







Language Ambassadors 1/2

For Aya Expanse, an additional set of Language Ambassadors supported in testing the model across their languages and raising awareness of the model across their communities.

 Mehmet Emre Akbulut
 Samer Attrah
 MUHDIN AWOL
 Kenza Benkirane
 Mann Bhanushali
  Isabella Bicalho Frazeto
 Danylo Boiko
 Sabrina Boumaiza
 Samuel Cahyawijaya
 Emirhan Çelik
 Ryan Chan
 Aurélien-Morgan CLAUDON
  Urszula Czerwinska

 Joana da Matta
 Nguyễn Đạt
 Manasvi Dawane
 Akanksha Devkar
 Sharad Duwal
 Abdeljalil EL MAJJODI
 Shafagh Fadaei
 Neil Fernandes
 Silvia Fernandez
  Hamidreza Ghader
 Manuel Goulão
 Bassam Gouti
 María Grandury
 Miguel Guerrero















 Siddhesh Gunjal
 Mohammed Hamdy
 Hafedh Hichri
  Nhu Hoang Anh Quynh
 Kyle Howard
 Jiwung Hyun
 Joseph Marvin Imperial
 Burin Intachuen
 Ryan Junejo
 Juan Junqueras
  Karthik Reddy Kanjula
 Albert Kao
 Morteza Kashani

















 Ahmed Khaled
 Niharika Khanna
  Dipika Khullar
 Christopher Klamm
 Nazar Kohut
  Alkis Koudounas
 Diana Kozachek
 Katrina Lawrence
 James León
  Jiazheng Li
 Nicolò Loddo
 Dante, Fu On Lok
  Iro Malta
 Harras Mansoor
 Bhavnick Minhas















Aya Expanse














Language Ambassadors 2/2

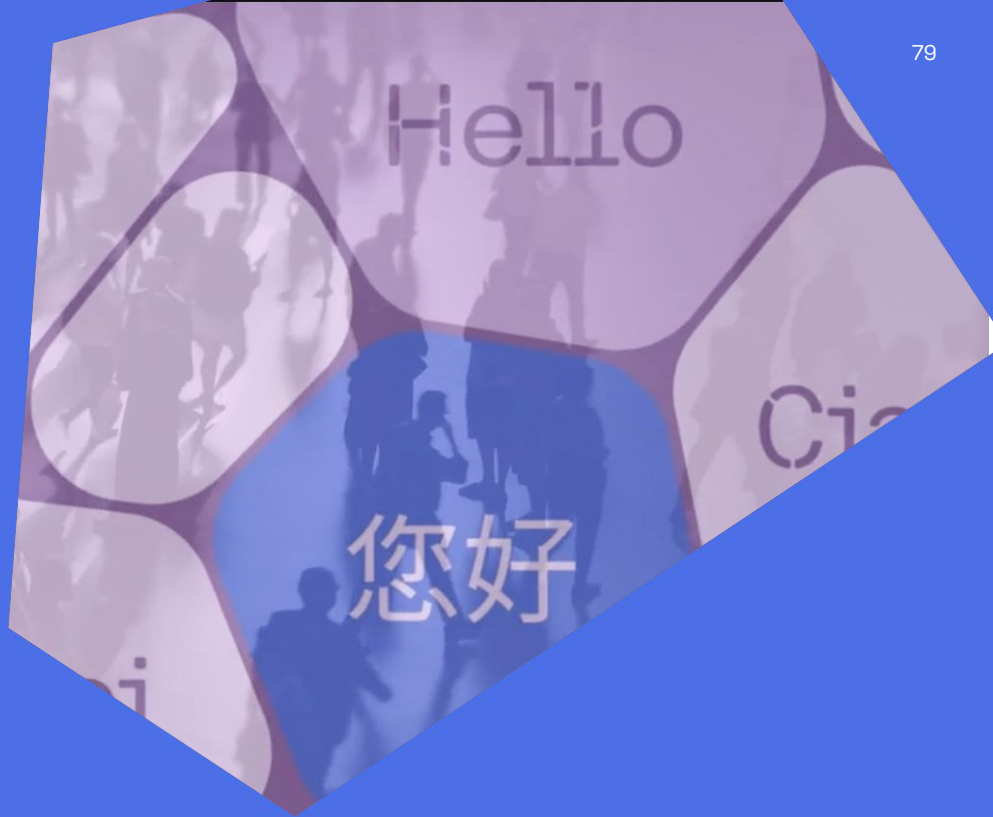
For Aya Expanse, an additional set of Language Ambassadors supported in testing the model across their languages and raising awareness of the model across their communities.

Shachar Mirkin
 Roa'a Mohammad
 Yiyang Nan
  Sree Harsha Nelaturu
 Jekaterina Novikova
 Roni Obaid
 Olympiah Otieno
 Enes Özgözler
 Yavuz Alp Sencer Ozturk
 Carlos Patiño
 Jebish Purbey
 Maria Quijano Jesurum
 Swati Rajwal
 Didi Ramsaran Chin

 Divyaraj Rana
 Aditya Retnanto
 Rodrigo Ribeiro Gomes
 Esra'a Saleh
  Roshan Santhosh
  Drishti Sharma
 Kinza Sheikh
  Aditya Shrivastava
 Vivek Silimkhan
 Marjana Skenduli
 Soham Sonar
 Gürkan Soykan
 David Styveen
 Anthony Susevski

 Adrian Szymczak
 Joanne Tan
 Quentin Tardif
  Ameer Taylor
 Yiorgos Tsalikidis
 Roman Tymtsiv
 Muhammad Saad Uddin
 Louis Ulmer
 Sundar Sripada V. S.
 Freddie Vargus
 Vlad Vasilescu
 Karan Verma
 Henry Vo

  Minh Chien Vu
 Hieu Vu
 Azmine Tushik Wasi
 Warren Williams
 Joseph Wilson
 Gusti Winata
 Ege Yakut
 Eray Yapağcı
 Taha Yassine
 Serhan YILMAZ
 Hanna Yukhymenko
 Mike Zhang



05

Responsibility

Safety for All Languages

The model may produce undesirable responses, such as toxic, biased, or harmful responses - but we want to ensure a safe and responsible use - across all languages.

Previous safety mitigations have predominantly focused on English, which can lead to safety oversights in other languages. This means models might produce safe outputs in English but unsafe ones when prompted in different languages.

With Aya, we focus on a wide, multilingual evaluation of biases, toxicity, and harmfulness, and we implement a multilingual safety measure to prevent misuse for potentially harmful user intentions.



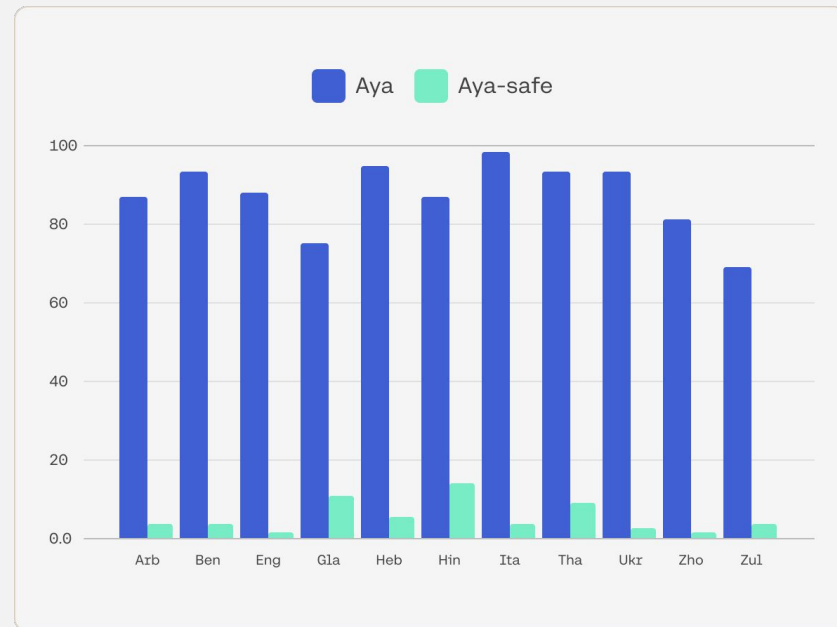


Multilingual Safety Context Distillation

First we define a set of unsafe contexts, where a user queries the model with an adversarial prompt and a harmful intention. We can then train the Aya Model to generate refusal messages for such use cases across all of its languages.

The refusal messages are obtained by querying a teacher model with a safety preamble that explicitly discourages harmful responses. By training on these responses, we distill concepts of safety into the Aya Model, achieving *more harmless responses*, and *maintaining open-ended generation quality*.

NOTE: The release of the Aya model will make community-based red-teaming efforts possible by exposing an open-source multilingual model for community research.

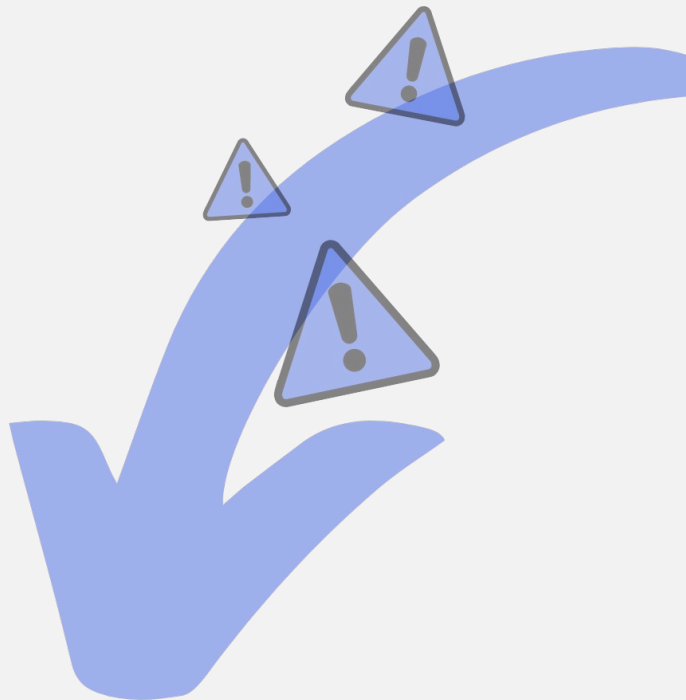


Measuring Toxicity and Bias

Benchmarking toxicity and bias in models helps us understand how often and how seriously the model might give responses that could be toxic or biased across languages.

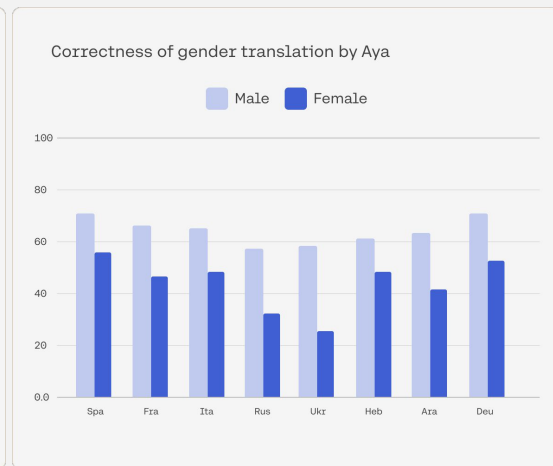
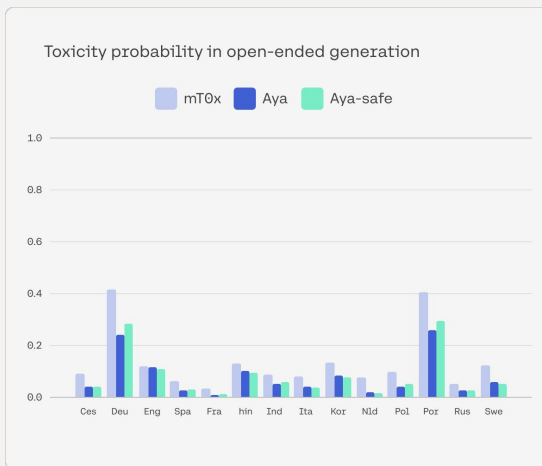
The Aya Model is tested on two evaluation scenarios:

- 1) Toxicity and bias in open-ended generation, across 14 languages.
- 2) Gender bias in machine translation, across 8 languages.



Results From Benchmarking Toxicity and Bias

1. Our findings show that instruction fine-tuning and safety mitigation reduce toxicity and bias.
2. Absolute tendencies towards toxic and bias outputs vary across languages.
3. The problem is not solved: especially racial and gender biases are still present.

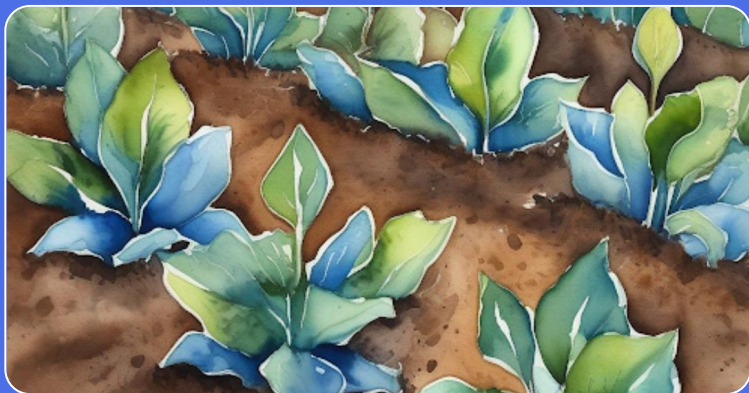


06

The Aya Movement



Read the Research

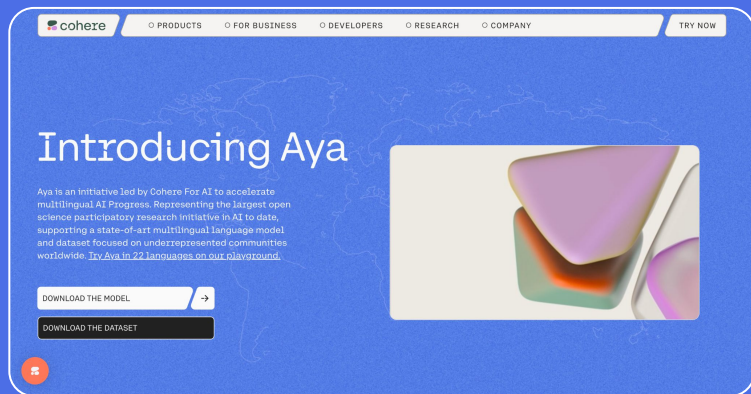


Read our research, [Aya Dataset: An Open-Access Collection for Multilingual Instruction Tuning.](#)



Read our research, [Aya Model: An Instruction Finetuned Open-Access Multilingual Language Model.](#)

Learn more



Visit the [Aya webpage](#) to download the model and dataset, see the latest Aya press coverage, and get to know some of our collaborators.



Read our [blog post](#) on Aya 101's release and on [Aya Expanse](#).

Dive Deeper



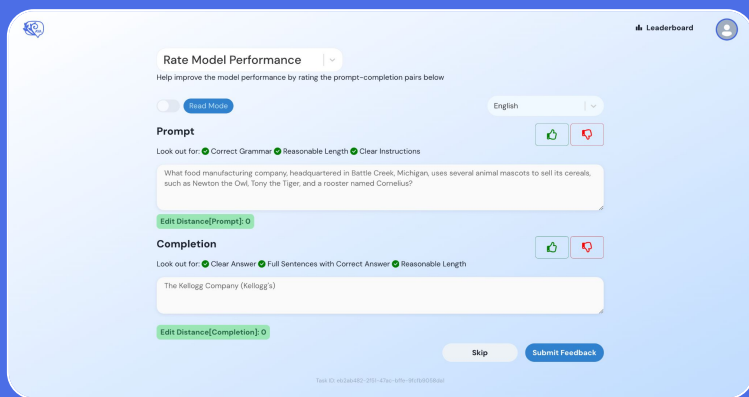
Watch [*The Journey of Aya*](#), a 20-minute documentary featuring many of our collaborators that highlights the importance of progress in multilingual ML, and showcases how this major research effort came together over the past year.



Use your own prompts to [Try Aya on the Cohere Playground](#) in 22 sample languages, or try Aya Expanse on [Hugging Face Spaces](#).

Join us

This is only the beginning. Aya will be a foundation for additional open science projects and we expect to continue to improve Aya capabilities.



Contribute to Aya. Share expertise in your language to be include. We will continue to release data every year or each time an additional 20,000 annotations are contributed (whichever comes first).



Join Aya Community - a space for ML researchers worldwide to connect, learn from one another, and work collaboratively to advance the field of ML research. We will continue to host open science initiatives.



cohere.com/research/aya



[@CohereForAI](https://twitter.com/CohereForAI)



[/showcase/cohere-for-ai](https://www.linkedin.com/showcase/cohere-for-ai)