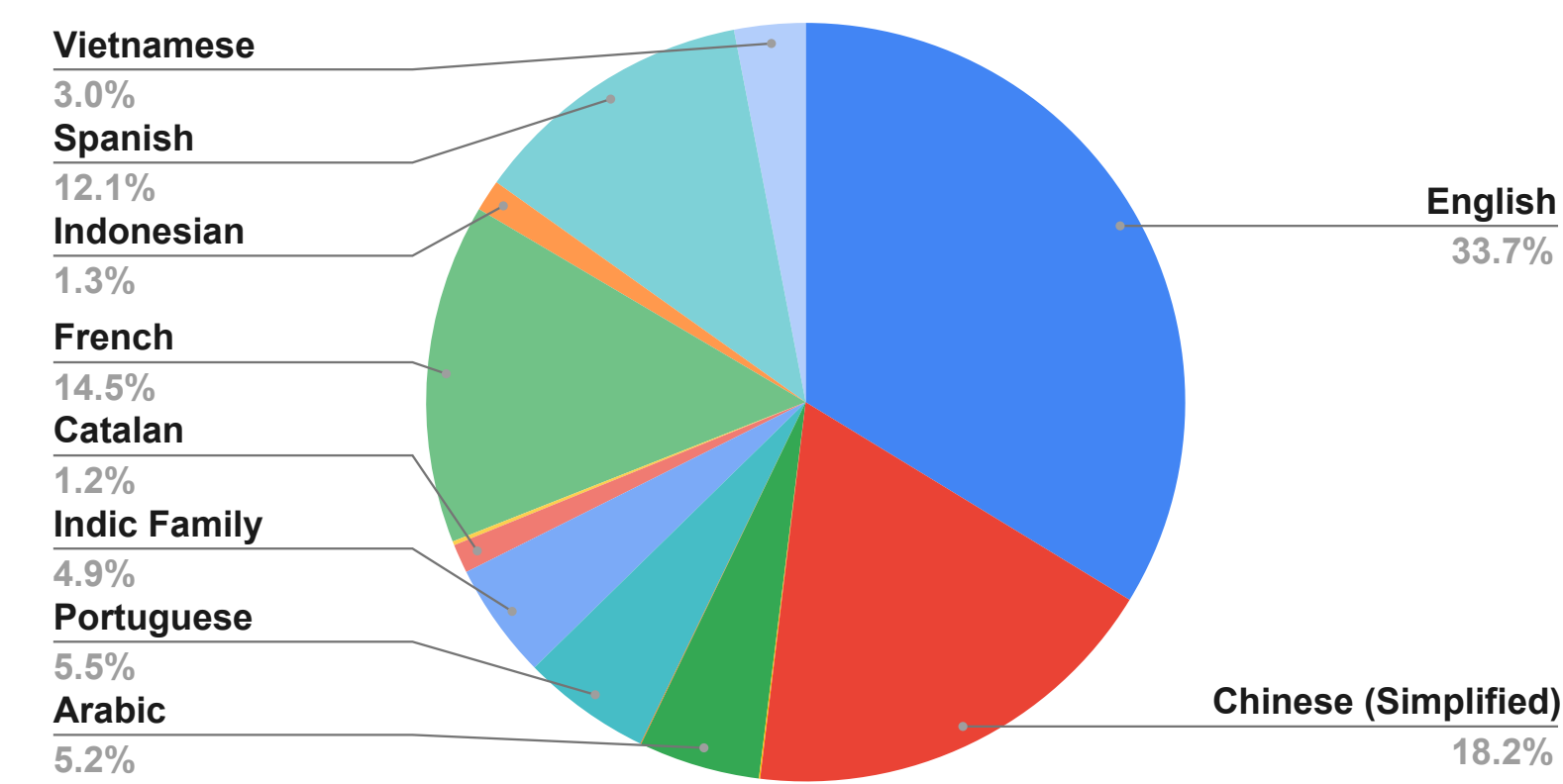


## Background

**Goal:** A LLM should be able to *understand* and *generate* a text in multiple natural languages.

- **Understanding Capability:** when the instructions for LLMs are expressed in different languages, LLMs should understand these instructions and generate a correct output.
- **Generating Capability:** LLMs should be able to generate the correct response in the target language and perform consistently well on (almost) all languages when a fixed language (e.g., English) is used as the instruction language.

### Language Distribution in BLOOM



### Multilingual capability of LLMs on XQuAD

LLM	Language of inputs and outputs						
	En	Zh	Vi	Tr	Ar	El	Hi
<i>Understanding capability (Instruction identical to inputs)</i>							
ChatGPT	<b>56.0</b>	20.5	26.8	18.3	24.1	17.7	0.6
LLaMA	<b>76.6</b>	27.2	36.6	27.8	11.8	22.3	14.3
BLOOM	<b>83.9</b>	83.0	79.9	27.4	79.2	22.8	82.7
<i>Generation capability (Instructions in English)</i>							
ChatGPT	<b>56.0</b>	37.1	36.1	34.5	32.0	29.7	17.5
LLaMA	<b>76.6</b>	66.3	42.9	38.1	24.2	40.7	30.8
BLOOM	<b>83.9</b>	81.8	79.2	27.6	77.2	49.2	80.8

## Key Contributions

- We construct two datasets, one of which contains a multilingual instructions in 100 languages, and the other one contains cross-lingual human preferences in 30 languages.
- We evaluate the multilingual capabilities of LLMs in two dimensions: understanding and generating capability. Unlike previous studies that assess these capabilities in isolation, we urge the community to consider both capabilities when evaluating the multilingual performance of LLMs.
- We scale the multilingual capabilities of LLMs to perform well across 100 languages.



Paper



Code

## Datasets

### Two datasets we constructed

- Multilingual Instruction Dataset:
  - Select 100 languages which covers in FLORES-101 dataset.
  - Instructions and Responses are obtained either in a translated way or generated from chatGPT.
- Cross-Lingual Human Feedback Dataset:
  - Instructions and outputs are in different languages;
  - Cover 30 languages, which simulate up to  $30 \times 29 = 870$  generation scenarios.

### Ablation Study

- Monolingual vs. Cross-Lingual Feedback

	Low		High	
	mono	cross	mono	cross
PAWS-X	-	-	58.43	<b>61.94</b>
XCOPA	47.26	<b>49.71</b>	52.15	<b>54.68</b>
Self-Instruct*	3.25	<b>9.16</b>	12.14	<b>14.71</b>
XL-Sum	3.38	<b>9.99</b>	12.52	<b>16.61</b>
FLORES(f)	0.85	<b>2.89</b>	4.57	<b>9.07</b>
FLORES(t)	1.55	<b>5.64</b>	8.45	<b>16.98</b>

### Response Quality

- Evalution on BLEU and Content Preservation

High			Low		
	BLEU	CP		BLEU	CP
Arabic	73.16	0.82	Armenian	47.16	0.64
Chinese	80.27	0.91	Gujarati	39.68	0.55
French	77.71	0.85	Kannada	41.72	0.57
German	75.50	0.84	Malayalam	45.24	0.62
Hindi	73.26	0.81	Marathi	41.37	0.56
<b>Avg.</b>	75.98	0.85	<b>Avg.</b>	43.03	0.59

### Cross-Lingual Human Feedback Dataset

#### Traditional Human Feedback

Given two texts, modify Text 1 to make it more concise and follow the structure of Text 2.

**Text 1:** I work for a small company with less than ten people. I was able to learn a lot on the job and I have obtained the skills I need to be a successful professional.  
**Text 2:** My current job is in a very small team. I have gained a very good set of skills that will help me move forward in my career.

I worked on a small team, less than 10 people, and gained valuable skills to be successful professionally.

I'm very lucky to have found a job with a micro team of less than 10 people. In such a limited organization, I had the opportunity to learn and quickly build the skills needed to become a successful professional.

#### Cross-Lingual Human Feedback

Given two texts, modify Text 1 to make it more concise and follow the structure of Text 2. Please return the results in Chinese.

**Text 1:** I work for a small company with less than ten people. I was able to learn a lot on the job and I have obtained the skills I need to be a successful professional.  
**Text 2:** My current job is in a very small team. I have gained a very good set of skills that will help me move forward in my career.

我在一个小团队工作，不到10人，获得了有价值的技能，可以在职业上获得成功。

我非常庆幸能找到一份只有10人不到的微型团队的工作。在这么有限的组织中，我有机会学习并且快速建立成为成功专业人士所需的技能。

## Experiments & Results

### Datasets and Tasks

- **Understanding Task:** PAWS-X
- **Generation Task:** FLORES-100 and XL-Sum
- **Reasoning Task:** XCOPA
- **Expert-written Task:** Self-Instruct\* (translation version)

### Baselines

- **Off-the-shelf LLMs:** LLaMA-2 and BLOOM
- **Publicly avaiable multilingual instruction-tuned models:**  $BX_{LLaMA}$  and  $BX_{BLOOM}$
- **Supervised Fine-Tuning (SFT) models:**  $SFT_{LLaMA}$  and  $SFT_{BLOOM}$

### Main Results

Understanding Capabilities												
	PAWS-X	XCOPA		Self-Instruct*		XL-Sum			FLORES(f)		FLORES(t)	
		low	high	low	high	low	mid	high	low	high	low	high
LLaMA	38.10	47.44	47.22	7.09	12.57	4.07	5.44	2.84	3.07	4.95	2.96	6.61
$BX_{LLaMA}$	37.28	49.53	49.00	6.31	11.88	2.17	5.52	7.89	2.69	2.38	3.15	5.31
$SFT_{LLaMA}$	42.32	50.19	49.86	7.32	12.72	4.70	7.34	7.55	3.13	3.93	3.16	6.92
xLLMs-100	<b>46.95</b>	<b>51.53</b>	<b>51.96</b>	<b>12.94</b>	<b>15.35</b>	<b>8.83</b>	<b>13.90</b>	<b>17.29</b>	<b>3.27</b>	<b>8.09</b>	<b>4.04</b>	<b>14.18</b>
BLOOM	36.47	44.27	49.14	7.56	8.67	9.03	14.06	16.80	2.54	2.04	2.05	2.56
$BX_{BLOOM}$	36.42	46.28	50.35	4.81	8.11	4.89	8.47	11.71	2.14	1.74	2.41	1.57
$SFT_{BLOOM}$	36.67	49.42	52.31	6.31	11.88	5.62	10.12	14.33	<b>3.12</b>	3.79	2.62	2.52
xLLMs-100	<b>39.83</b>	<b>52.50</b>	<b>55.59</b>	<b>7.94</b>	<b>13.35</b>	<b>12.87</b>	<b>15.23</b>	<b>18.38</b>	3.02	<b>4.71</b>	<b>3.94</b>	<b>6.54</b>
Generating Capabilities												
	PAWS-X	XCOPA		Self-Instruct*		XL-Sum			FLORES(f)		FLORES(t)	
		low	high	low	high	low	mid	high	low	high	low	high
LLaMA	50.22	49.33	51.52	5.38	8.81	6.26	5.80	8.08	1.35	3.90	2.11	4.95
$BX_{LLaMA}$	48.41	48.00	49.85	7.01	9.80	1.11	2.74	1.70	1.56	5.33	1.37	1.61
$SFT_{LLaMA}$	50.36	48.93	50.05	7.10	12.15	4.51	6.06	9.21	2.42	4.56	2.71	7.29
xLLMs-100	<b>61.94</b>	<b>49.71</b>	<b>54.68</b>	<b>9.16</b>	<b>14.71</b>	<b>9.99</b>	<b>13.57</b>	<b>16.61</b>	<b>2.89</b>	<b>9.07</b>	<b>5.64</b>	<b>16.98</b>
BLOOM	47.39	49.85	49.47	4.07	7.01	6.08	7.77	8.91	0.78	1.20	0.99	1.49
$BX_{BLOOM}$	47.26	47.72	49.98	5.88	8.21	1.98	3.59	4.58	0.47	0.82	1.95	2.33
$SFT_{BLOOM}$	48.50	49.13	49.28	7.78	11.51	3.89	8.87	10.89	2.59	3.12	2.05	2.56
xLLMs-100	<b>50.53</b>	<b>52.36</b>	<b>52.26</b>	<b>10.17</b>	<b>13.62</b>	<b>8.77</b>	<b>11.74</b>	<b>12.36</b>	<b>3.97</b>	<b>5.79</b>	<b>4.22</b>	<b>7.68</b>

### Different Dataset for Multilingual Tuning

	Low		High	
	para	instruct	para	instruct
PAWS-X	-	-	40.17	<b>50.36</b>
XCOPA	37.14	<b>48.93</b>	42.13	<b>50.05</b>
Self-Instruct*	2.63	<b>7.10</b>	5.48	<b>12.15</b>
XL-Sum	1.10	<b>4.51</b>	5.12	<b>9.21</b>
FLORES(f)	<b>5.06</b>	2.42	<b>13.27</b>	4.56
FLORES(t)	<b>12.36</b>	2.71	<b>18.27</b>	7.29

### Off-Target Analysis

	FLORES(f)		FLORES(t)	
	Low	High	Low	High
LLaMA	23.26	16.76	14.15	10.16
$BX_{LLaMA}$	14.13	8.32	12.17	8.24
$SFT_{LLaMA}$	10.26	6.34	8.72	6.23
xLLMs-100	<b>8.82</b>	<b>3.47</b>	<b>6.95</b>	<b>1.46</b>

### Language Democratization

	LLaMA	$BX_{LLaMA}$	$SFT_{LLaMA}$	xLLMs-100
PAWS-X	60.56	58.77	60.63	<b>66.43</b>
XCOPA	93.33	98.52	<b>99.31</b>	89.63
Self-Instruct*	57.85	68.68	62.63	<b>73.92</b>
XL-Sum	47.09	8.90	50.35	<b>67.21</b>
FLORES(f)	34.33	34.00	25.84	<b>34.68</b>
FLORES(t)	49.84	<b>58.28</b>	35.53	48.28