

# Tuning LLMs with Contrastive Alignment Instructions for Machine Translation in Unseen, Low-resource Languages

Anonymous ACL submission

## Abstract

This article introduces contrastive alignment instructions (**AlignInstruct**) to address two challenges in machine translation (MT) on large language models (LLMs). One is the expansion of supported languages to previously unseen ones. The second relates to the lack of data in low-resource languages. Model fine-tuning through MT instructions (**MTInstruct**) is a straightforward approach to the first challenge. However, MTInstruct is limited by weak cross-lingual signals inherent in the second challenge. AlignInstruct emphasizes cross-lingual supervision via a cross-lingual discriminator built using statistical word alignments. Our results based on fine-tuning the BLOOMZ models (1b1, 3b, and 7b1) in up to 24 unseen languages showed that: (1) LLMs can effectively translate unseen languages using MTInstruct; (2) AlignInstruct led to consistent improvements in translation quality across 48 translation directions involving English; (3) Discriminator-based instructions outperformed their generative counterparts as cross-lingual instructions; (4) AlignInstruct improved performance in 30 zero-shot directions.

## 1 Introduction

Large language models (LLMs) (Brown et al., 2020; Chowdhery et al., 2022; Scao et al., 2022; Touvron et al., 2023a; Muennighoff et al., 2023; OpenAI, 2023; Anil et al., 2023; Touvron et al., 2023b) achieved good performance for a wide range of NLP tasks for prevalent languages. However, insufficient coverage for low-resource languages remains to be one significant limitation. Low-resource languages are either not present, or orders of magnitude smaller in size than dominant languages in the pre-training dataset. This limitation is in part due to the prohibitive cost incurred by curating good quality and adequately sized datasets for pre-training. Incrementally adapting existing multilingual LLMs to incorporate an

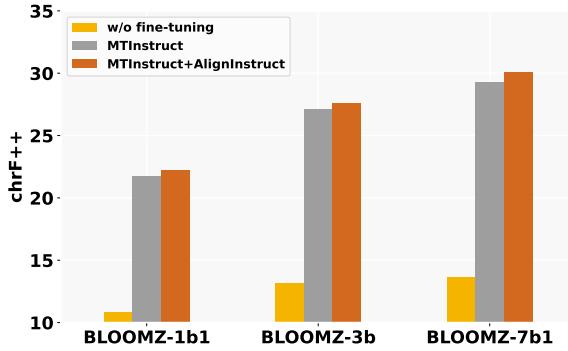


Figure 1: Average chrF++ scores of BLOOMZ models across 24 unseen languages, comparing settings of without fine-tuning, fine-tuning with MTInstruct, and fine-tuning that combines MTInstruct and AlignInstruct.

unseen, low-resource language thus becomes a cost-effective priority to address this limitation. Previous study (de la Rosa and Fernández, 2022; Müller and Laurent, 2022; Yong et al., 2023) explored extending language support using either continual pre-training (Neubig and Hu, 2018; Artetxe et al., 2020; Muller et al., 2021; Ebrahimi and Kann, 2021), or parameter efficient fine-tuning (PEFT) methods (Pfeiffer et al., 2020; Hu et al., 2022; Liu et al., 2022) on monolingual tasks. Extending language support for cross-lingual tasks remains underexplored due to the challenge of incrementally inducing cross-lingual understanding and generation abilities in LLMs (Yong et al., 2023).

This study focused on machine translation (MT) to highlight the cross-lingual LLM adaptation challenge. The challenge lies in enabling translation for low-resource languages that often lack robust cross-lingual signals. We first explored the efficacy of fine-tuning LLMs with MT instructions (MTInstruct) in unseen, low-resource languages. MTInstruct is a method previously shown to bolster the translation proficiency of LLMs for supported languages (Li et al., 2023). Subsequently, given that cross-lingual alignments are suboptimal

mal in LLMs as a result of data scarcity of low-resource languages, we proposed contrastive alignment instructions (AlignInstruct) to explicitly provide cross-lingual supervision during MT fine-tuning. AlignInstruct is a cross-lingual discriminator formulated using statistical word alignments. Our approach was inspired by prior studies (Lambert et al., 2012; Ren et al., 2019; Lin et al., 2020; Mao et al., 2022), which indicated the utility of word alignments in enhancing MT. In addition to AlignInstruct, we discussed two word-level cross-lingual instruction alternatives cast as generative tasks, for comparison with AlignInstruct.

Our experiments fine-tuned the BLOOMZ models (Muennighoff et al., 2023) of varying sizes (1b1, 3b, and 7b1) for 24 unseen, low-resource languages, and evaluated translation on OPUS-100 (Zhang et al., 2020) and Flores-200 (Costajussà et al., 2022). We first showed that MTInstruct effectively induced the translation capabilities of LLMs for these languages. Building on the MTInstruct baseline, the multi-task learning combining AlignInstruct and MTInstruct resulted in stronger translation performance without the need for additional training corpora. The performance improved with larger BLOOMZ models, as illustrated in Fig. 1, indicating that AlignInstruct is particularly beneficial for larger LLMs during MT fine-tuning. When compared with the generative variants of AlignInstruct, our results indicated that discriminator-style instructions better complemented MTInstruct. Furthermore, merging AlignInstruct with its generative counterparts did not further improve translation quality, underscoring the efficacy and sufficiency of AlignInstruct in leveraging word alignments for MT.

In zero-shot translation evaluations on the OPUS benchmark, AlignInstruct exhibited improvements over the MTInstruct baseline in 30 zero-shot directions not involving English, when exclusively fine-tuned with three unseen languages (German, Dutch, and Russian). However, when the fine-tuning data incorporated supported languages (Arabic, French, and Chinese), the benefits of AlignInstruct were only evident in zero-shot translations where the target language was a supported language. In addition, to interpret the inherent modifications within the BLOOMZ models after applying MTInstruct or AlignInstruct, we conducted a visualization of the layer-wise cross-lingual alignment capabilities of the model representations.

## 2 Methodology

This section presents MTInstruct as the baseline, and AlignInstruct. The MTInstruct baseline involved fine-tuning LLMs using MT instructions. AlignInstruct dealt with the lack of cross-lingual signals stemming from the limited parallel training data in low-resource languages. The expectation was enhanced cross-lingual supervision cast as a discriminative task without extra training corpora. Following this, we introduced two generative variants of AlignInstruct for comparison.<sup>1</sup>

### 2.1 Baseline: MTInstruct

Instruction tuning (Wang et al., 2022; Mishra et al., 2022; Chung et al., 2022; Ouyang et al., 2022; Sanh et al., 2022; Wei et al., 2022) has been shown to generalize LLMs’ ability to perform various downstream tasks, including MT (Li et al., 2023).

Given a pair of the parallel sentences,  $((x_i)_1^N, (y_j)_1^M)$ , where  $(x_i)_1^N := x_1 x_2 \dots x_N$ ,  $(y_j)_1^M := y_1 y_2 \dots y_M$ .  $x_i, y_j \in \mathcal{V}$  are members of the vocabulary  $\mathcal{V}$  containing unique tokens that accommodate languages  $X$  and  $Y$ . Li et al. (2023) showed that the following MT instructions (MTInstruct) can improve the translation ability in an LLM with a limited number of parallel sentences:

- **Input:** “Translate from  $Y$  to  $X$ .  
 $Y: y_1 y_2 \dots y_M$ .  
 $X:$ ”
- **Output:** “ $x_1 x_2 \dots x_N$ .”

Note that Li et al. (2023) demonstrated the utility of MTInstruct solely within the context of fine-tuning for languages acquired at pre-training phase. This study called for an assessment of MTInstruct on its efficacy for adapting to previously unsupported languages, denoted as  $X$ , accompanied by the parallel data in a supported language  $Y$ .

### 2.2 AlignInstruct

Word alignments have been demonstrated to enhance MT performance (Lambert et al., 2012; Ren et al., 2019; Lin et al., 2020; Mao et al., 2022), both in the fields of statistical machine translation (SMT) (Brown et al., 1993) and neural machine translation (NMT) (Sutskever et al., 2014; Bahdanau et al., 2015). Ren et al. (2019) and Mao et al. (2022) reported the utility of SMT-derived

<sup>1</sup>We also discussed monolingual instructions for MT fine-tuning in App. F.

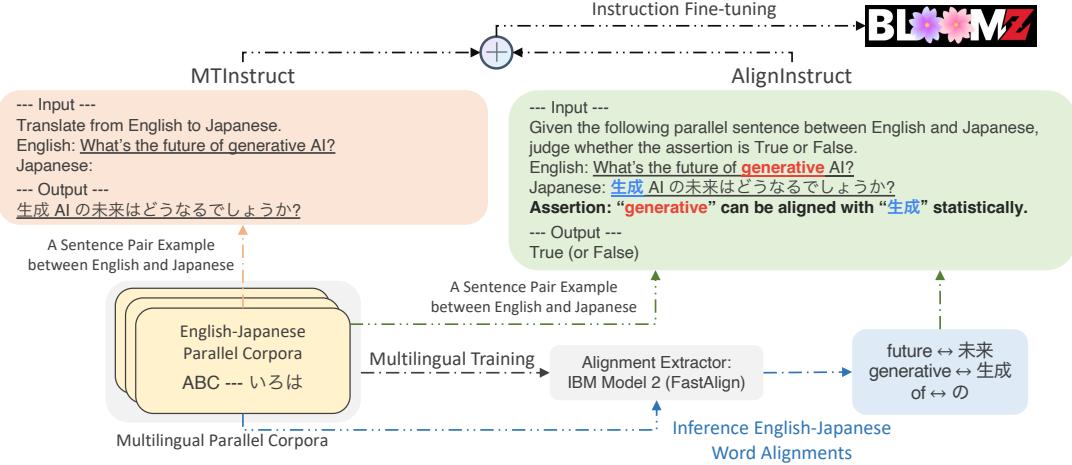


Figure 2: **Proposed instruction tuning methods combining MTInstruct (Sec. 2.1) and AlignInstruct (Sec. 2.2) for LLMs in MT tasks.**  $\oplus$  denotes combining multiple instruction patterns with a specific fine-tuning curriculum (Sec. 3.2). IBM Model 2 indicates word alignment model of statistical machine translation (Brown et al., 1993).

contrastive word alignments in guiding encoder-decoder NMT model training. Built upon their findings, we introduced AlignInstruct for bolstering cross-lingual alignments in LLMs. We expected AlignInstruct to enhance translation performance particularly for languages with no pre-training data and limited fine-tuning data.

As shown in Fig. 2, we employed FastAlign (Dyer et al., 2013) to extract statistical word alignments from parallel corpora. Our approach depended on a trained FastAlign model (IBM Model 2, Brown et al., 1993) to ensure the quality of the extracted word pairs. These high-quality word alignment pairs were regarded as “gold” word pairs for constructing AlignInstruct instructions.<sup>2</sup> Assuming one gold word pair  $(x_k x_{k+1}, y_l y_{l+1} y_{l+2})$  was provided for the sentence pair  $((x_i)_1^N, (y_j)_1^M)$ , the AlignInstruct instruction reads:

- **Input:** “Given the following parallel sentence between  $Y$  and  $X$ , judge whether the assertion is True or False.  
 $Y: y_1 y_2 \dots y_M$ .  
 $X: x_1 x_2 \dots x_N$ .  
Assertion: “ $y_l y_{l+1} y_{l+2}$ ” can be aligned with “ $x_k x_{k+1}$ ” statistically.”
- **Output:** “True” (or “False”)

Instructions with the “False” output were constructed by uniformly swapping out part of the word pair to create misalignment. We anticipated

<sup>2</sup>Note that these word pairs may not necessarily represent direct translations of each other; instead, they are word pairs identified based on their co-occurrence probability within the similar context. Refer to IBM model 2 in SMT.

that this treatment forced the model to learn to infer the output by recognizing true alignment-enriched instructions. This would require the model to encode word-level cross-lingual representation, a crucial characteristic for MT tasks.

### 2.3 Generative Counterparts of AlignInstruct

Previous studies (Liang et al., 2022; Yu et al., 2023) have suggested the importance of both discriminative and generative tasks in fine-tuning LLMs. We accordingly considered two generative variants of AlignInstruct. We then compared them with AlignInstruct to determine the most effective training task. As detailed in Sec. 4, our results indicated that these variants underperformed AlignInstruct when applied to unseen, low-resource languages.

#### 2.3.1 HintInstruct

HintInstruct as a generative variant of AlignInstruct was instructions containing word alignment hints. It was inspired by Ghazvininejad et al. (2023), where dictionary hints were shown to improve few-shot in-context learning. Instead of relying on additional dictionaries, we used the same word alignments described in Sec. 2.2, which were motivated by the common unavailability of high-quality dictionaries for unseen, low-resource languages. Let  $\{(x_{k_s} x_{k_s+1} \dots x_{k_s+n_s}, y_{l_s} y_{l_s+1} \dots y_{l_s+m_s})\}_{s=1}^S$  be  $S$  word pairs extracted from the sentence pair  $((x_i)_1^N, (y_j)_1^M)$ . HintInstruct follows the instruction pattern:

- **Input:** “Use the following alignment hints and translate from  $Y$  to  $X$ .

162  
163  
164  
165  
166  
167  
168  
169  
170  
171  
172  
173  
174  
175  
176  
177  
178  
179

180  
181  
182  
183  
184  
185  
186  
187  
188  
189  
190

191  
192  
193  
194  
195  
196  
197  
198  
199  
200  
201  
202  
203  
204  
205  
206  
207  
208  
209  
210  
211  
212  
213  
214  
215  
216  
217  
218  
219  
220  
221  
222

223 Alignments between  $X$  and  $Y$ :  
 224  $-(x_{k_1}x_{k_1+1}\dots x_{k_1+n_1}, y_{l_1}y_{l_1+1}\dots y_{l_1+m_1}),$   
 225  $-(x_{k_2}x_{k_2+1}\dots x_{k_2+n_2}, y_{l_2}y_{l_2+1}\dots y_{l_2+m_2}),$   
 226  $\dots,$   
 227  $-(x_{k_S}x_{k_S+1}\dots x_{k_S+n_S}, y_{l_S}y_{l_S+1}\dots y_{l_S+m_S}),$   
 228  $Y: y_1y_2\dots y_M.$   
 229  $X:$   
 230 • **Output:** “ $x_1x_2\dots x_N$ .”

231 where  $S$  denotes the number of the word alignment  
 232 pairs used to compose the instructions. Different  
 233 from AlignInstruct, HintInstruct expects the trans-  
 234 lation targets to be generated.

### 235 2.3.2 ReviseInstruct

236 ReviseInstruct was inspired by Ren et al. (2019)  
 237 and Liu et al. (2020) for the notion of generating  
 238 parallel words or phrases, thereby encouraging a  
 239 model to encode cross-lingual alignments. A Re-  
 240 viseInstruct instruction contained a partially cor-  
 241 rupted translation target, as well as a directive to  
 242 identify and revise these erroneous tokens. To-  
 243 kens are intentionally corrupted at the granularity  
 244 of individual words, aligning with the word-level  
 245 granularity in AlignInstruct and HintInstruct. Revi-  
 246 seInstruct follows the instruction pattern:<sup>3</sup>

- 247 • **Input:** “Given the following translation of  $X$   
 248 from  $Y$ , output the incorrectly translated word  
 249 and correct it.  
 250  $Y: y_1y_2\dots y_M.$   
 251  $X: x_1x_2\dots x_kx_{k+1}\dots x_{k+n}\dots x_N.$ ”
- 252 • **Output:** “The incorrectly translated  
 253 word is “ $x_kx_{k+1}\dots x_{k+n}$ ”. It should be  
 254 “ $x_jx_{j+1}\dots x_{j+m}$ ”.”

## 255 3 Experimental Settings

### 256 3.1 Backbone Models and Unseen Languages

257 Our experiments fine-tuned the BLOOMZ mod-  
 258 els (Muennighoff et al., 2023) for MT in un-  
 259 seen, low-resource languages. BLOOMZ is an  
 260 instruction fine-tuned multilingual LLM from  
 261 BLOOM (Scao et al., 2022) that supports trans-  
 262 lation across 46 languages. Two lines of experiments  
 263 evaluated the effectiveness of the MTInstruct base-  
 264 line and AlignInstruct:

265 **BLOOMZ+24** Tuning BLOOMZ-7b1, BLOOMZ-  
 266 3b, and BLOOMZ-1b1<sup>4</sup> for 24 unseen, low-  
 267 resource languages. These experiments aimed to:

268 (1) assess the effectiveness of AlignInstruct in mul-  
 269 tilingual, low-resource scenarios; (2) offer compari-  
 270 son across various model sizes. We used the OPUS-  
 271 100 (Zhang et al., 2020)<sup>5</sup> datasets as training data.  
 272 OPUS-100 is an English-centric parallel corpora,  
 273 with around 4.5M parallel sentences in total for 24  
 274 selected languages, averaging 187k sentence pairs  
 275 for each language and English. Refer to App. A  
 276 for training data statistics. We used OPUS-100  
 277 and Flores-200 (Costa-jussà et al., 2022)<sup>6</sup> for eval-  
 278 uating translation between English and 24 unseen  
 279 languages (48 directions in total) on in-domain and  
 280 out-of-domain test sets, respectively. The identical  
 281 prompt as introduced in Sec. 2.1 was employed for  
 282 inference. Inferences using alternative MT prompts  
 283 are discussed in App. G.

284 **BLOOMZ+3** Tuning BLOOMZ-7b1 with three  
 285 unseen languages, German, Dutch, and Russian,  
 286 or a combination of these three unseen languages  
 287 and another three seen (Arabic, French, and Chi-  
 288 nese). We denote the respective setting as **de-nl-ru**  
 289 and **ar-de-fr-nl-ru-zh**. These experiments as-  
 290 sessed the efficacy of AlignInstruct in zero-shot  
 291 translation scenarios, where translation directions  
 292 were not presented during fine-tuning, as well as  
 293 the translation performance when incorporating  
 294 supported languages as either source or target lan-  
 295 guages. To simulate the low-resource fine-tuning  
 296 scenario, we randomly sampled 200k parallel sen-  
 297 tences for each language. For evaluation, we used  
 298 the OPUS-100 supervised and zero-shot test sets,  
 299 comprising 12 supervised directions involving Eng-  
 300 lish and 30 zero-shot directions without English  
 301 among six languages.

302 Notably, BLOOMZ’s pre-training data includes  
 303 the English portion of the Flores-200 dataset, po-  
 304 tentially leading to data leakage during evalua-  
 305 tion (Muennighoff et al., 2023; Zhu et al., 2023a).  
 306 To mitigate this, our evaluation also compared  
 307 translation quality before and after fine-tuning,  
 308 thereby distinguishing the genuine improvements  
 309 in translation capability attributable to the fine-  
 310 tuning process (refer to the results in Sec. 4).

### 311 3.2 Training Details and Curricula

312 The PEFT method, LoRA (Hu et al., 2022), was  
 313 chosen to satisfy the parameter efficiency require-  
 314 ment for low-resource languages, as full-parameter  
 315 fine-tuning would likely under-specify the mod-

<sup>3</sup>We illustrated examples of HintInstruct and ReviseInstruct in App. E for reference.

<sup>4</sup><https://huggingface.co/bigscience/bloomz>

<sup>5</sup><https://opus.nlpl.eu/opus-100.php>

<sup>6</sup><https://github.com/facebookresearch/flores/blob/main/flores200/README.md>

BLOOMZ model	Objective	OPUS en→xx			OPUS xx→en			Flores en→xx			Flores xx→en		
		BLEU	chrF++	COMET									
	w/o fine-tuning	3.61	8.82	47.81	6.70	18.49	51.68	2.00	9.35	36.54	9.95	24.47	52.05
<i>Individual objectives</i>													
BLOOMZ-7b1	MTInstruct	11.54	25.33	64.54	18.59	33.25	69.84	3.30	17.10	40.58	11.37	27.14	56.33
	AlignInstruct	4.73	9.23	49.85	5.32	12.90	53.26	1.97	8.90	42.35	3.47	11.93	39.58
	<i>Multiple objectives with different curricula</i>												
	MT+Align	<b>12.28</b>	<b>26.17</b>	<b>65.54</b>	<b>18.72</b>	<b>34.02</b>	<b>70.69</b>	3.26	<b>17.20</b>	<b>41.07</b>	<b>11.60</b>	<b>27.38</b>	<b>56.98</b>
	Align→MT	<b>11.73</b>	<b>25.48</b>	64.54	17.54	32.62	69.76	<b>3.35</b>	<b>17.21</b>	<b>40.85</b>	11.32	<b>27.21</b>	<b>56.50</b>
	MT+Align→MT	<b>12.10</b>	<b>26.16</b>	<b>65.43</b>	18.23	<b>33.54</b>	<b>70.60</b>	3.28	<b>17.26</b>	<b>41.13</b>	<b>11.48</b>	<b>27.34</b>	<b>56.78</b>
BLOOMZ-3b	w/o fine-tuning	4.63	9.93	48.53	5.90	16.38	48.05	2.00	9.09	39.52	5.86	18.56	47.03
	<i>Individual objectives</i>												
	MTInstruct	10.40	23.08	62.28	16.10	31.15	68.36	2.85	16.23	39.21	8.92	24.57	53.33
BLOOMZ-1b1	AlignInstruct	1.70	4.05	43.89	0.87	3.20	41.93	0.16	3.09	31.10	0.10	1.80	29.46
	<i>Multiple objectives with different curricula</i>												
	MT+Align	<b>10.61</b>	<b>23.64</b>	<b>62.84</b>	<b>16.73</b>	<b>31.51</b>	<b>68.52</b>	<b>2.95</b>	<b>16.62</b>	<b>39.83</b>	<b>9.50</b>	<b>25.16</b>	<b>54.35</b>
	Align→MT	10.22	22.53	61.99	15.90	30.31	67.79	<b>3.02</b>	<b>16.43</b>	<b>39.46</b>	<b>9.07</b>	<b>24.70</b>	<b>53.71</b>
	MT+Align→MT	<b>10.60</b>	<b>23.35</b>	<b>62.69</b>	<b>16.58</b>	<b>31.64</b>	<b>68.98</b>	<b>2.93</b>	<b>16.57</b>	<b>39.78</b>	<b>9.41</b>	<b>25.08</b>	<b>54.13</b>
	w/o fine-tuning	3.76	7.57	46.98	4.78	14.11	49.34	1.24	6.93	38.13	3.49	14.56	43.26
BLOOMZ-1b1	<i>Individual objectives</i>												
	MTInstruct	7.42	17.85	57.53	11.99	25.59	63.93	2.11	14.40	36.35	5.33	20.65	48.83
	AlignInstruct	2.51	5.29	45.17	3.13	8.92	48.48	0.35	3.79	31.70	1.35	6.43	33.63
	<i>Multiple objectives with different curricula</i>												
	MT+Align	<b>7.80</b>	<b>18.48</b>	<b>57.77</b>	<b>12.57</b>	<b>25.92</b>	<b>64.03</b>	<b>2.16</b>	<b>14.54</b>	<b>37.05</b>	<b>5.46</b>	<b>20.90</b>	<b>49.31</b>
	Align→MT	<b>7.49</b>	<b>18.09</b>	<b>57.67</b>	11.80	24.70	63.29	2.08	14.28	<b>36.61</b>	5.24	20.53	48.76
	MT+Align→MT	<b>7.98</b>	<b>18.61</b>	<b>57.94</b>	<b>12.43</b>	<b>25.78</b>	63.93	<b>2.16</b>	<b>14.46</b>	<b>37.02</b>	<b>5.37</b>	<b>20.67</b>	<b>49.01</b>

Table 1: **Results of BLOOMZ+24 fine-tuned with MTInstruct and AlignInstruct on different curricula** as described in 3.2. Scores that surpass the MTInstruct baseline are marked in **bold**.

els. See App. B for implementation details. How AlignInstruct and MTInstruct are integrated into training remained undetermined. To that end, we investigated three training curricula:

**Multi-task Fine-tuning** combined multiple tasks in a single training session (Caruana, 1997). This was realized by joining MTInstruct and AlignInstruct training data, denoted as **MT+Align**.<sup>7</sup>

**Pre-fine-tuning & Fine-tuning** arranges fine-tuning in a two-stage curriculum (Bengio et al., 2009), first with AlignInstruct, then with MTInstruct.<sup>8</sup> This configuration, denoted as **Align→MT**, validates whether AlignInstruct should precede MTInstruct.

**Mixed Fine-tuning** (Chu et al., 2017) arranged the above curricula to start with MT+Align, followed by MTInstruct, denoted as **MT+Align→MT**.

## 4 Evaluation and Analysis

This section reports BLEU (Papineni et al., 2002; Post, 2018), chrF++ (Popović, 2015), and COMET (Rei et al., 2020)<sup>9</sup> scores for respective experimental configurations. We further character-

ized of the degree to which intermediate embeddings were language-agnostic after fine-tuning.

### 4.1 BLOOMZ+24 Results

Tab. 1 shows the scores for the unmodified BLOOMZ models, as well as BLOOMZ+24 under MTInstruct, AlignInstruct, and the three distinct curricula. Non-trivial improvements in all metrics were evident for BLOOMZ+24 under MTInstruct. This suggests that MTInstruct can induce translation capabilities in unseen languages. Applying AlignInstruct and MTInstruct via the curricula further showed better scores than the baselines, suggesting the role of AlignInstruct as complementing MTInstruct. Align→MT was an exception, performing similarly to MTInstruct. This may indicate AlignInstruct’s complementarity depends on its cadence relative to MTInstruct in a curriculum.

Superior OPUS and Flores scores under the xx→en direction were evident, compared to the reverse direction, en→xx. This suggests that our treatments induced understanding capabilities more than generative ones. This may be attributed to the fact that BLOOMZ had significant exposure to English, and that we used English-centric corpora. Finally, we noted the inferior performance of Flores than OPUS. This speaks to the challenge of instilling translation abilities in unseen languages for the out-of-domain MT task. Our future work will focus on enhancing the domain generalization

<sup>7</sup>Note that AlignInstruct and MTInstruct were derived from the same parallel corpora.

<sup>8</sup>An effective curriculum often starts with a simple and general task, followed by a task-specific task.

<sup>9</sup>COMET scores do not currently support Limburgish (li), Occitan (oc), Tajik (tg), Turkmen (tk), and Tatar (tt) among the 24 languages in the BLOOMZ+24 setting. Thus, we report the average COMET scores for the remaining 19 languages.

Objective	en-af	af-en	en-am	am-en	en-be	be-en	en-cy	cy-en	en-ga	ga-en	en-gd	gd-en
MTInstruct	25.0	38.5	3.0	3.4	8.9	14.0	20.2	33.2	15.6	29.2	13.1	66.0
MT+Align	25.0	36.9	<b>3.4</b>	<b>4.9</b>	8.3	13.9	<b>20.6</b>	<b>33.8</b>	<b>17.6</b>	<b>32.6</b>	<b>15.6</b>	48.1
Objective	en-gl	gl-en	en-ha	ha-en	en-ka	ka-en	en-kk	kk-en	en-km	km-en	en-ky	ky-en
MTInstruct	16.9	24.7	12.3	10.0	4.6	10.0	12.6	14.6	19.7	13.9	16.0	21.1
MT+Align	17.1	24.4	<b>14.6</b>	<b>11.4</b>	4.9	<b>10.5</b>	12.3	<b>15.6</b>	<b>20.4</b>	<b>14.4</b>	15.8	<b>23.3</b>
Objective	en-li	li-en	en-my	my-en	en-nb	nb-en	en-nn	nn-en	en-oc	oc-en	en-si	si-en
MTInstruct	13.5	21.3	6.2	5.2	12.7	22.2	18.3	27.1	10.0	13.4	5.2	11.5
MT+Align	13.2	<b>22.3</b>	<b>7.6</b>	<b>6.3</b>	<b>13.5</b>	<b>24.2</b>	<b>19.0</b>	<b>28.5</b>	9.1	13.5	5.1	<b>13.9</b>
Objective	en-tg	tg-en	en-tk	tk-en	en-tt	tt-en	en-ug	ug-en	en-uz	uz-en	en-yi	yi-en
MTInstruct	5.5	8.0	24.4	30.4	1.9	3.6	1.2	4.2	3.1	5.7	7.1	14.9
MT+Align	<b>6.6</b>	<b>8.8</b>	<b>27.2</b>	<b>31.2</b>	2.1	<b>5.0</b>	1.1	<b>5.5</b>	<b>3.5</b>	<b>7.4</b>	<b>11.1</b>	12.8

Table 2: Language-wise BLEU results on BLOOMZ-7b1 for BLOOMZ+24 fine-tuned using MTInstruct or MT+Align. Scores significantly (Koehn, 2004) outperforming the MTInstruct baseline are emphasized in **bold** while those decreased significantly (Koehn, 2004) are marked in *italics*.

BLOOMZ model	Objective	OPUS en→xx			OPUS xx→en			Flores en→xx			Flores xx→en		
		BLEU	chrF++	COMET									
BLOOMZ-7b1	MTInstruct	11.54	25.33	64.54	18.59	33.25	69.84	3.30	17.10	40.58	11.37	27.14	56.33
	MT+Align	<b>12.28</b>	<b>26.17</b>	<b>65.54</b>	<b>18.72</b>	<b>34.02</b>	<b>70.69</b>	3.26	<b>17.20</b>	<b>41.07</b>	<b>11.60</b>	<b>27.38</b>	<b>56.98</b>
	MT+Hint	<b>12.12</b>	<b>25.92</b>	<b>64.60</b>	18.25	33.18	<b>70.31</b>	<b>3.34</b>	<b>17.13</b>	<b>41.10</b>	<b>11.45</b>	<b>27.37</b>	<b>56.86</b>
	MT+Revise	<b>11.96</b>	<b>25.73</b>	<b>64.73</b>	<b>18.69</b>	<b>33.74</b>	<b>70.32</b>	<b>3.34</b>	17.10	<b>41.07</b>	<b>11.44</b>	<b>27.37</b>	<b>56.73</b>
BLOOMZ-3b	MTInstruct	10.40	23.08	62.28	16.10	31.15	68.36	2.85	16.23	39.21	8.92	24.57	53.33
	MT+Align	<b>10.61</b>	<b>23.64</b>	<b>62.84</b>	<b>16.73</b>	<b>31.51</b>	<b>68.52</b>	<b>2.95</b>	<b>16.62</b>	<b>39.83</b>	<b>9.50</b>	<b>25.16</b>	<b>54.35</b>
	MT+Hint	<b>10.49</b>	<b>23.34</b>	<b>62.65</b>	<b>16.29</b>	<b>31.43</b>	<b>68.83</b>	<b>3.11</b>	<b>16.95</b>	<b>39.91</b>	<b>9.52</b>	<b>25.25</b>	<b>54.28</b>
	MT+Revise	<b>10.52</b>	23.03	62.04	<b>16.22</b>	30.98	68.28	<b>2.99</b>	<b>16.83</b>	<b>39.52</b>	<b>9.47</b>	<b>25.21</b>	<b>53.91</b>
BLOOMZ-1b1	MTInstruct	7.42	17.85	57.53	11.99	25.59	63.93	2.11	14.40	36.35	5.33	20.65	48.83
	MT+Align	<b>7.80</b>	<b>18.48</b>	<b>57.77</b>	<b>12.57</b>	<b>25.92</b>	<b>64.03</b>	<b>2.16</b>	<b>14.54</b>	<b>37.05</b>	<b>5.46</b>	<b>20.90</b>	<b>49.31</b>
	MT+Hint	<b>7.71</b>	<b>18.15</b>	<b>57.76</b>	11.52	24.88	63.63	<b>2.21</b>	<b>14.61</b>	<b>37.24</b>	<b>5.47</b>	<b>20.78</b>	<b>48.97</b>
	MT+Revise	7.31	<b>17.99</b>	57.45	<b>12.00</b>	25.33	63.81	2.07	14.32	<b>36.68</b>	<b>5.41</b>	<b>20.91</b>	<b>49.09</b>

Table 3: Results of BLOOMZ+24 fine-tuned combining MTInstruct with AlignInstruct (or its generative variants). Scores that surpass the MTInstruct baseline are marked in **bold**.

capabilities of LLM fine-tuning in MT tasks.

Moreover, we reported the language-wise scores in Tab. 2. Specifically, in the “en-xx” direction, 11 languages showed statistically significant (Koehn, 2004) improvements, and only 2 decreased significantly. In the “xx-en” direction, the improvements were more pronounced, with 18 languages improving significantly (most by over 1 BLEU point) and 3 decreasing significantly. The average improvement for “en-xx” was 0.74, which was substantial, especially given the limited volume of parallel data available for each language. The smaller average increase in “xx-en” can be attributed to a large decrease in one language (gd), likely due to limited training data (which can be potentially addressed with oversampling). The significantly enhanced performance in most individual languages underscores the effectiveness of our proposed methods.

#### 4.2 Assessing AlignInstruct Variants

From the results reported in Tab. 3, we observed the objectives with AlignInstruct consistently outperformed those with HintInstruct or ReviseInstruct

across metrics and model sizes. Namely, easy, discriminative instructions, rather than hard, generative ones, may be preferred for experiments under similar data constraints. The low-resource constraint likely made MTInstruct more sensitive to the difficulty of its accompanying tasks.

Further, combining more than two instruction tuning tasks simultaneously did not guarantee consistent improvements, see Tab. 4. Notably, MT+Align either outperformed or matched the performance of other objective configurations. While merging multiple instruction tuning tasks occasionally resulted in superior BLEU and chrF++ scores for OPUS xx→en, it fell short in COMET scores compared to MT+Align. This indicated that while such configurations might enhance word-level translation quality, as reflected by BLEU and chrF++ scores, due to increased exposure to cross-lingual word alignments, MT+Align better captured the context of the source sentence as reflected by COMET scores. Overall, these instruction tuning tasks did not demonstrate significant synergistic effects for fine-tuning for unseen languages.

Objective	OPUS en→xx			OPUS xx→en			Flores en→xx			Flores xx→en		
	BLEU	chrF++	COMET									
MTInstruct	11.54	25.33	64.54	18.59	33.25	69.84	3.30	17.10	40.58	11.37	27.14	56.33
MT+Align	<b>12.28</b>	<b>26.17</b>	<b>65.54</b>	<b>18.72</b>	<b>34.02</b>	<b>70.69</b>	3.26	<b>17.20</b>	<b>41.07</b>	<b>11.60</b>	<b>27.38</b>	<b>56.98</b>
MT+Align+Revise	<b>12.08</b>	<b>25.73</b>	<b>64.55</b>	<b>19.23</b>	<b>34.32</b>	<b>70.60</b>	<b>3.33</b>	<b>17.25</b>	<b>41.17</b>	<b>11.60</b>	<b>27.61</b>	<b>57.22</b>
MT+Align+Hint	<b>12.02</b>	<b>25.51</b>	<b>64.58</b>	<b>19.40</b>	<b>34.44</b>	<b>70.65</b>	3.25	16.87	<b>41.13</b>	<b>11.58</b>	<b>27.48</b>	<b>56.93</b>
MT+Hint+Revise	<b>12.10</b>	<b>25.69</b>	<b>64.68</b>	<b>19.58</b>	<b>34.49</b>	<b>70.55</b>	<b>3.34</b>	<b>17.24</b>	<b>41.13</b>	<b>11.70</b>	<b>27.62</b>	<b>57.19</b>
MT+Align+Hint+Revise	<b>12.00</b>	<b>25.39</b>	<b>64.55</b>	<b>19.68</b>	<b>34.48</b>	<b>70.64</b>	<b>3.40</b>	<b>17.17</b>	<b>41.21</b>	<b>11.67</b>	<b>27.54</b>	<b>57.16</b>

Table 4: **Results of BLOOMZ+24 combining MTInstruct with multiple objectives among AlignInstruct, HintInstruct, and ReviseInstruct on BLOOMZ-7b1.** Scores that surpass MTInstruct are marked in **bold**.

Fine-tuned Languages	Objective	Zero-shot Directions				Supervised Directions			
		Directions	BLEU	chrF++	COMET	Directions	BLEU	chrF++	COMET
-	w/o fine-tuning	overall	6.89	19.14	57.95	en→xx	13.38	26.65	64.28
		seen→seen	16.95	30.78	74.58	xx→en	21.70	42.05	72.72
		seen→unseen	2.30	13.31	49.98	en→seen	20.13	32.87	76.99
		unseen→seen	7.78	20.07	62.74	en→unseen	6.63	20.43	51.56
		unseen→unseen	2.37	14.83	46.06	seen→en	26.30	48.70	78.22
						unseen→en	17.10	35.40	67.23
de-nl-ru	MTInstruct	overall	8.38	22.75	59.93	en→xx	17.05	32.02	69.26
		seen→seen	14.52	27.25	70.48	xx→en	25.13	45.02	76.29
		seen→unseen	6.14	22.82	54.75	en→seen	17.60	29.87	73.81
		unseen→seen	7.56	19.22	61.99	en→unseen	16.50	34.17	64.70
		unseen→unseen	6.85	23.45	54.07	seen→en	25.73	47.07	77.52
	MT+Align	overall	<b>8.86</b>	<b>23.30</b>	<b>60.70</b>	unseen→en	24.53	42.97	75.06
		seen→seen	<b>14.77</b>	<b>27.80</b>	<b>71.07</b>	en→xx	16.63	31.73	68.79
		seen→unseen	<b>6.31</b>	<b>23.08</b>	<b>54.81</b>	xx→en	<b>25.62</b>	<b>45.37</b>	<b>76.45</b>
		unseen→seen	<b>8.61</b>	<b>20.24</b>	<b>63.81</b>	en→seen	15.80	28.47	72.35
		unseen→unseen	<b>7.15</b>	<b>23.70</b>	<b>54.51</b>	en→unseen	<b>17.47</b>	<b>35.00</b>	<b>65.24</b>
ar-de-fr-nl-ru-zh	MTInstruct	overall	11.79	26.36	63.22	seen→en	31.97	52.93	79.72
		seen→seen	22.68	35.32	76.39	unseen→en	26.20	37.77	78.22
		seen→unseen	7.10	24.50	55.18	en→unseen	16.17	33.27	63.50
		unseen→seen	12.56	24.74	68.83	seen→en	24.73	43.07	74.88
		unseen→unseen	6.78	22.62	53.69	unseen→en	21.18	35.52	70.86
	MT+Align	overall	<b>12.13</b>	<b>26.65</b>	<b>63.23</b>	en→xx	<b>21.33</b>	<b>35.65</b>	<b>70.99</b>
		seen→seen	<b>23.67</b>	<b>36.53</b>	<b>76.89</b>	xx→en	<b>28.60</b>	<b>48.27</b>	<b>77.49</b>
		seen→unseen	<b>7.27</b>	24.32	54.96	en→seen	<b>26.30</b>	37.63	<b>78.25</b>
		unseen→seen	<b>12.92</b>	<b>25.29</b>	<b>69.10</b>	en→unseen	<b>16.37</b>	<b>33.67</b>	<b>63.73</b>
		unseen→unseen	6.68	22.30	53.19	seen→en	<b>32.03</b>	<b>53.07</b>	<b>79.93</b>

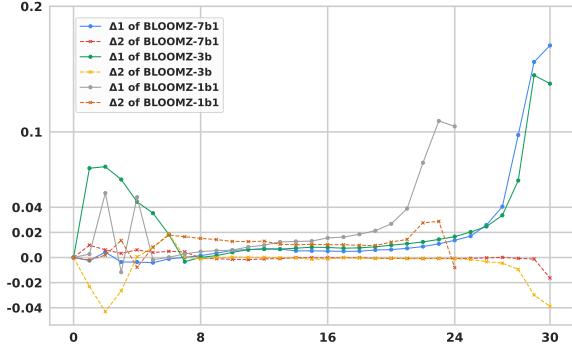
Table 5: **Results of BLOOMZ+3 without fine-tuning or fine-tuned with MTInstruct, or MT+Align.** Scores that surpass the MTInstruct baseline are marked in **bold**. xx includes seen and unseen languages.

### 4.3 BLOOMZ+3 Zero-shot Evaluation

Tab. 5 reports the results of the two settings, de-nl-ru and ar-de-fr-nl-ru-zh. Results of MT+Align+Hint+Revise and pivot-based translation are reported in App. C and H. In the de-nl-ru setting, where BLOOMZ was fine-tuned with the three unseen languages, we noticed MT+Align consistently outperformed the MTInstruct baseline across all evaluated zero-shot directions. Notably, MT+Align enhanced the translation quality for unseen→seen and seen→unseen directions compared to w/o fine-tuning and MTInstruct, given that the model was solely fine-tuned on de, nl, and ru data. This suggested AlignInstruct not only benefits the languages supplied in the data but also has a positive impact on other languages through

cross-lingual alignment supervision. In terms of supervised directions involving English, we noticed performance improvements associated with unseen languages, and regression in seen ones. The regression may be attributed to forgetting for the absence of seen languages in fine-tuning data. Indeed, continuous exposure to English maintained the translation quality for seen→en. As LoRA is modular, the regression can be mitigated by detaching the LoRA parameters for seen languages.

The ar-de-fr-nl-ru-zh setting yielded a consistently higher translation quality across all directions when compared with the de-nl-ru setting. This improvement was expected, as all the six languages were included. Translation quality improved for when generating seen languages under



**Figure 3: Differences in cosine similarity of layer-wise embeddings for BLOOMZ+24.**  $\Delta 1$  represents the changes from the unmodified BLOOMZ to the one on MTInstruct, and  $\Delta 2$  from MTInstruct to MT+Align.

the zero-shot scenario. However, the same observation cannot be made for unseen languages. This phenomenon underscored the effectiveness of AlignInstruct in enhancing translation quality for BLOOMZ’s supported languages, but suggested limitations for unseen languages when mixed with supported languages in zero-shot scenarios. In the supervised directions, we found all translation directions surpassed the performance of the MTInstruct baseline. This highlighted the overall effectiveness of AlignInstruct in enhancing translation quality across a range of supervised directions.

#### 4.4 How did MTInstruct and AlignInstruct Impact BLOOMZ’s Representations?

This section analyzed the layer-wise cosine similarities between the embeddings of parallel sentences to understand the changes in internal representations after fine-tuning. The parallel sentences were prepared from the English-centric validation datasets. We then mean-pool the outputs at each layer as sentence embeddings and compute the cosine similarities, as illustrated in Fig. 3. Results for BLOOMZ+3 are discussed in App. D.

We observed that, after MTInstruct fine-tuning, the cosine similarities rose in nearly all layers ( $\Delta 1$ , Fig. 3). This may be interpreted as enhanced cross-lingual alignment, and as indicating the acquisition of translation capabilities. Upon further combination with AlignInstruct ( $\Delta 2$ , Fig. 3), the degree of cross-lingual alignment rose in the early layers (layers 4 - 7) then diminished in the final layers (layers 29 & 30). This pattern aligned with the characteristics of encoder-decoder multilingual NMT models, where language-agnostic encoder representations with language-specific decoder representations im-

prove multilingual NMT performance (Liu et al., 2021; Wu et al., 2021; Mao et al., 2023). This highlights the beneficial impact of AlignInstruct.

## 5 Related Work

**Prompting LLMs for MT** LLMs have shown good performance for multilingual MT through few-shot in-context learning (ICL) (Jiao et al., 2023). Agrawal et al. (2023) and Zhang et al. (2023a) explored strategies to compose better examples for ICL for XGLM-7.5B (Lin et al., 2022) and GLM-130B (Zeng et al., 2023). Ghazvininejad et al. (2023), Peng et al. (2023), and Moslem et al. (2023) claimed that dictionary-based hints and domain-specific style information can improve prompting OPT (Zhang et al., 2022), GPT-3.5 (Brown et al., 2020), and BLOOM (Scao et al., 2022) for MT. He et al. (2023) used LLMs to mine useful knowledge for prompting GPT-3.5 for MT.

**Fine-tuning LLMs for MT** ICL-based methods do not support languages unseen during pre-training. Current approaches address this issue via fine-tuning. Zhang et al. (2023b) explored adding new languages to LLaMA (Touvron et al., 2023a) with interactive translation task for unseen high-resource languages. However, similar task datasets are usually not available for most unseen, low-resource languages. Li et al. (2023) and Xu et al. (2023a) showed multilingual fine-tuning with translation instructions can improve the translation ability in supported languages. Our study extended their finding to apply in the context of unseen, low-resource languages. In parallel research, Yang et al. (2023) undertook MT instruction fine-tuning in a massively multilingual context for unseen languages. However, their emphasis was on fine-tuning curriculum based on resource availability of languages, whereas we exclusively centered on low-resource languages and instruction tuning tasks.

## 6 Conclusion

In this study, we introduced AlignInstruct for enhancing the fine-tuning of LLMs for MT in unseen, low-resource languages while limiting the use of additional training corpora. Our multilingual and zero-shot findings demonstrated the strength of AlignInstruct over the MTInstruct baseline and other instruction variants. Our future work pertains to exploring using large monolingual corpora of unseen languages for MT and refining the model capability to generalize across diverse MT prompts.

## 528 Limitations

529 **Multilingual LLMs** In this study, our investigations  
530 were confined to the fine-tuning of BLOOMZ  
531 models with sizes of 1.1B, 3B, and 7.1B. We did  
532 not experiment with the 175B BLOOMZ model  
533 due to computational resource constraints. How-  
534 ever, examining this model could provide valuable  
535 insights into the efficacy of our proposed tech-  
536 niques. Additionally, it would be instructive to  
537 experiment with other recent open-source multilin-  
538 gual LLMs, such as mGPT (Shliazhko et al., 2022)  
539 and LLaMa2 (Touvron et al., 2023b).

540 **PEFT Methods and Adapters** As discussed in the  
541 BLOOM+1 paper (Yong et al., 2023), alternative  
542 PEFT techniques, such as (IA)<sup>3</sup> (Liu et al., 2022),  
543 have the potential to enhance the adaptation perfor-  
544 mance of LLM pre-training for previously unseen  
545 languages. These approaches are worth exploring  
546 for MT fine-tuning in such languages, in addition to  
547 the LoRA methods employed in this study. Further-  
548 more, our exploration was limited to fine-tuning  
549 multiple languages using shared additional parame-  
550 ters. Investigating efficient adaptation through the  
551 use of the mixture of experts (MoE) approach for  
552 MT tasks (Fan et al., 2021; Costa-jussà et al., 2022;  
553 Mohammadshahi et al., 2022; Koishkenov et al.,  
554 2023; Xu et al., 2023b) presents another intriguing  
555 avenue for LLM fine-tuning.

556 **Instruction Fine-tuning Data** Another limitation  
557 of our study is that we exclusively explored MT  
558 instruction fine-tuning using fixed templates to cre-  
559 ate MT and alignment instructions. Investigat-  
560 ing varied templates (either manually (Yang et al.,  
561 2023) or automatically constructed (Zhou et al.,  
562 2023)) might enhance the fine-tuned MT model’s  
563 ability to generalize across different MT task de-  
564 scriptions. Additionally, leveraging large monolin-  
565 gual corpora in unseen languages could potentially  
566 enhance the effectiveness of monolingual instruc-  
567 tions for MT downstream tasks, offering further  
568 insights beyond the resource-constrained scenar-  
569 os examined in this work. Furthermore, the cre-  
570 ation and utilization of instruction tuning datasets,  
571 akin to xP3 (Muennighoff et al., 2023), for unseen,  
572 low-resource languages could potentially amplify  
573 LLMs’ proficiency in following instructions in such  
574 languages. Zhu et al. (2023b) has investigated mul-  
575 tilingual instruction tuning datasets. However, the  
576 scalability of such high-quality datasets to thou-  
577 sandes of low-resource languages still remains to be  
578 addressed.

## 579 Comparison with the State-of-the-art Multilin- 580 gual NMT Models

581 In this study, we refrained  
582 from contrasting translations in low-resource lan-  
583 guages with best-performing multilingual NMT  
584 models like NLLB-200 (Costa-jussà et al., 2022),  
585 as our primary objective centered on enhancing  
586 the MTInstruct baseline through improved cross-  
587 lingual alignment within LLMs, rather than delv-  
588 ing into the best combination of techniques for MT  
589 fine-tuning in LLMs. In future exploration, our  
590 methods can potentially be integrated with the MT  
591 fine-tuning paradigm proposed by the concurrent  
592 work of Xu et al. (2023a), paving the way for ele-  
593 vating the state-of-the-art translation quality using  
594 LLMs.

## 594 References

595 Sweta Agrawal, Chunting Zhou, Mike Lewis, Luke  
596 Zettlemoyer, and Marjan Ghazvininejad. 2023. In  
597 context examples selection for machine translation.  
598 In *Findings of the Association for Computational  
599 Linguistics: ACL 2023*, pages 8857–8873, Toronto,  
600 Canada. Association for Computational Linguistics.

601 Rohan Anil, Andrew M. Dai, Orhan Firat, Melvin John-  
602 son, Dmitry Lepikhin, Alexandre Passos, Siamak  
603 Shakeri, Emanuel Taropa, Paige Bailey, Zhifeng  
604 Chen, Eric Chu, Jonathan H. Clark, Laurent El  
605 Shafey, Yanping Huang, Kathy Meier-Hellstern, Gau-  
606 rav Mishra, Erica Moreira, Mark Omernick, Kevin  
607 Robinson, Sebastian Ruder, Yi Tay, Kefan Xiao,  
608 Yuanzhong Xu, Yujing Zhang, Gustavo Hernández  
609 Ábreo, Junwhan Ahn, Jacob Austin, Paul Barham,  
610 Jan A. Botha, James Bradbury, Siddhartha Brahma,  
611 Kevin Brooks, Michele Catasta, Yong Cheng, Colin  
612 Cherry, Christopher A. Choquette-Choo, Aakanksha  
613 Chowdhery, Clément Crepy, Shachi Dave, Mostafa  
614 Dehghani, Sunipa Dev, Jacob Devlin, Mark Díaz,  
615 Nan Du, Ethan Dyer, Vladimir Feinberg, Fangxi-  
616 aoyu Feng, Vlad Fienber, Markus Freitag, Xavier  
617 Garcia, Sebastian Gehrmann, Lucas Gonzalez, and  
618 et al. 2023. *Palm 2 technical report*. *CoRR*,  
619 abs/2305.10403.

620 Mikel Artetxe, Sebastian Ruder, and Dani Yogatama.  
621 2020. On the cross-lingual transferability of mono-  
622 lingual representations. In *Proceedings of the 58th  
623 Annual Meeting of the Association for Computational  
624 Linguistics*, pages 4623–4637, Online. Association  
625 for Computational Linguistics.

626 Dzmitry Bahdanau, Kyunghyun Cho, and Yoshua Ben-  
627 gio. 2015. Neural machine translation by jointly  
628 learning to align and translate. In *3rd International  
629 Conference on Learning Representations, ICLR 2015,  
630 San Diego, CA, USA, May 7-9, 2015, Conference  
631 Track Proceedings*.

632 Yoshua Bengio, Jérôme Louradour, Ronan Collobert,  
633 and Jason Weston. 2009. Curriculum learning. In

634	<i>Proceedings of the 26th Annual International Conference on Machine Learning, ICML 2009, Montreal, Quebec, Canada, June 14–18, 2009, volume 382 of ACM International Conference Proceeding Series, pages 41–48.</i> ACM.	693
635		694
636		695
637		696
638		697
639	Peter F. Brown, Stephen A. Della Pietra, Vincent J. Della Pietra, and Robert L. Mercer. 1993. <b>The mathematics of statistical machine translation: Parameter estimation.</b> <i>Computational Linguistics</i> , 19(2):263–311.	698
640		699
641		700
642		701
643		
644	Tom B. Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared Kaplan, Prafulla Dhariwal, Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, Sandhini Agarwal, Ariel Herbert-Voss, Gretchen Krueger, Tom Henighan, Rewon Child, Aditya Ramesh, Daniel M. Ziegler, Jeffrey Wu, Clemens Winter, Christopher Hesse, Mark Chen, Eric Sigler, Mateusz Litwin, Scott Gray, Benjamin Chess, Jack Clark, Christopher Berner, Sam McCandlish, Alec Radford, Ilya Sutskever, and Dario Amodei. 2020. <b>Language models are few-shot learners.</b> In <i>Advances in Neural Information Processing Systems 33: Annual Conference on Neural Information Processing Systems 2020, NeurIPS 2020, December 6–12, 2020, virtual.</i>	702
645		703
646		704
647		705
648		706
649		707
650		708
651		709
652		710
653		711
654		712
655		713
656		714
657		715
658		716
659	Rich Caruana. 1997. <b>Multitask learning.</b> <i>Machine Learning</i> , 28(1):41–75.	717
660		
661	Aakanksha Chowdhery, Sharan Narang, Jacob Devlin, Maarten Bosma, Gaurav Mishra, Adam Roberts, Paul Barham, Hyung Won Chung, Charles Sutton, Sebastian Gehrmann, Parker Schuh, Kensen Shi, Sasha Tsvyashchenko, Joshua Maynez, Abhishek Rao, Parker Barnes, Yi Tay, Noam Shazeer, Vinodkumar Prabhakaran, Emily Reif, Nan Du, Ben Hutchinson, Reiner Pope, James Bradbury, Jacob Austin, Michael Isard, Guy Gur-Ari, Pengcheng Yin, Toju Duke, Anselm Levskaya, Sanjay Ghemawat, Sunipa Dev, Henryk Michalewski, Xavier Garcia, Vedant Misra, Kevin Robinson, Liam Fedus, Denny Zhou, Daphne Ippolito, David Luan, Hyeontaek Lim, Barret Zoph, Alexander Spiridonov, Ryan Sepassi, David Dohan, Shivani Agrawal, Mark Omernick, Andrew M. Dai, Thanumalayan Sankaranarayana Pillai, Marie Pellat, Aitor Lewkowycz, Erica Moreira, Rewon Child, Oleksandr Polozov, Katherine Lee, Zongwei Zhou, Xuezhi Wang, Brennan Saeta, Mark Diaz, Orhan Firat, Michele Catasta, Jason Wei, Kathy Meier-Hellstern, Douglas Eck, Jeff Dean, Slav Petrov, and Noah Fiedel. 2022. <b>Palm: Scaling language modeling with pathways.</b> <i>CoRR</i> , abs/2204.02311.	718
662		719
663		720
664		721
665		722
666		723
667		724
668		725
669		726
670		
671	Javier de la Rosa and Andrés Fernández. 2022. <b>Zero-shot reading comprehension and reasoning for spanish with BERTIN GPT-J-6B.</b> In <i>Proceedings of the Iberian Languages Evaluation Forum (IberLEF 2022) co-located with the Conference of the Spanish Society for Natural Language Processing (SEPLN 2022), A Coruña, Spain, September 20, 2022, volume 3202 of CEUR Workshop Proceedings.</i> CEUR-WS.org.	727
672		728
673		729
674		730
675		731
676		732
677		733
678		
679	Chris Dyer, Victor Chahuneau, and Noah A. Smith. 2013. <b>A simple, fast, and effective reparameterization of IBM model 2.</b> In <i>Proceedings of the 2013 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies</i> , pages 644–648, Atlanta, Georgia. Association for Computational Linguistics.	727
680		728
681		729
682		730
683		731
684	Abteen Ebrahimi and Katharina Kann. 2021. <b>How to adapt your pretrained multilingual model to 1600 languages.</b> In <i>Proceedings of the 59th Annual Meeting of the Association for Computational Linguistics and the 11th International Joint Conference on Natural Language Processing (Volume 1: Long Papers)</i> , pages 4555–4567, Online. Association for Computational Linguistics.	732
685		733
686		734
687		735
688		736
689		737
690		738
691		739
692		740
693		741
694	Angela Fan, Shruti Bhosale, Holger Schwenk, Zhiyi Ma, Ahmed El-Kishky, Siddharth Goyal, Mandeep Baines, Onur Celebi, Guillaume Wenzek, Vishrav Chaudhary, Naman Goyal, Tom Birch, Vitaliy Liptchinsky, Sergey Edunov, Michael Auli, and Armand Joulin. 2021. <b>Beyond english-centric multilingual machine translation.</b> <i>J. Mach. Learn. Res.</i> , 22:107:1–107:48.	742
695		743
696		744
697		745
698		746
699		747
700		748
701		749
702	Marjan Ghazvininejad, Hila Gonen, and Luke Zettlemoyer. 2023. <b>Dictionary-based phrase-level prompt-</b>	750
703		751

752	ing of large language models for machine translation.	810
753	<i>CoRR</i> , abs/2302.07856.	811
754	Zhiwei He, Tian Liang, Wenxiang Jiao, Zhuosheng	812
755	Zhang, Yujiu Yang, Rui Wang, Zhaopeng Tu, Shum-	813
756	ing Shi, and Xing Wang. 2023. Exploring human-	814
757	like translation strategy with large language models.	815
758	<i>CoRR</i> , abs/2305.04118.	816
759	Edward J. Hu, Yelong Shen, Phillip Wallis, Zeyuan	817
760	Allen-Zhu, Yuanzhi Li, Shean Wang, Lu Wang, and	
761	Weizhu Chen. 2022. <b>Lora: Low-rank adaptation of</b>	
762	<b>large language models</b> . In <i>The Tenth International</i>	
763	<i>Conference on Learning Representations, ICLR 2022,</i>	
764	<i>Virtual Event, April 25-29, 2022</i> . OpenReview.net.	
765	Wenxiang Jiao, Wenxuan Wang, Jen-tse Huang, Xing	818
766	Wang, and Zhaopeng Tu. 2023. <b>Is chatgpt A</b>	819
767	<b>good translator? A preliminary study.</b> <i>CoRR</i> ,	820
768	abs/2301.08745.	821
769	Philipp Koehn. 2004. <b>Statistical significance tests for</b>	822
770	<b>machine translation evaluation.</b> In <i>Proceedings of the</i>	823
771	<i>2004 Conference on Empirical Methods in Natural</i>	824
772	<i>Language Processing</i> , pages 388–395, Barcelona,	825
773	Spain. Association for Computational Linguistics.	826
774	Yeskendir Koishkenov, Alexandre Berard, and Vas-	827
775	silina Nikoulina. 2023. <b>Memory-efficient NLLB-200:</b>	828
776	<b>Language-specific expert pruning of a massively mul-</b>	829
777	<b>tilingual machine translation model.</b> In <i>Proceedings</i>	830
778	<i>of the 61st Annual Meeting of the Association for</i>	
779	<i>Computational Linguistics (Volume 1: Long Papers)</i> ,	
780	pages 3567–3585, Toronto, Canada. Association for	
781	Computational Linguistics.	
782	Patrik Lambert, Simon Petitrenaud, Yanjun Ma, and	831
783	Andy Way. 2012. <b>What types of word alignment im-</b>	832
784	<b>prove statistical machine translation?</b> <i>Mach. Transl.</i>	833
785	26(4):289–323.	834
786	Jiahuan Li, Hao Zhou, Shujian Huang, Shanbo Chen,	835
787	and Jiajun Chen. 2023. <b>Eliciting the translation</b>	836
788	<b>ability of large language models via multilingual</b>	
789	<b>finetuning with translation instructions.</b> <i>CoRR</i> ,	
790	abs/2305.15083.	
791	Xiaozhuan Liang, Ningyu Zhang, Siyuan Cheng,	837
792	Zhenru Zhang, Chuanqi Tan, and Huajun Chen. 2022.	838
793	<b>Contrastive demonstration tuning for pre-trained lan-</b>	839
794	<b>guage models.</b> In <i>Findings of the Association for</i>	840
795	<i>Computational Linguistics: EMNLP 2022</i> , pages	841
796	799–811, Abu Dhabi, United Arab Emirates. Associa-	842
797	tion for Computational Linguistics.	843
798	Xi Victoria Lin, Todor Mihaylov, Mikel Artetxe, Tianlu	844
799	Wang, Shuhui Chen, Daniel Simig, Myle Ott, Na-	845
800	man Goyal, Shruti Bhosale, Jingfei Du, Ramakanth	846
801	Pasunuru, Sam Shleifer, Punit Singh Koura, Vishrav	847
802	Chaudhary, Brian O’Horo, Jeff Wang, Luke Zettle-	848
803	moyer, Zornitsa Kozareva, Mona Diab, Veselin Stoy-	849
804	anov, and Xian Li. 2022. <b>Few-shot learning with</b>	850
805	<b>multilingual generative language models.</b> In <i>Proce-</i>	851
806	<i>edings of the 2022 Conference on Empirical Methods</i>	
807	<i>in Natural Language Processing</i> , pages 9019–9052,	
808	Abu Dhabi, United Arab Emirates. Association for	
809	Computational Linguistics.	
810	Zehui Lin, Xiao Pan, Mingxuan Wang, Xipeng Qiu,	852
811	Jiangtao Feng, Hao Zhou, and Lei Li. 2020. <b>Pre-</b>	853
812	<b>training multilingual neural machine translation by</b>	854
813	<b>leveraging alignment information.</b> In <i>Proceedings</i>	855
814	<i>of the 2020 Conference on Empirical Methods in</i>	856
815	<i>Natural Language Processing (EMNLP)</i> , pages 2649–	857
816	2663, Online. Association for Computational Lin-	858
817	guistics.	859
818	Danni Liu, Jan Niehues, James Cross, Francisco	860
819	Guzmán, and Xian Li. 2021. <b>Improving zero-shot</b>	861
820	<b>translation by disentangling positional information.</b>	862
821	In <i>Proceedings of the 59th Annual Meeting of the</i>	863
822	<i>Association for Computational Linguistics and the</i>	864
823	<i>11th International Joint Conference on Natural Lan-</i>	865
824	<i>guage Processing (Volume 1: Long Papers)</i> , pages	866
825	1259–1273, Online. Association for Computational	
826	Linguistics.	
827	Haokun Liu, Derek Tam, Mohammed Muqeeth, Jay Mo-	867
828	hta, Tenghao Huang, Mohit Bansal, and Colin Raffel.	868
829	2022. <b>Few-shot parameter-efficient fine-tuning is bet-</b>	869
830	<b>ter and cheaper than in-context learning.</b> In <i>NeurIPS</i> .	870
831	Yinhan Liu, Jiatao Gu, Naman Goyal, Xian Li, Sergey	871
832	Edunov, Marjan Ghazvininejad, Mike Lewis, and Luke	872
833	Zettlemoyer. 2020. <b>Multilingual denoising pre-</b>	873
834	<b>training for neural machine translation.</b> <i>Transac-</i>	874
835	<i>tions of the Association for Computational Linguis-</i>	875
836	<i>tics</i> , 8:726–742.	876
837	Zhuoyuan Mao, Chenhui Chu, Raj Dabre, Haiyue Song,	877
838	Zhen Wan, and Sadao Kurohashi. 2022. <b>When do</b>	878
839	<b>contrastive word alignments improve many-to-many</b>	879
840	<b>neural machine translation?</b> In <i>Findings of the Asso-</i>	880
841	<i>ciation for Computational Linguistics: NAACL 2022</i> ,	881
842	pages 1766–1775, Seattle, United States. Associa-	882
843	tion for Computational Linguistics.	883
844	Zhuoyuan Mao, Raj Dabre, Qianying Liu, Haiyue Song,	884
845	Chenhui Chu, and Sadao Kurohashi. 2023. <b>Explor-</b>	885
846	<b>ing the impact of layer normalization for zero-shot</b>	886
847	<b>neural machine translation.</b> In <i>Proceedings of the</i>	887
848	<i>61st Annual Meeting of the Association for Compu-</i>	888
849	<i>tational Linguistics (Volume 2: Short Papers)</i> , pages	889
850	1300–1316, Toronto, Canada. Association for Com-	890
851	putational Linguistics.	891
852	Paulius Micikevicius, Sharan Narang, Jonah Alben,	892
853	Gregory F. Diamos, Erich Elsen, David García,	893
854	Boris Ginsburg, Michael Houston, Oleksii Kuchaiev,	894
855	Ganesh Venkatesh, and Hao Wu. 2018. <b>Mixed pre-</b>	895
856	<b>cision training.</b> In <i>6th International Conference on</i>	896
857	<i>Learning Representations, ICLR 2018, Vancouver,</i>	897
858	<i>BC, Canada, April 30 - May 3, 2018, Conference</i>	898
859	<i>Track Proceedings.</i> OpenReview.net.	899
860	Swaroop Mishra, Daniel Khashabi, Chitta Baral, and	900
861	Hannaneh Hajishirzi. 2022. <b>Cross-task generaliza-</b>	901
862	<b>tion via natural language crowdsourcing instructions.</b>	902
863	In <i>Proceedings of the 60th Annual Meeting of the</i>	903
864	<i>Association for Computational Linguistics (Volume</i>	904
865	<i>1: Long Papers)</i> , pages 3470–3487, Dublin, Ireland.	905
866	Association for Computational Linguistics.	906

867	Alireza Mohammadshahi, Vassilina Nikoulina, Alexandre Berard, Caroline Brun, James Henderson, and Laurent Besacier. 2022. <b>SMaLL-100: Introducing shallow multilingual machine translation model for low-resource languages.</b> In <i>Proceedings of the 2022 Conference on Empirical Methods in Natural Language Processing</i> , pages 8348–8359, Abu Dhabi, United Arab Emirates. Association for Computational Linguistics.	923
868		924
869		925
870		926
871		927
872		928
873		929
874		
875		
876	Yasmin Moslem, Rejwanul Haque, John D. Kelleher, and Andy Way. 2023. <b>Adaptive machine translation with large language models.</b> In <i>Proceedings of the 24th Annual Conference of the European Association for Machine Translation</i> , pages 227–237, Tampere, Finland. European Association for Machine Translation.	930
877		931
878		932
879		933
880		934
881		
882		
883	Niklas Muennighoff, Thomas Wang, Lintang Sutawika, Adam Roberts, Stella Biderman, Teven Le Scao, M Saiful Bari, Sheng Shen, Zheng Xin Yong, Hayley Schoelkopf, Xiangru Tang, Dragomir Radev, Alham Fikri Aji, Khalid Almubarak, Samuel Albanie, Zaid Alyafeai, Albert Webson, Edward Raff, and Colin Raffel. 2023. <b>Crosslingual generalization through multitask finetuning.</b> In <i>Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)</i> , pages 15991–16111, Toronto, Canada. Association for Computational Linguistics.	935
884		936
885		937
886		938
887		939
888		940
889		941
890		
891		
892		
893		
894		
895	Benjamin Muller, Antonios Anastasopoulos, Benoît Sagot, and Djamé Seddah. 2021. <b>When being unseen from mBERT is just the beginning: Handling new languages with multilingual language models.</b> In <i>Proceedings of the 2021 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies</i> , pages 448–462, Online. Association for Computational Linguistics.	942
896		943
897		944
898		945
899		946
900		
901		
902		
903		
904	Martin Müller and Florian Laurent. 2022. <b>Cedille: A large autoregressive french language model.</b> <i>CoRR</i> , abs/2202.03371.	947
905		948
906		949
907	Graham Neubig and Junjie Hu. 2018. <b>Rapid adaptation of neural machine translation to new languages.</b> In <i>Proceedings of the 2018 Conference on Empirical Methods in Natural Language Processing</i> , pages 875–880, Brussels, Belgium. Association for Computational Linguistics.	950
908		951
909		
910		
911		
912		
913	OpenAI. 2023. <b>GPT-4 technical report.</b> <i>CoRR</i> , abs/2303.08774.	952
914		953
915	Long Ouyang, Jeffrey Wu, Xu Jiang, Diogo Almeida, Carroll L. Wainwright, Pamela Mishkin, Chong Zhang, Sandhini Agarwal, Katarina Slama, Alex Ray, John Schulman, Jacob Hilton, Fraser Kelton, Luke Miller, Maddie Simens, Amanda Askell, Peter Welinder, Paul F. Christiano, Jan Leike, and Ryan Lowe. 2022. <b>Training language models to follow instructions with human feedback.</b> In <i>NeurIPS</i> .	954
916		
917		
918		
919		
920		
921		
922		
914	Kishore Papineni, Salim Roukos, Todd Ward, and Wei Jing Zhu. 2002. <b>Bleu: a method for automatic evaluation of machine translation.</b> In <i>Proceedings of the 40th Annual Meeting of the Association for Computational Linguistics</i> , pages 311–318, Philadelphia, Pennsylvania, USA. Association for Computational Linguistics.	962
915		963
916		964
917		965
918		966
919		967
920		
921		
922		
923	Keqin Peng, Liang Ding, Qihuang Zhong, Li Shen, Xuebo Liu, Min Zhang, Yuanxin Ouyang, and Dacheng Tao. 2023. <b>Towards making the most of chatgpt for machine translation.</b> <i>CoRR</i> , abs/2303.13780.	968
924		969
925		970
926		971
927		972
928		973
929		974
930		975
931		976
932		977
933		978
934		979

980	M Saiful Bari, Canwen Xu, Urmish Thakker, Shanya Sharma Sharma, Eliza Szczechla, Taewoon Kim, Gunjan Chhablani, Nihal V. Nayak, Debajyoti Datta, Jonathan Chang, Mike Tian-Jian Jiang, Han Wang, Matteo Manica, Sheng Shen, Zheng Xin Yong, Harshit Pandey, Rachel Bawden, Thomas Wang, Trishala Neeraj, Jos Rozen, Abheesht Sharma, Andrea Santilli, Thibault Févry, Jason Alan Fries, Ryan Teehan, Teven Le Scao, Stella Biderman, Leo Gao, Thomas Wolf, and Alexander M. Rush. 2022. <b>Multitask prompted training enables zero-shot task generalization.</b> In <i>The Tenth International Conference on Learning Representations, ICLR 2022, Virtual Event, April 25-29, 2022</i> . OpenReview.net.	1039
981	Bashlykov, Soumya Batra, Prajwal Bhargava, Shruti Bhosale, Dan Bikel, Lukas Blecher, Cristian Canton-Ferrer, Moya Chen, Guillem Cucurull, David Esiobu, Jude Fernandes, Jeremy Fu, Wenyin Fu, Brian Fuller, Cynthia Gao, Vedanuj Goswami, Naman Goyal, Anthony Hartshorn, Saghar Hosseini, Rui Hou, Hakan Inan, Marcin Kardas, Viktor Kerkez, Madian Khabsa, Isabel Kloumann, Artem Korenev, Punit Singh Koura, Marie-Anne Lachaux, Thibaut Lavril, Jenya Lee, Diana Liskovich, Yinghai Lu, Yuning Mao, Xavier Martinet, Todor Miaylov, Pushkar Mishra, Igor Molybog, Yixin Nie, Andrew Poulton, Jeremy Reizenstein, Rashi Rungta, Kalyan Saladi, Alan Schelten, Ruan Silva, Eric Michael Smith, Ranjan Subramanian, Xiaoqing Ellen Tan, Binh Tang, Ross Taylor, Adina Williams, Jian Xiang Kuan, Puxin Xu, Zheng Yan, Iliyan Zarov, Yuchen Zhang, Angela Fan, Melanie Kambadur, Sharan Narang, Aurélien Rodriguez, Robert Stojnic, Sergey Edunov, and Thomas Scialom. 2023b. <b>Llama 2: Open foundation and fine-tuned chat models.</b> <i>CoRR</i> , abs/2307.09288.	1040
982		1041
983		1042
984		1043
985		1044
986		1045
987		1046
988		1047
989		1048
990		1049
991		1050
992		1051
993		1052
994	Teven Le Scao, Angela Fan, Christopher Akiki, Elie Pavlick, Suzana Ilic, Daniel Hesslow, Roman Castagné, Alexandra Sasha Luccioni, François Yvon, Matthias Gallé, Jonathan Tow, Alexander M. Rush, Stella Biderman, Albert Webson, Pawan Sasanka Ammanamanchi, Thomas Wang, Benoît Sagot, Niklas Muennighoff, Albert Villanova del Moral, Olatunji Ruwase, Rachel Bawden, Stas Bekman, Angelina McMillan-Major, Iz Beltagy, Huu Nguyen, Lucile Saulnier, Samson Tan, Pedro Ortiz Suarez, Victor Sanh, Hugo Laurençon, Yacine Jernite, Julien Launay, Margaret Mitchell, Colin Raffel, Aaron Gokaslan, Adi Simhi, Aitor Soroa, Alham Fikri Aji, Amit Alfassy, Anna Rogers, Ariel Kreisberg Nitzav, Canwen Xu, Chenghao Mou, Chris Emezue, Christopher Klamm, Colin Leong, Daniel van Strien, David Ifeoluwa Adelani, and et al. 2022. <b>BLOOM: A 176b-parameter open-access multilingual language model.</b> <i>CoRR</i> , abs/2211.05100.	1053
995		1054
996		1055
997		1056
998		1057
999		1058
1000		1059
1001		
1002		
1003		
1004		
1005		
1006		
1007		
1008		
1009		
1010		
1011		
1012		
1013	Oleh Shliazhko, Alena Fenogenova, Maria Tikhonova, Vladislav Mikhailov, Anastasia Kozlova, and Tatiana Shavrina. 2022. <b>mgpt: Few-shot learners go multilingual.</b> <i>CoRR</i> , abs/2204.07580.	1060
1014		1061
1015		1062
1016		1063
1017		1064
1018		1065
1019		
1020		
1021		
1022		
1023		
1024	Kaitao Song, Xu Tan, Tao Qin, Jianfeng Lu, and Tie-Yan Liu. 2019. <b>MASS: masked sequence to sequence pre-training for language generation.</b> In <i>Proceedings of the 36th International Conference on Machine Learning, ICML 2019, 9-15 June 2019, Long Beach, California, USA</i> , volume 97 of <i>Proceedings of Machine Learning Research</i> , pages 5926–5936. PMLR.	1066
1025		1067
1026		1068
1027		1069
1028		1070
1029		1071
1030		1072
1031		1073
1032		1074
1033		1075
1034		1076
1035		1077
1036		1078
1037	Ilya Sutskever, Oriol Vinyals, and Quoc V. Le. 2014. <b>Sequence to sequence learning with neural networks.</b> In <i>Advances in Neural Information Processing Systems 27: Annual Conference on Neural Information Processing Systems 2014, December 8-13 2014, Montreal, Quebec, Canada</i> , pages 3104–3112.	1079
1038		1080
1039		1081
1040		1082
1041		1083
1042		
1043		
1044		
1045		
1046		
1047		
1048		
1049		
1050		
1051		
1052		
1053		
1054		
1055		
1056		
1057		
1058		
1059		
1060		
1061		
1062		
1063		
1064		
1065		
1066		
1067		
1068		
1069		
1070		
1071		
1072		
1073		
1074		
1075		
1076		
1077		
1078		
1079		
1080		
1081		
1082		
1083		
1084		
1085		
1086		
1087		
1088		
1089		
1090		
1091		
1092		
1093		
1094		
1095		
1096		
1097		
1098		

1099	translation: Boosting translation performance of large language models. <i>CoRR</i> , abs/2309.11674.	1156
1100		1157
1101		1158
1102		1159
1103		1160
1104		1161
1105		
1106		
1107		
1108		
1109		
1110	Haoran Xu, Weiting Tan, Shuyue Stella Li, Yunmo Chen, Benjamin Van Durme, Philipp Koehn, and Kenton Murray. 2023b. Condensing multilingual knowledge with lightweight language-specific modules. <i>CoRR</i> , abs/2305.13993.	1162
1111		1163
1112		1164
1113		1165
1114		1166
1115		1167
1116		
1117		
1118		
1119		
1120		
1121		
1122		
1123	Zheng Xin Yong, Hailey Schoelkopf, Niklas Muenighoff, Alham Fikri Aji, David Ifeoluwa Adelani, Khalid Almubarak, M Saiful Bari, Lintang Sutawika, Jungo Kasai, Ahmed Baruwa, Genta Winata, Stella Biderman, Edward Raff, Dragomir Radev, and Vasilina Nikoulina. 2023. <b>BLOOM+1: Adding language support to BLOOM for zero-shot prompting</b> . In <i>Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)</i> , pages 11682–11703, Toronto, Canada. Association for Computational Linguistics.	1168
1124		1169
1125		1170
1126		1171
1127		1172
1128		
1129		
1130		
1131		
1132		
1133		
1134		
1135		
1136		
1137		
1138		
1139	Zhang Ze Yu, Lau Jia Jaw, Wong Qin Jiang, and Zhang Hui. 2023. Fine-tuning language models with generative adversarial feedback. <i>CoRR</i> , abs/2305.06176.	1173
1140		1174
1141		1175
1142		1176
1143		1177
1144		
1145		
1146		
1147		
1148		
1149		
1150		
1151		
1152		
1153		
1154		
1155		
1156		
1157		
1158		
1159		
1160		
1161		
1162		
1163		
1164		
1165		
1166		
1167		
1168		
1169		
1170		
1171		
1172		
1173		
1174		
1175		
1176		
1177		
1178		
1179		
1180		
1181		
1182		
1183		
1184		
1185		
1186		
1187		
1188		
1189		
1190		
1191		
1192		
1193		
1194		
1195		
1196		
1197		
1198		
1199		
1200		
1201		
1202		
1203		
1204		
1205		
1206		

Language	ISO 639-1	Language Family	Subgrouping	Script	Seen Script	#sent.
Afrikaans	af	Indo-European	Germanic	Latin	✓	275,512
Amharic	am	Afro-Asiatic	Semitic	Ge'ez	✗	89,027
Belarusian	be	Indo-European	Balto-Slavic	Cyrillic	✗	67,312
Welsh	cy	Indo-European	Celtic	Latin	✓	289,521
Irish	ga	Indo-European	Celtic	Latin	✓	289,524
Scottish Gaelic	gd	Indo-European	Celtic	Latin	✓	16,316
Galician	gl	Indo-European	Italic	Latin	✓	515,344
Hausa	ha	Afro-Asiatic	Chadic	Latin	✓	97,983
Georgian	ka	Kartvelian	Georgian-Zan	Georgian	✗	377,306
Kazakh	kk	Turkic	Common Turkic	Cyrillic	✗	79,927
Khmer	km	Austroasiatic	Khmeric	Khmer	✗	111,483
Kyrgyz	ky	Turkic	Common Turkic	Cyrillic	✗	27,215
Limburgish	li	Indo-European	Germanic	Latin	✓	25,535
Burmese	my	Sino-Tibetan	Burmo-Qiangic	Myanmar	✗	24,594
Norwegian Bokmål	nb	Indo-European	Germanic	Latin	✓	142,906
Norwegian Nynorsk	nn	Indo-European	Germanic	Latin	✓	486,055
Occitan	oc	Indo-European	Italic	Latin	✓	35,791
Sinhala	si	Indo-European	Indo-Aryan	Sinhala	✗	979,109
Tajik	tg	Indo-European	Iranian	Cyrillic	✗	193,882
Turkmen	tk	Turkic	Common Turkic	Latin	✓	13,110
Tatar	tt	Turkic	Common Turkic	Cyrillic	✗	100,843
Uyghur	ug	Turkic	Common Turkic	Arabic	✓	72,170
Northern Uzbek	uz	Turkic	Common Turkic	Latin	✓	173,157
Eastern Yiddish	yi	Indo-European	Germanic	Hebrew	✗	15,010
Total						4,498,632

Table 6: **Statistics of training data for BLOOMZ+24:** 24 unseen, low-resource languages for BLOOMZ. ✓ and ✗ indicate whether script is seen or unseen.

et al., 2020). We tuned the optimal learning rate for each individual experiment according to validation loss. We conducted all experiments once due to computational resource constraints and reported the average scores across all languages.

## C Results of MT+Align+Hint+Revise for BLOOMZ+3

We present the results in Tab. 7. Co-referencing the results in Tab. 5, compared with MT+Align, we observed a clear advantage for the MT+Align+Hint+Revise setting in supervised directions involving English ( $\text{en} \rightarrow \text{seen}$  and  $\text{seen} \rightarrow \text{en}$ ) in the ar-fr-de-nl-ru-zh setting. This result suggested that AlignInstruct's variants played a crucial role in preserving the BLOOMZ's capabilities for supported languages. However, in all other scenarios, AlignInstruct alone proved sufficient to enhance the performance beyond the MTInstruct baseline, but hard to achieve further improvements with additional instructions.

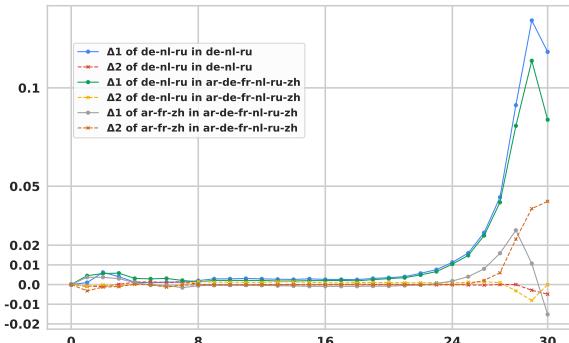


Figure 4: **Differences in cosine similarity of layer-wise embeddings for BLOOMZ+3.**  $\Delta 1$  represents the changes from the unmodified BLOOMZ to the one on MTInstruct, and  $\Delta 2$  from MTInstruct to MT+Align.

## D Representation Change of BLOOMZ+3

The representation change observed in de-nl-ru was consistent with the findings presented in Sec. 4.4, which highlighted an initial increase in cross-lingual alignment in the early layers, followed by a decrease in the final layers. When mixing fine-tuning data with supported languages, the changes

Languages	Directions	Zero-shot Directions			Supervised Directions			
		BLEU	chrF++	COMET	Directions	BLEU	chrF++	COMET
de-nl-ru	overall	<b>8.94</b>	<b>23.53</b>	<b>60.67</b>	en→xx	16.70	31.83	68.98
	seen→seen	14.00	<b>27.58</b>	<b>70.59</b>	xx→en	<b>25.18</b>	45.00	<b>76.45</b>
	seen→unseen	<b>6.49</b>	<b>23.01</b>	<b>54.92</b>	en→seen	15.97	28.53	<b>72.69</b>
	unseen→seen	<b>9.50</b>	<b>21.90</b>	<b>64.69</b>	en→unseen	<b>17.43</b>	<b>35.13</b>	<b>65.27</b>
	unseen→unseen	6.73	22.70	53.34	seen→en	25.33	46.70	77.51
ar-de-fr-nl-ru-zh	overall	<b>12.07</b>	<b>26.67</b>	63.13	unseen→en	<b>25.03</b>	<b>43.30</b>	<b>75.39</b>
	seen→seen	<b>23.52</b>	<b>36.13</b>	<b>76.62</b>	en→xx	<b>21.62</b>	<b>36.12</b>	<b>70.94</b>
	seen→unseen	<b>7.16</b>	24.48	55.02	xx→en	<b>28.92</b>	<b>48.60</b>	<b>77.50</b>
	unseen→seen	<b>12.91</b>	<b>25.23</b>	<b>68.91</b>	en→seen	<b>26.87</b>	<b>38.40</b>	<b>78.40</b>
	unseen→unseen	6.73	<b>22.65</b>	53.12	en→unseen	<b>16.37</b>	<b>33.83</b>	63.49

Table 7: **Results of BLOOMZ+3 with MT+Align+Hint+Revise.** Co-referencing Tab. 5, scores that surpass the MTInstruct baseline are marked in **bold**.

exhibited more intricate patterns. As illustrated by ar-fr-zh in ar-de-fr-nl-ru-zh in Fig. 4, sentence alignment declined after MTInstruct fine-tuning but elevated after further combining with AlignInstruct. We leave the interpretation of this nuanced behavior in future work.

## E Examples of HintInstruct and ReviseInstruct

We illustrated examples of HintInstruct and ReviseInstruct in Fig. 5.

## F Assessing Monolingual Instructions

New language capabilities may be induced through continual pre-training on monolingual next-word prediction tasks (Yong et al., 2023). The coherence of the generated sentences is crucial in MT (Wang et al., 2020; Liu et al., 2020), especially when the target languages are unseen and low-resource. We examined the significance of this approach in fostering the translation quality. We reused the same parallel corpora to avoid introducing additional monolingual datasets.

Given a monolingual sentence,  $(x_i)_1^N$ , with length  $N$  in an unseen language  $X$ . The LLM is incrementally trained on the following task:

- **Input:** “Given the context, complete the following sentence:  $x_1 x_2 \dots x_{l < N}$ ,”
- **Output:** “ $x_{l+1} x_{l+2} \dots x_N$ .”

We conducted experiments with two MonoInstruct settings: **MonoInstruct-full**, an objective to

### HintInstruct

--- Input ---  
Use the following alignment hints and translate from English to Japanese.  
Alignments between English and Japanese:  
- (future, 未来),  
- (generative, 生成),  
- (of, の),  
English: What’s the future of generative AI?  
Japanese:  
--- Output ---  
生成 AI の未来はどうなるでしょうか?

### ReviseInstruct

--- Input ---  
Given the following translation of Japanese from English, output the incorrectly translated word and correct it.  
English: What’s the future of generative AI?  
Japanese: 生成 AI の過去はどうなるでしょうか?  
--- Output ---  
The incorrectly translated word is “過去”. It should be “未来”.

Figure 5: Examples of HintInstruct and ReviseInstruct.

generate the entire sentence, and **MonoInstruct-half** for generating the latter half of the sentence given the first half, inspired by GPT (Radford et al., 2018) and MASS (Song et al., 2019), respectively. We reported the MonoInstruct results in Tab. 8. Firstly, we observed that fine-tuning MTInstruct in conjunction with either MonoInstruct-full or MonoInstruct-half harms the MT performance, which could be attributed to the inherent difficulty of monolingual instruction tasks and the limited amount of monolingual data. We found that the

Objective	OPUS en→xx			OPUS xx→en			Flores en→xx			Flores xx→en		
	BLEU	chrF++	COMET	BLEU	chrF++	COMET	BLEU	chrF++	COMET	BLEU	chrF++	COMET
MTInstruct	11.54	25.33	64.54	18.59	33.25	69.84	3.30	17.10	40.58	11.37	27.14	56.33
MT+Mono-full	9.89	22.42	62.52	15.43	29.04	66.64	3.00	16.68	40.49	10.26	25.15	54.17
MT+Mono-half	10.23	22.45	62.22	15.51	29.65	67.29	3.18	16.91	40.57	10.66	26.15	54.80
MT+Mono-full+Align	10.15	22.35	62.22	15.72	29.86	67.70	3.07	16.59	<b>40.78</b>	10.61	25.58	55.17
MT+Mono-half+Align	10.09	22.61	62.98	16.00	30.34	67.96	3.10	16.75	<b>40.70</b>	10.79	26.27	55.40
MT+Mono-full+Align+Hint+Revise	10.33	23.04	63.19	17.16	31.61	68.26	3.23	16.70	<b>40.90</b>	10.98	26.18	55.50
MT+Mono-half+Align+Hint+Revise	10.62	23.10	62.92	17.32	31.80	68.56	3.20	16.93	<b>41.00</b>	11.09	26.77	55.99

Table 8: **Results of BLOOMZ+24 fine-tuned incorporating monolingual instructions on BLOOMZ-7b1.** Scores that surpass the MTInstruct baseline are marked in **bold**.

simpler MT+Mono-half yielded better results than MT+Mono-full as richer contexts were provided. However, MonoInstruct still did not improve the MTInstruct baseline. Secondly, further combining MonoInstrcut with AlignInstruct variants yielded improvements compared with MT+Mono-full (or half), but underperformed the MTInstruct baseline. This suggested that improving MT performance with monolingual instructions is challenging without access to additional monolingual data.

## G Inference using Different MT Prompts

We investigated the performance of fine-tuned models when using various MT prompts during inference, aiming to understand models’ generalization capabilities with different test prompts. We examined five MT prompts for the fine-tuned models of BLOOMZ-7b1, following Zhang et al. (2023a), which are presented in Tab. 9. The results, showcased in Tab. 10, revealed that in comparison to the default prompt used during fine-tuning, the translation performance tended to decline when using other MT prompts. We observed that MT+Align consistently surpasses MTInstruct for xx→en translations, though the results were mixed for en→xx directions. Certain prompts, such as PROMPT-3 and PROMPT-4, exhibited a minor performance drop, while others significantly impacted translation quality. These findings underscored the need for enhancing the models’ ability to generalize across diverse MT prompts, potentially by incorporating a range of MT prompt templates during the fine-tuning process, as stated in the Limitations section.

## H Zero-shot Translation using English as Pivot

Pivot translation serves as a robust technique for zero-shot translation, especially given that we used English-centric data during fine-tuning. In Tab. 11, we present results that utilize English as an inter-

Prompt	Definition
PROMPT-default	Translate from $Y$ to $X$ . $Y: y_1 y_2 \dots y_M$ . $X:$
PROMPT-1	$Y: y_1 y_2 \dots y_M$ . $X:$
PROMPT-2	$y_1 y_2 \dots y_M$ . $X:$
PROMPT-3	Translate to $X$ . $Y: y_1 y_2 \dots y_M$ . $X:$
PROMPT-4	Translate from $Y$ to $X$ . $y_1 y_2 \dots y_M$ . $X:$
PROMPT-5	Translate to $X$ . $y_1 y_2 \dots y_M$ . $X:$

Table 9: **MT prompt variants investigated for fine-tuned models.** These MT prompts are following the design in Zhang et al. (2023a).

mediary pivot for translations between non-English language pairs. Our findings indicated that employing the English pivot typically yielded an enhancement of approximately 1.1 - 1.2 BLEU points compared to direct translations in zero-shot directions when fine-tuning BLOOMZ. When contrasting the MTInstruct baseline with our proposed MT+Align, we observed that combining AlignInstruct consistently boosted performance in pivot translation scenarios.

## I Per Language Result Details of BLOOMZ+24 and BLOOMZ+3

We present per language detailed results of original BLOOMZ-7b1 and fine-tuned BLOOMZ-7b1 models in Tab. 12, 13, 14, 15, 16, 17, 18, 19, respectively for the BLOOMZ+24 and BLOOMZ+3 settings.

Prompt	Objective	en→xx			xx→en		
		BLEU	chrF++	COMET	BLEU	chrF++	COMET
PROMPT-default	MTInstruct	11.54	25.33	64.54	18.59	33.25	69.84
	MT+Align	<b>12.28</b>	<b>26.17</b>	<b>65.54</b>	<b>18.72</b>	<b>34.02</b>	<b>70.69</b>
PROMPT-1	MTInstruct	5.29	11.31	50.20	7.87	20.08	57.46
	MT+Align	<b>5.30</b>	<b>11.38</b>	<b>50.95</b>	<b>8.93</b>	<b>20.77</b>	<b>58.38</b>
PROMPT-2	MTInstruct	2.20	6.68	45.56	7.15	19.08	57.22
	MT+Align	1.91	5.35	43.84	<b>7.61</b>	18.80	56.76
PROMPT-3	MTInstruct	10.59	22.69	62.65	15.85	29.93	67.59
	MT+Align	9.20	20.80	60.96	<b>16.17</b>	<b>30.58</b>	<b>68.70</b>
PROMPT-4	MTInstruct	8.67	20.73	61.50	15.20	28.95	66.61
	MT+Align	<b>8.91</b>	20.53	<b>61.64</b>	<b>16.25</b>	<b>30.67</b>	<b>67.94</b>
PROMPT-5	MTInstruct	6.61	14.55	55.99	10.88	22.41	61.40
	MT+Align	6.02	12.28	52.42	<b>11.83</b>	<b>23.85</b>	<b>62.09</b>

Table 10: **Results of using different MT prompts for BLOOMZ-7b1 fine-tuned models during inference.** Refer to Tab. 9 for details about definitions of different MT prompts. We report the average results for the BLOOMZ+24 setting. Results better than the MTInstruct baseline are marked in **bold**.

MTInstruct	BLEU	chrF++	COMET	MT+Align	BLEU	chrF++	COMET
overall	11.79	26.36	63.22	overall	<b>12.13</b>	<b>26.65</b>	<b>63.23</b>
seen→seen	22.68	35.32	76.39	seen→seen	<b>23.67</b>	<b>36.53</b>	<b>76.89</b>
seen→unseen	7.10	24.50	55.18	seen→unseen	<b>7.27</b>	24.32	54.96
unseen→seen	12.56	24.74	68.83	unseen→seen	<b>12.92</b>	<b>25.29</b>	<b>69.10</b>
unseen→unseen	6.78	22.62	53.69	unseen→unseen	6.68	22.30	53.19
MTInstruct with English pivot	BLEU	chrF++	COMET	MT+Align with English pivot	BLEU	chrF++	COMET
overall	12.99	28.01	65.38	overall	<b>13.25</b>	<b>28.30</b>	<b>65.57</b>
seen→seen	23.10	35.30	76.30	seen→seen	<b>23.48</b>	<b>35.57</b>	<b>76.43</b>
seen→unseen	9.00	27.67	59.54	seen→unseen	<b>9.28</b>	<b>28.03</b>	<b>59.73</b>
unseen→seen	13.18	24.98	68.77	unseen→seen	<b>13.36</b>	<b>25.22</b>	<b>68.94</b>
unseen→unseen	8.57	25.77	58.17	unseen→unseen	<b>8.83</b>	<b>26.07</b>	<b>58.42</b>

Table 11: **Results of BLOOMZ+3 using English as a pivot language for zero-shot translation evaluation.** Results of MT+Align surpassing corresponding those of MTInstruct are marked in **bold**.

Language	OPUS en→xx			OPUS xx→en			Flores en→xx			Flores xx→en		
	BLEU	chrF++	COMET	BLEU	chrF++	COMET	BLEU	chrF++	COMET	BLEU	chrF++	COMET
af	3.8	13.2	56.38	7.6	22.0	59.14	2.6	14.9	33.60	20.1	38.0	65.61
am	0.1	0.3	33.17	0.5	8.3	43.57	0.3	0.6	30.65	1.9	12.6	46.24
be	4.2	5.1	47.26	7.3	17.5	48.57	0.4	3.3	31.58	4.2	22.3	49.27
cy	2.7	10.5	53.21	6.2	16.0	53.25	1.2	11.2	34.17	6.0	20.3	53.45
ga	1.2	10.6	42.85	4.0	16.4	46.05	1.2	11.6	33.94	5.5	19.6	46.97
gd	9.3	16.0	51.40	47.6	55.9	59.30	1.2	11.2	36.28	4.2	18.8	43.73
gl	4.5	25.6	64.93	17.2	36.7	66.07	13.4	38.5	74.77	51.0	67.8	85.77
ha	0.1	5.4	38.42	0.3	11.2	42.58	1.5	10.2	35.77	6.9	18.9	47.37
ka	0.3	1.9	31.97	0.6	9.2	44.48	0.4	1.4	28.81	2.4	17.0	47.57
kk	4.3	4.9	50.51	5.1	14.2	51.51	0.5	1.6	33.66	5.1	19.8	51.40
km	2.8	4.5	51.68	3.9	11.1	50.40	0.8	2.9	39.56	5.6	16.2	50.42
ky	10.0	10.6	54.23	10.3	24.0	55.99	0.6	1.6	30.19	3.8	17.9	48.05
li	6.6	16.2	-	5.9	24.8	-	2.0	14.9	-	9.8	29.8	-
my	1.8	2.4	45.44	3.0	5.0	48.33	0.4	0.8	29.58	1.0	3.7	44.15
nb	5.8	18.2	57.01	13.9	33.0	56.37	3.9	19.3	46.74	19.8	40.3	63.56
nn	6.3	18.6	62.33	8.9	25.3	56.28	3.7	19.7	41.75	16.9	37.5	62.37
oc	6.0	13.6	-	5.1	18.6	-	9.6	33.6	-	53.0	68.5	-
si	0.6	2.0	41.84	1.6	9.4	48.58	0.5	1.4	28.08	1.6	9.1	42.67
tg	0.4	1.4	-	1.1	11.8	-	0.4	1.5	-	3.3	18.0	-
tk	7.9	10.6	-	5.3	13.0	-	0.7	8.7	-	4.2	20.1	-
tt	0.0	1.0	-	0.2	13.3	-	0.3	1.4	-	4.2	20.2	-
ug	0.0	0.4	32.44	0.3	11.2	45.69	0.3	0.9	31.34	3.0	16.5	48.99
uz	0.7	2.1	35.94	1.0	12.8	41.86	1.5	11.5	40.65	3.1	18.7	49.43
yi	7.3	16.5	57.47	4.0	23.0	63.91	0.7	1.7	33.22	2.1	15.6	41.87
avg.	3.61	8.82	47.81	6.70	18.49	51.68	2.00	9.35	36.54	9.95	24.47	52.05

Table 12: Detailed results of BLOOMZ-7b1 without fine-tuning.

Language	OPUS en→xx			OPUS xx→en			Flores en→xx			Flores xx→en		
	BLEU	chrF++	COMET	BLEU	chrF++	COMET	BLEU	chrF++	COMET	BLEU	chrF++	COMET
af	25.0	41.4	71.05	38.5	52.3	78.94	10.1	31.0	45.42	33.9	51.1	72.66
am	3.0	12.8	59.55	3.4	19.8	59.71	0.2	5.2	42.97	1.4	16.0	49.47
be	8.9	14.9	55.16	14.0	24.9	62.37	0.7	12.3	30.90	3.7	21.0	49.99
cy	20.2	38.0	71.55	33.2	49.3	77.72	5.0	20.3	38.38	13.1	30.2	57.47
ga	15.6	37.1	63.87	29.2	49.1	75.94	3.7	21.2	39.17	12.5	30.3	57.53
gd	13.1	24.7	62.14	66.0	69.6	77.70	2.2	19.6	40.75	7.1	22.3	50.05
gl	16.9	37.6	70.62	24.7	43.6	75.62	21.9	45.2	77.26	46.6	64.5	86.86
ha	12.3	32.7	71.75	10.0	29.8	64.51	1.9	17.1	49.24	6.8	22.1	48.81
ka	4.6	18.1	67.39	10.0	24.3	60.50	0.3	6.8	27.46	1.5	14.9	46.10
kk	12.6	19.5	66.07	14.6	28.2	71.80	0.8	13.0	35.76	3.9	19.7	52.24
km	19.7	25.2	63.24	13.9	32.1	75.02	0.5	12.3	35.60	6.2	22.4	56.45
ky	16.0	20.5	66.27	21.1	33.8	73.06	0.9	12.7	36.10	3.0	17.5	50.40
li	13.5	32.8	-	21.3	35.7	-	3.3	19.9	-	14.6	31.4	-
my	6.2	14.3	58.04	5.2	15.6	63.65	0.2	12.9	40.37	1.3	12.7	48.38
nb	12.7	30.4	63.27	22.2	42.1	76.74	7.9	28.4	44.15	25.6	44.3	72.56
nn	18.3	38.0	77.18	27.1	47.7	81.80	7.3	25.7	45.35	24.3	42.9	70.06
oc	10.0	20.0	-	13.4	27.1	-	8.0	27.5	-	46.9	63.5	-
si	5.2	21.4	68.16	11.5	26.4	70.79	0.9	12.9	41.73	3.7	19.2	57.41
tg	5.5	22.0	-	8.0	25.9	-	1.1	15.8	-	3.1	19.6	-
tk	24.4	26.7	-	30.4	37.8	-	0.7	10.8	-	3.9	18.8	-
tt	1.9	17.6	-	3.6	19.6	-	0.4	13.7	-	1.6	14.3	-
ug	1.2	19.7	49.76	4.2	21.2	61.34	0.4	12.9	35.88	1.7	16.7	50.29
uz	3.1	18.2	62.12	5.7	22.0	61.12	0.5	3.6	34.67	3.9	18.8	50.32
yi	7.1	24.3	59.13	14.9	20.2	58.66	0.3	9.5	29.77	2.5	17.2	43.27
avg.	11.54	25.33	64.54	18.6	33.25	68.84	3.30	17.10	40.58	11.37	27.14	56.33

Table 13: Detailed results of BLOOMZ-7b1 fine-tuned with MTInstruct for BLOOMZ+24.

Language	OPUS en→xx			OPUS xx→en			Flores en→xx			Flores xx→en		
	BLEU	chrF++	COMET	BLEU	chrF++	COMET	BLEU	chrF++	COMET	BLEU	chrF++	COMET
af	25.0	41.9	70.72	36.9	52.2	78.68	10.6	31.9	45.84	33.5	51.1	72.84
am	3.4	13.2	60.62	4.9	22.8	62.43	0.3	5.4	44.20	1.4	16.4	51.05
be	8.3	14.5	55.23	13.9	25.1	62.72	0.8	12.5	30.93	3.6	20.6	49.14
cy	20.6	39.0	71.73	33.8	49.4	77.55	4.7	20.3	38.70	14.6	31.5	58.34
ga	17.6	39.3	65.76	32.6	52.7	77.49	3.4	21.4	39.99	13.6	31.6	58.73
gd	15.6	27.2	62.09	48.1	55.4	75.90	2.3	20.3	40.81	7.4	22.0	49.99
gl	17.1	37.2	70.85	24.4	43.3	75.90	21.7	44.9	77.09	45.6	63.5	86.60
ha	14.6	35.0	73.34	11.4	31.3	65.69	1.9	17.3	50.88	7.4	22.5	49.57
ka	4.9	18.9	67.54	10.5	25.3	61.27	0.3	6.9	27.61	2.1	16.0	47.04
kk	12.3	19.3	65.73	15.6	28.0	71.01	0.9	13.0	35.86	4.1	19.8	52.43
km	20.4	26.5	63.38	14.4	35.2	75.62	0.6	12.5	35.44	7.1	22.9	57.81
ky	15.8	19.6	64.74	23.3	35.8	74.70	0.9	13.3	36.71	2.9	17.4	50.06
li	13.2	29.4	-	22.3	38.2	-	3.1	19.7	-	12.5	28.7	-
my	7.6	15.4	58.84	6.3	18.0	66.45	0.3	13.3	40.97	1.2	14.4	50.79
nb	13.5	31.4	64.08	24.2	44.2	77.58	7.9	28.7	44.12	25.5	44.9	72.72
nn	19.0	38.0	77.61	28.5	47.7	81.68	7.0	26.7	46.14	25.8	44.1	70.55
oc	9.1	19.3	-	13.5	27.5	-	7.5	25.9	-	47.3	63.8	-
si	5.1	22.1	69.60	13.9	29.1	72.51	1.1	13.1	43.01	5.6	22.7	61.89
tg	6.6	23.7	-	8.8	27.2	-	0.9	15.6	-	3.4	19.9	-
tk	27.2	26.2	-	31.2	38.7	-	0.7	11.4	-	3.8	18.2	-
tt	2.1	18.6	-	5.0	21.5	-	0.4	13.3	-	1.5	13.7	-
ug	1.1	20.7	51.12	5.5	23.4	63.42	0.4	13.8	37.51	2.1	16.3	50.45
uz	3.5	18.6	62.09	7.4	23.3	62.01	0.2	1.9	34.50	3.7	18.2	50.09
yi	11.1	33.1	70.13	12.8	21.2	60.47	0.4	9.8	30.08	2.6	17.0	42.57
avg.	12.28	26.17	65.54	18.72	34.02	70.69	3.26	17.20	41.07	11.60	27.38	56.98

Table 14: Detailed results of BLOOMZ-7b1 fine-tuned with MT+Align for BLOOMZ+24.

Zero-shot	BLEU	chrF++	COMET	Supervised	BLEU	chrF++	COMET
ar-de	1.4	14.8	56.19	en-ar	11.1	32.4	75.66
ar-fr	21.9	46.1	74.19	en-de	12.2	29.2	59.16
ar-nl	0.6	11.2	56.59	en-fr	26.8	49.2	77.42
ar-ru	3.1	6.2	48.41	en-nl	2.0	16.0	46.52
ar-zh	18.4	14.4	73.65	en-ru	5.7	16.1	49.00
de-ar	2.0	17.8	64.91	en-zh	22.5	17.0	77.90
de-fr	12.0	33.4	63.45	avg.	13.38	26.65	64.28
de-nl	3.7	17.9	47.30				
de-ru	1.3	11.8	45.53				
de-zh	8.9	7.6	61.52				
fr-ar	11.2	33.4	74.20		BLEU	chrF++	COMET
fr-de	4.6	23.4	48.83	ar-en	26.7	48.4	78.12
fr-nl	2.8	17.2	52.14	de-en	21.1	38.5	71.99
fr-ru	3.1	10.4	45.12	fr-en	27.7	49.8	79.46
fr-zh	20.9	17.0	76.20	nl-en	12.3	31.1	61.29
nl-ar	1.3	13.2	59.46	ru-en	17.9	36.6	68.40
nl-de	5.9	22.8	46.49	zh-en	24.5	47.9	77.08
nl-fr	9.6	29.6	58.30	avg.	21.70	42.05	72.72
nl-ru	0.8	9.0	42.83				
nl-zh	3.3	3.7	53.96				
ru-ar	6.5	25.3	68.38				
ru-de	2.0	17.0	48.06				
ru-fr	15.7	38.7	67.54				
ru-nl	0.5	10.5	46.14				
ru-zh	10.7	11.3	67.18				
zh-ar	8.6	29.7	73.47				
zh-de	1.6	17.6	49.90				
zh-fr	20.7	44.1	75.79				
zh-nl	0.6	10.4	48.53				
zh-ru	2.9	8.6	44.13				
avg.	6.89	19.14	57.95				
seen→seen	16.95	30.78	74.58	en→seen	20.13	32.87	76.99
seen→unseen	2.30	13.31	49.98	en→unseen	6.63	20.43	51.56
unseen→seen	7.78	20.07	62.74	seen→en	26.30	48.70	78.22
unseen→unseen	2.37	14.83	46.06	unseen→en	17.10	35.40	67.23

Table 15: Detailed results of BLOOMZ-7b1 without fine-tuning.

Zero-shot	BLEU	chrF++	COMET	Supervised	BLEU	chrF++	COMET
ar-de	4.7	20.9	56.43	en-ar	9.1	27.2	71.47
ar-fr	20.8	42.5	71.47	en-de	19.8	36.1	66.53
ar-nl	7.2	22.9	58.29	en-fr	23.0	44.5	74.98
ar-ru	5.0	21.0	54.73	en-nl	15.5	36.1	64.76
ar-zh	14.0	12.4	67.94	en-ru	14.2	30.3	62.82
de-ar	2.4	16.2	64.53	en-zh	20.7	17.9	74.97
de-fr	11.9	31.2	64.44	avg.	17.05	32.02	69.26
de-nl	9.4	28.1	54.22				
de-ru	5.1	19.6	55.41				
de-zh	4.2	5.8	55.26				
fr-ar	10.1	29.1	70.72		BLEU	chrF++	COMET
fr-de	8.6	27.7	53.77	ar-en	26.5	46.9	76.92
fr-nl	10.3	30.1	57.55	de-en	27.0	44.0	76.97
fr-ru	7.9	26.0	56.82	fr-en	27.5	49.0	78.80
fr-zh	18.1	18.5	72.24	nl-en	21.8	41.3	73.99
nl-ar	2.0	15.1	63.73	ru-en	24.8	43.6	74.23
nl-de	9.7	28.1	52.58	zh-en	23.2	45.3	76.83
nl-fr	13.2	32.3	65.17	avg.	25.13	45.02	76.29
nl-ru	5.1	18.6	55.13				
nl-zh	3.0	5.4	54.34				
ru-ar	5.9	15.0	60.36				
ru-de	5.6	23.8	52.66				
ru-fr	17.9	38.4	68.66				
ru-nl	6.2	22.5	54.41				
ru-zh	7.5	13.6	61.40				
zh-ar	6.7	22.1	67.48				
zh-de	3.3	19.6	51.75				
zh-fr	17.4	38.9	73.00				
zh-nl	4.8	19.3	54.41				
zh-ru	3.5	17.9	49.02				
avg.	8.38	22.75	59.93				
seen→seen	14.52	27.25	70.48	en→seen	17.60	29.87	73.81
seen→unseen	6.14	22.82	54.75	en→unseen	16.50	34.17	64.70
unseen→seen	7.56	19.22	61.99	seen→en	25.73	47.07	77.52
unseen→unseen	6.85	23.45	54.07	unseen→en	24.53	42.97	75.06

Table 16: Detailed results of BLOOMZ-7b1 fine-tuned with MTInstruct for BLOOMZ+3 de-nl-ru.

Zero-shot	BLEU	chrF++	COMET	Supervised	BLEU	chrF++	COMET
ar-de	5.1	20.8	55.25	en-ar	8.4	26.0	70.45
ar-fr	20.3	42.5	71.78	en-de	21.1	36.7	67.15
ar-nl	6.4	21.6	57.48	en-fr	22.9	44.4	74.67
ar-ru	5.2	21.5	55.51	en-nl	16.1	36.8	65.26
ar-zh	16.0	14.1	69.55	en-ru	15.2	31.5	63.30
de-ar	2.4	16.3	64.01	en-zh	16.1	15.0	71.93
de-fr	13.5	34.3	66.25	avg.	16.63	31.73	68.79
de-nl	9.7	28.0	55.00				
de-ru	5.3	19.6	55.61				
de-zh	7.2	7.3	60.64				
fr-ar	10.0	28.2	69.86		BLEU	chrF++	COMET
fr-de	9.2	27.8	54.03	ar-en	27.1	47.0	76.54
fr-nl	10.8	31.0	58.50	de-en	27.8	44.4	77.57
fr-ru	8.6	26.7	57.07	fr-en	27.1	48.7	78.82
fr-zh	15.9	15.8	70.78	nl-en	22.6	42.2	74.25
nl-ar	2.2	15.4	63.47	ru-en	25.6	44.2	74.46
nl-de	10.2	28.5	53.65	zh-en	23.5	45.7	77.04
nl-fr	14.4	34.4	66.55	avg.	25.62	45.37	76.45
nl-ru	5.3	19.3	55.53				
nl-zh	5.5	6.2	58.77				
ru-ar	6.5	16.0	62.69				
ru-de	6.1	24.3	52.89				
ru-fr	18.2	39.0	69.95				
ru-nl	6.3	22.5	54.36				
ru-zh	7.6	13.3	61.94				
zh-ar	8.7	26.5	70.88				
zh-de	3.0	19.5	50.82				
zh-fr	17.7	39.7	73.56				
zh-nl	4.4	19.3	54.20				
zh-ru	4.1	19.5	50.47				
avg.	8.86	23.30	60.70				
seen→seen	14.77	27.80	71.07	en→seen	15.80	28.47	72.35
seen→unseen	6.31	23.08	54.81	en→unseen	17.47	35.00	65.24
unseen→seen	8.61	20.24	63.81	seen→en	25.90	47.13	77.47
unseen→unseen	7.15	23.70	54.51	unseen→en	25.33	43.60	75.43

Table 17: Detailed results of BLOOMZ-7b1 fine-tuned with MT+Align for BLOOMZ+3 de-nl-ru.

Zero-shot	BLEU	chrF++	COMET	Supervised	BLEU	chrF++	COMET
ar-de	6.9	24.7	58.10	en-ar	14.6	35.6	76.70
ar-fr	26.2	48.2	74.96	en-de	20.4	36.0	65.96
ar-nl	8.8	24.7	59.53	en-fr	27.9	50.0	77.65
ar-ru	6.5	22.7	55.33	en-nl	14.8	34.8	63.11
ar-zh	28.6	22.3	77.64	en-ru	13.3	29.0	61.43
de-ar	3.3	19.8	68.27	en-zh	36.1	27.7	80.31
de-fr	15.2	35.8	67.05	avg.	21.18	35.52	70.86
de-nl	8.2	26.0	53.35				
de-ru	4.4	17.9	54.79				
de-zh	12.0	9.9	65.20				
fr-ar	14.2	35.2	74.84		BLEU	chrF++	COMET
fr-de	8.9	28.4	53.81	ar-en	33.7	53.5	79.81
fr-nl	10.1	29.9	56.92	de-en	27.1	43.9	77.04
fr-ru	8.1	26.0	55.96	fr-en	29.6	51.0	79.60
fr-zh	30.2	25.6	79.43	nl-en	22.0	41.4	73.54
nl-ar	3.1	18.2	67.72	ru-en	25.1	43.9	74.05
nl-de	10.4	27.7	52.67	zh-en	32.6	54.3	79.75
nl-fr	16.9	37.3	68.46	avg.	28.35	48.00	77.30
nl-ru	4.8	17.8	54.71				
nl-zh	8.1	7.0	63.96				
ru-ar	11.9	31.5	72.45				
ru-de	6.1	23.7	52.74				
ru-fr	21.2	42.5	71.71				
ru-nl	6.8	22.6	53.91				
ru-zh	21.3	20.7	74.63				
zh-ar	13.1	34.1	74.92				
zh-de	4.1	22.3	52.13				
zh-fr	23.8	46.5	76.54				
zh-nl	4.8	19.9	54.26				
zh-ru	5.7	21.9	50.60				
avg.	11.79	26.36	63.22				
seen→seen	22.68	35.32	76.39	en→seen	26.20	37.77	78.22
seen→unseen	7.10	24.50	55.18	en→unseen	16.17	33.27	63.50
unseen→seen	12.56	24.74	68.83	seen→en	31.97	52.93	79.72
unseen→unseen	6.78	22.62	53.69	unseen→en	24.73	43.07	74.88

Table 18: Detailed results of BLOOMZ-7b1 fine-tuned with MTInstruct for BLOOMZ+3 ar-de-fr-nl-ru-zh.

Zero-shot	BLEU	chrF++	COMET	Supervised	BLEU	chrF++	COMET
ar-de	6.7	24.2	57.45	en-ar	15.1	35.8	76.76
ar-fr	27.5	49.2	75.21	en-de	20.6	35.9	65.88
ar-nl	8.7	24.8	59.14	en-fr	27.5	49.4	77.46
ar-ru	6.7	21.6	55.04	en-nl	15.0	35.6	63.70
ar-zh	30.1	24.4	78.54	en-ru	13.5	29.5	61.62
de-ar	3.5	19.7	68.39	en-zh	36.3	27.7	80.52
de-fr	15.4	35.8	67.81	avg.	21.33	35.65	70.99
de-nl	9.6	27.3	53.74				
de-ru	4.7	17.9	54.23				
de-zh	12.0	9.9	65.40				
fr-ar	14.9	36.3	74.98		BLEU	chrF++	COMET
fr-de	9.2	28.3	52.96	ar-en	33.9	53.7	79.74
fr-nl	11.3	31.1	57.62	de-en	27.1	43.6	77.13
fr-ru	8.8	26.2	56.31	fr-en	29.7	51.0	80.03
fr-zh	31.1	26.9	79.93	nl-en	22.6	42.3	73.94
nl-ar	3.3	18.5	68.02	ru-en	25.8	44.5	74.07
nl-de	9.4	26.5	52.33	zh-en	32.5	54.5	80.01
nl-fr	17.2	37.3	68.38	avg.	28.60	48.27	77.49
nl-ru	4.4	17.1	53.63				
nl-zh	8.3	7.0	64.08				
ru-ar	12.4	32.1	72.40				
ru-de	5.7	22.9	51.90				
ru-fr	21.5	42.7	72.08				
ru-nl	6.3	22.1	53.32				
ru-zh	22.7	24.6	75.36				
zh-ar	13.9	35.4	75.68				
zh-de	3.6	21.3	51.32				
zh-fr	24.5	47.0	76.98				
zh-nl	4.9	20.3	54.30				
zh-ru	5.5	21.1	50.49				
avg.	12.13	26.65	63.23				
seen→seen	23.67	36.53	76.89	en→seen	26.30	37.63	78.25
seen→unseen	7.27	24.32	54.96	en→unseen	16.37	33.67	63.73
unseen→seen	12.92	25.29	69.10	seen→en	32.03	53.07	79.93
unseen→unseen	6.68	22.30	53.19	unseen→en	25.17	43.47	75.05

Table 19: Detailed results of BLOOMZ-7b1 fine-tuned with MT+Align for BLOOMZ+3 ar-de-fr-nl-ru-zh.