Direct Preference Optimization for Neural Machine Translation with Minimum Bayes Risk Decoding

Anonymous ACL submission

Abstract

Minimum Bayes Risk (MBR) decoding can significantly improve translation performance of Multilingual Large Language Models (MLLMs). However, MBR decoding is computationally expensive. We show how the recently developed Reinforcement Learning technique, Direct Preference Optimization (DPO), can fine-tune MLLMs to get the gains of MBR without any additional computation in inference. Our method uses only a small monolingual fine-tuning set and yields significantly improved performance on multiple NMT test sets compared to MLLMs without DPO.

1 Introduction

002

007

011

017

036

037

MBR decoding (Kumar and Byrne, 2004; Eikema and Aziz, 2022; Suzgun et al., 2023) is a twopass procedure that generates multiple translation hypotheses and selects a hypothesis based on Bayesian risk. Recent work (Garcia et al., 2023; Suzgun et al., 2023; REDACTED, 2023) has shown that MBR decoding can significantly boost the translation performance of MLLMs (Lin et al., 2022; Muennighoff et al., 2023; Zeng et al., 2023a), outperforming greedy decoding and beam search. However, MBR decoding is expensive, both in computation and in latency.

Our goal is to fine-tune a base MLLM so that it has the same single-pass decoding performance as MBR decoding. We propose¹ a novel self-supervised fine-tuning method based on DPO (Rafailov et al., 2023). Our method uses MBR decoding on an MLLM to produce a preference dataset consisting of pairs of ranked translations. The DPO algorithm is used to fine-tune the MLLM to prefer the higher-ranked translations over lower-ranked ones. MLLMs optimized for MBR preference achieve significantly better translation performance when decoded with beam search, achieving translation quality on par with MBR decoding of the original model.

038

040

041

043

045

047

051

053

054

060

061

062

063

064

065

067

068

2 MBR and DPO

We follow the expectation-by-sampling approach to MBR (Eikema and Aziz, 2022). Given a set of sampled translations $H(\mathbf{x}) = \{\mathbf{y}' \sim P(\cdot|\mathbf{x})\}$ and a loss (or utility) function $L(\cdot, \cdot)$, the score (negative Bayes risk) of each translation is found as

$$S(\mathbf{y}) = -\frac{1}{|H(x)|} \sum_{\mathbf{y}' \in H(\mathbf{x})} L(\mathbf{y}', \mathbf{y}) \qquad (1)$$

and the MBR hypothesis is then computed as

$$\mathbf{y}^* = \underset{\mathbf{y} \in H(\mathbf{x})}{\operatorname{arg\,max}} S(\mathbf{y}) \tag{2}$$

This is simple but expensive. Our goal is to train a model that produces translations with scores consistent with MBR, but without multi-step decoding.

2.1 DPO Fine-Tuning Objective

DPO (Rafailov et al., 2023) reformulates the usual approach to Reinforcement Learning from Human Feedback (RLHF) so as to avoid a distinct reward modelling step. The typical RLHF criteria is

$$\max_{\pi_{\theta}} \mathbb{E}_{\mathbf{x} \sim D, \mathbf{y} \sim \pi_{\theta}(\mathbf{y}|\mathbf{x})} [r_{\phi}(\mathbf{x}, \mathbf{y})]$$
(3)

$$-\beta \mathbb{D}_{KL} \left[\pi_{\theta}(\mathbf{y}|\mathbf{x}) \parallel \pi_{\text{ref}}(\mathbf{y}|\mathbf{x}) \right]$$

where r_{ϕ} is a reward model trained from human feedback, π_{θ} is the model being trained, and π_{ref} is the reference model. DPO effectively replaces the reward model with a preference distribution based on π_{θ} , the model being trained; DPO also retains the KL regularization term with weighting β .

The preference dataset D for DPO consists of triplets $(\mathbf{x}, \mathbf{y}_w, \mathbf{y}_l)$ where \mathbf{x} is the input prompt, \mathbf{y}_w is the winnng (prefered) response, and \mathbf{y}_l is the

¹Reviewers: This work originally appeared in an MPhil dissertation submitted in August 2023, redacted for anonymity.

144

145

146

147

148

149

150

151

152

153

154

155

156

157

158

159

161

114

115

116

117

losing (disprefered) response. DPO uses the language model likelihood to approximate the reward as $\beta \log \frac{\pi_{\theta}(\mathbf{y}|\mathbf{x})}{\pi_{\text{ref}}(\mathbf{y}|\mathbf{x})}$. During training, with π_{θ} typically initialized from π_{ref} , the objective is to maximize the expected reward margin between \mathbf{y}_w and \mathbf{y}_l :

$$L_{\text{DPO}} = -\mathbb{E}_{(\mathbf{x}, \mathbf{y}_w, \mathbf{y}_l) \sim D}[\log \sigma(M(\mathbf{y}_w, \mathbf{y}_l, \mathbf{x}, \theta))]$$
(4)

where the reward margin $M(\mathbf{y}_w, \mathbf{y}_l, \mathbf{x}, \theta)$ is

$$\beta \left(\log \frac{\pi_{\theta}(\mathbf{y}_{w}|\mathbf{x})}{\pi_{\text{ref}}(\mathbf{y}_{w}|\mathbf{x})} - \log \frac{\pi_{\theta}(\mathbf{y}_{l}|\mathbf{x})}{\pi_{\text{ref}}(\mathbf{y}_{l}|\mathbf{x})} \right)$$
(5)

2.2 Related Work in Translation

Previous work has explored the effectiveness of enhancing the translation performance of LLMs via Reinforcement Learning (RL) algorithms or supervised fine-tuning. Dong et al. (2023) proposed RAFT that iteratively generates samples and fine-tunes the model on the filtered samples ranked by a reward model. Gulcehre et al. (2023) proposed ReST that uses similar method for translation task, where they apply several fine-tuning steps on a sampled dataset, each time higher ranked samples.

Similar to our pairwise preference learning, Zeng et al. (2023b) introduced a framework TIM to enhance the translation performance of LLMs by learning to compare good translations and bad translations via a preference learning loss.

Contemporaneous with this work, Finkelstein et al. (2023) proposed MBR fine-tuning, which finetunes an NMT model on the MBR decoding outputs generated either by the model itself or by an LLM. However, their MBR fine-tuning utilizes only the final translations of MBR decoding whereas our fine-tuning method uses sets of sampled translations ranked by MBR, thus enabling the model to learn the same ranking preferences as MBR.

3 Methodology

094

097

098

100

101

102

Our method combines MBR decoding and DPO 103 fine-tuning (REDACTED, 2023). We use the MBR procedure to calculate a score (Equation 1) for each 105 of a set of translation hypothesis generated by the 106 base model. We then fine-tune the base model us-107 ing the DPO objective (Equations 4,5) where the 108 109 winning and losing hypotheses provided to DPO are chosen based on their relative MBR scores. If 110 successful, the fine-tuned model will have learned 111 to rank translations consistently with MBR decod-112 ing under the base model. 113

3.1 Creation of the DPO Preference Sets

Following Eikema and Aziz (2022), we use sampling to generate the translation hypotheses that will be used in DPO. For a source sentence \mathbf{x} we use simple ancestral sampling with a temperature of 0.7 to create a set of translations $H(x) = \{\mathbf{y} \sim \pi_{base}(\mathbf{y}|\mathbf{x})\}$ of size |H(x)|. We use this collection as both the MBR evidence and hypothesis spaces (Goel and Byrne, 2000).

The hypotheses in H(x) are ordered by their MBR scores as $\mathbf{y}_1, \mathbf{y}_2, ..., \mathbf{y}_{|H|}$ with the BLEURT metric (Sellam et al., 2020) as the utility function. The ordering reflects the MBR preference, i.e. \mathbf{y}_1 would be the most preferred MBR hypothesis.

Preference Selection Strategies DPO requires a set of preference triplets $\mathcal{D} = \{(\mathbf{x}, \mathbf{y}_w, \mathbf{y}_l)\}$ where \mathbf{y}_w has better MBR score than \mathbf{y}_l and both of the hypotheses are selected from the hypothesis set H(x). There are numerous strategies for selecting the preference pairs $(\mathbf{y}_w, \mathbf{y}_l)$ from the hypothesis set. We experimented with four selection schemes:

- BW is a simple strategy that selects the best and worst translation hypotheses from the ranked sets. For each source sentence x, we only have one preference triplet (x, y1, y|H(x)|).
- BMW adds the middle hypothesis y_m from the ranked lists with index m = ⌈|H(x)|/2⌉. This gives two triplets per source sentence: (x, y₁, y_m) and (x, y_m, y_{|H(x)|}).
- 3. CP selects consecutive pairs from the ranked list, yielding |H(x)| 1 triplets per source sentence, as $(\mathbf{x}, \mathbf{y}_1, \mathbf{y}_2), (\mathbf{x}, \mathbf{y}_2, \mathbf{y}_3), \ldots$
- 4. **CPS** introduces a **stride** into the CP selection strategy so as to avoid requiring DPO to learn distinctions between translations that are similarly ranked. For example, with a stride of 2 we select triplets $(\mathbf{x}, \mathbf{y}_1, \mathbf{y}_3), (\mathbf{x}, \mathbf{y}_3, \mathbf{y}_5), \dots$

3.2 DPO Fine-Tuning

With a set of preference triplets \mathcal{D} selected by one of the schemes above, DPO fine-tuning proceeds as described in Section 2.1 and by Rafailov et al. (2023). The base model serves as the reference model in Equation 4. The base model is also used to initialise π_{θ} , which is the model being fine-tuned. The only DPO hyper-parameter we tune is β , which regulates how the fine-tuned model departs from the reference model Rafailov et al. (2023).



Figure 1: Reward margins for DPO MBR fine-tuning of BLOOMZ and BLOOMZ-mt with BMW and CPS (stride of 2) selection strategies. Margins are calculated on the Zh-En fine-tuning set (WMT20 test set) as finetuning proceeds over one epoch. Results are plotted as moving averages with a window size of 20. CPS yields more preference pairs than BMW.

4 DPO MBR Fine-Tuning and MT

162

163

164

165

166

167

168

171

172

173

174

Datasets: We evaluate translation on the WMT21 news translation test sets (Akhbardeh et al., 2021) and the WMT22 general translation for Chinese-English (Kocmi et al., 2022), and the IWSLT 2017 test set for French-English (Cettolo et al., 2017). For DPO fine-tuning we use the source language text in the WMT20 test sets for Chinese-English (Barrault et al., 2020) and IWSLT 2017 validation sets for French-English. We do not use the corresponding reference translations, as DPO MBR fine-tuning is unsupervised. The fine-tuning and test sets are distinct and do not overlap.

Models: We use the BLOOMZ and BLOOMZ-mt 175 models (Muennighoff et al., 2023) with 7.1 bil-176 lion parameters as our base model. BLOOMZ-mt 177 was pre-trained on 366 billion tokens from mono-178 lingual texts and was fine-tuned for translation 179 task on Flores-200 (NLLB Team et al., 2022) and 180 Tatoeba (Tiedemann, 2020) datasets. To prompt the 181 model for translation, we include two randomly selected translation examples from the fine-tuning set 183 into the input prompt as demonstration examples; these prompts are kept fixed throughout.

Evaluation Metrics: We use two evaluation metrics: BLEU (Papineni et al., 2002) and BLEURT
(Sellam et al., 2020). We follow Garcia et al. (2023)
and use BLEURT as the main evaluation metric.
BLEU serves only as a safety check: DPO finetuning should increase BLEURT with no decrease
in BLEU.

Baselines and Targets: We take the base model
and evaluate it on all the test sets with both beam
search and MBR decoding. Our fine-tuned models,



Figure 2: Reward margin distributions over all preference pairs extracted via the BMW scheme from a held-out dataset (WMT18 Zh-En test). Distributions are gathered over the entire held-out set at model checkpoints at the beginning, a quarter, middle, three quarters, and end of one epoch of DPO fine-tuning. |H| = 8and $\beta = 0.7$. DPO fine-tuning generalises beyond its fine-tuning set and yields improved reward margins on held-out data.

when decoded with beam search, should achieve similar performance as MBR decoding under the base model and show improvement over the base model. We investigate two questions: 196

197

198

199

200

201

202

203

204

205

206

207

208

210

211

212

213

214

215

216

217

218

219

220

221

222

224

225

226

(1) Can DPO teach MLLMs to learn their MBR translation preferences?

(2) Does preference learning with DPO lead to improved translation?

4.1 DPO Fine-Tuning Teaches a MLLM to Learn Its MBR Preferences

Figure 1 shows that the reward margins remain positive and, with some fluctuations, increase as fine-tuning proceeds, for all three models. This suggests that DPO MBR fine-tuned models learn to put more probability mass on the winning hypotheses. The larger the margins, the more the models prefer the winning over the losing hypotheses.

To further investigate DPO MBR fine-tuning, we plot the distribution of reward margins on a heldout set, shown in Figure 2. The median of the distributions increase consistently as fine-tuning proceeds, indicating that the MBR preferences learned in fine-tuning also generalize to unseen data.

4.2 DPO MBR Translation

Table 1 gives our main translation results. Comparing Rows 3 & 4 and 7 & 8, we can see that DPO MBR fine-tuned models, when decoded with beam search, achieve similar performance in BLEURT ($\approx \pm 1$) as the base model decoded with MBR. Both configurations outperform the base model's beam search results by ≈ 3 BLEURT.

#	Model (Decoding)	WMT21		WMT22		IWSLT17	
		zh-en	en-zh	zh-en	en-zh	fr-en	en-fr
1	BLOOMZ (Beam)	59.6 16.4	59.2 22.3	59.9 14.0	55.9 22.2	72.7 38.1	69.3 37.6
2	BLOOMZ (MBR $ H = 8$)	61.9 14.3	62.5 19.7	62.1 11.6	62.7 20.3	73.6 34.2	70.4 32.6
3	BLOOMZ (MBR $ H = 32$)	63.5 15.0	64.7 20.2	64.0 12.4	64.9 21.2	74.8 36.3	72.6 34.1
4	BLOOMZ-DPO-MBR (Beam)	62.3 17.2	62.5 23.7	64.0 15.6	64.2 26.5	76.5 40.6	72.2 38.9
5	BLOOMZ-mt (Beam)	60.3 16.4	59.2 22.5	60.9 14.7	59.0 26.2	74.8 38.7	70.3 37.8
6	BLOOMZ-mt (MBR $ H = 8$)	61.6 13.5	62.6 20.2	63.0 12.2	64.7 23.3	75.4 35.2	71.0 31.8
7	BLOOMZ-mt (MBR $ H = 32$)	63.4 14.3	64.9 20.8	64.8 13.0	66.8 24.0	76.3 36.9	73.2 33.8
8	BLOOMZ-mt-DPO-MBR (Beam)	63.9 18.0	64.0 22.7	65.1 15.9	67.6 26.9	76.5 40.4	71.9 38.3

Table 1: Translation performance in BLEURT and BLEU (BLEURT | BLUE) for models with beam search and MBR decoding on two language pairs from WMT21 news translation test sets, WMT22 general translation test sets, and IWSLT 2017 test sets. DPO-MBR indicates our translation performance with our fine-tuning method. All the DPO MBR models were fine-tuned using the BMW strategy and $\beta = 0.7$. We set |H| = 32 to fine-tune BLOOMZ-mt-DPO-MBR on English-Chinese direction, |H| = 16 on the French-English direction, and set |H| = 8 to fine-tune other DPO MBR models. DPO-MBR improves BLEURT and BLEU in all conditions.

#	β	BLEU	BLEURT	
1	(Baseline)	16.4	60.3	
2	0.1	9.9	64.5	
3	0.3	11.8	64.8	
4	0.5	14.3	64.0	
5	0.7	16.4	63.3	
6	0.9	17.6	61.8	

Selection Strategy	H =8	H =16	H =32
BW	63.3	63.9	63.9
BMW	63.9	64.2	63.6
CPS (strides of 2, 4, and 8)	62.3	63.5	62.9

Table 2: Effect of regularization parameter β for DPO MBR fine-tuning of BLOOMZ using CPS. |H| = 8.

DPO MBR improves the translation ability of BLOOMZ and BLOOMZ-mt across a range of test sets. BLOOMZ-mt shows a notable improvement in BLEURT after DPO MBR fine-tuning, achieving the best performance on four out of six test sets.

4.2.1 KL-Divergence Regularization

227

229

230

237

241

242

243

We investigated the role of β , the KL-divergence regularization factor, in DPO. Table 2 shows that fine-tuning with small β values yields high BLEURT score (exceeding 64), but also a degradation in BLEU (from 16.4 to less than 12). Anecdotally, we find that small values of β lead to repetitive outputs that are penalised heavily under BLEU. Gains in both BLEU and BLEURT are readily found, but we conclude that DPO MBR fine-tuning requires some care in regularization.

4.2.2 Effects of Pair Selection Strategy

244Table 3 shows that models trained on preference245datasets constructed with three different pair se-246lection strategy achieve similar performance on247WMT21 Zh-En, with BLEURT scores in the range24862.9-63.9. DPO MBR appears robust to the se-249lection of preference pairs. However, in terms of

Table 3: WMT21 Zh-En BLEURT scores for BLOOMZ with DPO MBR fine-tuning with different preference pair selection strategies and hypothesis set sizes.

training efficiency, the BW and BMW strategies require fewer preference pairs (1 and 2 per source sentence, resp.) compared to the CP strategy. 250

251

253

254

255

256

257

259

260

261

262

263

264

265

266

267

270

4.2.3 Effects of Size of Hypothesis Set

Table 3 shows that the number of hypotheses needed in the training preference dataset is less than that needed for MBR decoding (Rows 3 & 7 in Table 1). The best performance (BLEURT of 63.9) can be achieved with 16 hypotheses for the BW strategy and 8 hypotheses for the BMW strategy, an improvement over MBR decoding of the base model with |H| = 8 (Row 2 & 6 in Table 1).

5 Conclusion

We introduce DPO MBR fine-tuning, an unsupervised preference optimization algorithm that leverages the ranked lists from MBR decoding to teach MLLMs the preference of MBR decoding. Our method enables MLLMs to achieve significant performance improvement when decoded with beam search in one pass, on par with the performance gained from two-pass MBR decoding².

²Open source code for replicating these experiments will be released via GitHub if accepted

271

274

276

277

282

286

287

290

291

294

298

303

310

311

312

313

315

316

317

319

6 Limitations

Our method was evaluated on WMT 2021 and WMT 2022 and IWSLT 2017 test sets, with highresource languages only (English, Chinese, and French). While our fine-tuned models performed well on these diverse test sets, behaviour may be different on medium-resource or low-resource languages or on other domains.

Our experiments focus on BLOOMZ and BLOOMZ-mt due to the ease of working with them and because BLOOMZ-mt is fine-tuned for translation. Other (M)LLMs may yield different results.

We report MBR results using simple ancestral sampling. Other work (Freitag et al., 2023) has found that there may be advantages in using other sampling schemes, such as epsilon sampling, for MBR. Those other sampling methods potentially offer further gains beyond what we have already shown.

We do not report human assessments of translation quality to verify improvements, but we note that Freitag et al. (2022) have reported extensive results showing that MBR decoding under BLEURT leads to improvements in translation quality as assessed by human judges. We therefore take improvement in BLEURT as our main measurement of improved translation quality.

7 Risks

Our unsupervised fine-tuning technique could potentially amplify undesirable biases or language already present in the baseline systems. This could possibly happen if the MBR utility function, in our case BLEURT, somehow encourages consensus amongst similar translations that are also undesirable. Mitigation should be straightforward, in that any monitoring of the baseline models could also be applied after DPO MBR fine-tuning to reject fine-tuned models that exhibit any increase in bad behaviour. Although it is not a focus of this work, DPO MBR could possibly be used as a strategy for risk mitigation by penalizing undesirable behaviour through introduction of specific penalties into the MBR utility function.

314 References

Farhad Akhbardeh, Arkady Arkhangorodsky, Magdalena Biesialska, Ondřej Bojar, Rajen Chatterjee, Vishrav Chaudhary, Marta R. Costa-jussa, Cristina España-Bonet, Angela Fan, Christian Federmann, Markus Freitag, Yvette Graham, Roman Grundkiewicz, Barry Haddow, Leonie Harter, Kenneth Heafield, Christopher Homan, Matthias Huck, Kwabena Amponsah-Kaakyire, Jungo Kasai, Daniel Khashabi, Kevin Knight, Tom Kocmi, Philipp Koehn, Nicholas Lourie, Christof Monz, Makoto Morishita, Masaaki Nagata, Ajay Nagesh, Toshiaki Nakazawa, Matteo Negri, Santanu Pal, Allahsera Auguste Tapo, Marco Turchi, Valentin Vydrin, and Marcos Zampieri. 2021. Findings of the 2021 conference on machine translation (WMT21). In *Proceedings of the Sixth Conference on Machine Translation*, pages 1–88, Online. Association for Computational Linguistics. 320

321

322

323

324

325

327

328

329

330

331

332

333

334

335

338

339

341

342

343

344

345

346

347

348

349

350

351

352

353

354

355

356

357

358

359

360

361

362

363

364

365

366

367

368

369

370

371

372

373

374

375

376

377

- Loïc Barrault, Magdalena Biesialska, Ondřej Bojar, Marta R. Costa-jussà, Christian Federmann, Yvette Graham, Roman Grundkiewicz, Barry Haddow, Matthias Huck, Eric Joanis, Tom Kocmi, Philipp Koehn, Chi-kiu Lo, Nikola Ljubešić, Christof Monz, Makoto Morishita, Masaaki Nagata, Toshiaki Nakazawa, Santanu Pal, Matt Post, and Marcos Zampieri. 2020. Findings of the 2020 conference on machine translation (WMT20). In *Proceedings of the Fifth Conference on Machine Translation*, pages 1–55, Online. Association for Computational Linguistics.
- Mauro Cettolo, Marcello Federico, Luisa Bentivogli, Jan Niehues, Sebastian Stüker, Katsuhito Sudoh, Koichiro Yoshino, and Christian Federmann. 2017. Overview of the IWSLT 2017 evaluation campaign. In *Proceedings of the 14th International Conference on Spoken Language Translation*, pages 2–14, Tokyo, Japan. International Workshop on Spoken Language Translation.
- Hanze Dong, Wei Xiong, Deepanshu Goyal, Yihan Zhang, Winnie Chow, Rui Pan, Shizhe Diao, Jipeng Zhang, Kashun Shum, and Tong Zhang. 2023. RAFT: Reward ranked finetuning for generative foundation model alignment.
- Bryan Eikema and Wilker Aziz. 2022. Sampling-based approximations to minimum Bayes risk decoding for neural machine translation. In *Proceedings of the 2022 Conference on Empirical Methods in Natural Language Processing*, pages 10978–10993, Abu Dhabi, United Arab Emirates. Association for Computational Linguistics.
- Mara Finkelstein, Subhajit Naskar, Mehdi Mirzazadeh, Apurva Shah, and Markus Freitag. 2023. MBR and QE finetuning: Training-time distillation of the best and most expensive decoding methods.
- Markus Freitag, Behrooz Ghorbani, and Patrick Fernandes. 2023. Epsilon sampling rocks: Investigating sampling strategies for minimum Bayes risk decoding for machine translation.
- Markus Freitag, David Grangier, Qijun Tan, and Bowen Liang. 2022. High quality rather than high model probability: Minimum Bayes risk decoding with neural metrics. *Transactions of the Association for Computational Linguistics*, 10:811–825.

- 378 379

- 398
- 400
- 401 402
- 403 404
- 405
- 406
- 407 408
- 409 410 411
- 412 413 414 415
- 417 418

416

419

- 420 421
- 422 423
- 424
- 425 426

427

- 428
- 429

434 435 436

- Xavier Garcia, Yamini Bansal, Colin Cherry, George Foster, Maxim Krikun, Fangxiaoyu Feng, Melvin Johnson, and Orhan Firat. 2023. The unreasonable effectiveness of few-shot learning for machine translation.
- Vaibhava Goel and William J Byrne. 2000. Minimum Bayes-risk automatic speech recognition. Computer Speech & Language, 14(2):115–135.
- Caglar Gulcehre, Tom Le Paine, Srivatsan Srinivasan, Ksenia Konyushkova, Lotte Weerts, Abhishek Sharma, Aditya Siddhant, Alex Ahern, Miaosen Wang, Chenjie Gu, Wolfgang Macherey, Arnaud Doucet, Orhan Firat, and Nando de Freitas. 2023. Reinforced self-training (ReST) for language modeling.
- Tom Kocmi, Rachel Bawden, Ondřej Bojar, Anton Dvorkovich, Christian Federmann, Mark Fishel, Thamme Gowda, Yvette Graham, Roman Grundkiewicz, Barry Haddow, Rebecca Knowles, Philipp Koehn, Christof Monz, Makoto Morishita, Masaaki Nagata, Toshiaki Nakazawa, Michal Novák, Martin Popel, and Maja Popović. 2022. Findings of the 2022 conference on machine translation (WMT22). In Proceedings of the Seventh Conference on Machine Translation (WMT), pages 1-45, Abu Dhabi, United Arab Emirates (Hybrid). Association for Computational Linguistics.
- Shankar Kumar and William Byrne. 2004. Minimum Bayes-risk decoding for statistical machine translation. In Proceedings of the Human Language Technology Conference of the North American Chapter of the Association for Computational Linguistics: HLT-NAACL 2004, pages 169-176, Boston, Massachusetts, USA. Association for Computational Linguistics.
- Xi Victoria Lin, Todor Mihaylov, Mikel Artetxe, Tianlu Wang, Shuohui Chen, Daniel Simig, Myle Ott, Naman Goyal, Shruti Bhosale, Jingfei Du, Ramakanth Pasunuru, Sam Shleifer, Punit Singh Koura, Vishrav Chaudhary, Brian O'Horo, Jeff Wang, Luke Zettlemoyer, Zornitsa Kozareva, Mona Diab, Veselin Stoyanov, and Xian Li. 2022. Few-shot learning with multilingual generative language models. In Proceedings of the 2022 Conference on Empirical Methods in Natural Language Processing, pages 9019-9052, Abu Dhabi, United Arab Emirates. Association for Computational Linguistics.
- Niklas Muennighoff, Thomas Wang, Lintang Sutawika, Adam Roberts, Stella Biderman, Teven Le Scao, M Saiful Bari, Sheng Shen, Zheng Xin Yong, Hailey Schoelkopf, Xiangru Tang, Dragomir Radev, Alham Fikri Aji, Khalid Almubarak, Samuel Albanie, Zaid Alyafeai, Albert Webson, Edward Raff, and Colin Raffel. 2023. Crosslingual generalization through multitask finetuning. In Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers), pages 15991-16111, Toronto, Canada. Association for Computational Linguistics.

NLLB Team, Marta R. Costa-jussà, James Cross, Onur Celebi, Maha Elbayad, Kenneth Heafield, Kevin Heffernan, Elahe Kalbassi, Janice Lam, Daniel Licht, Jean Maillard, Anna Sun, Skyler Wang, Guillaume Wenzek, Al Youngblood, Bapi Akula, Loic Barrault, Gabriel Mejia Gonzalez, Prangthip Hansanti, John Hoffman, Semarley Jarrett, Kaushik Ram Sadagopan, Dirk Rowe, Shannon Spruit, Chau Tran, Pierre Andrews, Necip Fazil Ayan, Shruti Bhosale, Sergey Edunov, Angela Fan, Cynthia Gao, Vedanuj Goswami, Francisco Guzmán, Philipp Koehn, Alexandre Mourachko, Christophe Ropers, Safiyyah Saleem, Holger Schwenk, and Jeff Wang. 2022. No language left behind: Scaling humancentered machine translation.

437

438

439

440

441

442

443

444

445

446

447

448

449

450

451

452

453

454

455

456

457

458

459

460

461

462

463

464

465

466

467

468

469

470

471

472

473

474

475

476

477

478

479

480

481

482

483

484

485

486

487

488

489

490

491

492

- Kishore Papineni, Salim Roukos, Todd Ward, and Wei-Jing Zhu. 2002. Bleu: a method for automatic evaluation of machine translation. In Proceedings of the 40th Annual Meeting of the Association for Computational Linguistics, pages 311-318, Philadelphia, Pennsylvania, USA. Association for Computational Linguistics. Implemented in SacreBLEU: https://github.com/mjpost/sacrebleu.
- Rafael Rafailov, Archit Sharma, Eric Mitchell, Stefano Ermon, Christopher D. Manning, and Chelsea Finn. 2023. Direct preference optimization: Your language model is secretly a reward model.
- REDACTED. 2023. Redacted. MPhil Dissertation, submitted August 2023.
- Thibault Sellam, Dipanjan Das, and Ankur Parikh. 2020. BLEURT: Learning robust metrics for text generation. In Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics, pages 7881–7892, Online. Association for Computational Linguistics. Implemented in https://github.com/ google-research/bleurt.
- Mirac Suzgun, Luke Melas-Kyriazi, and Dan Jurafsky. 2023. Follow the wisdom of the crowd: Effective text generation via minimum Bayes risk decoding. In Findings of the Association for Computational Linguistics: ACL 2023, pages 4265-4293, Toronto, Canada. Association for Computational Linguistics.
- Jörg Tiedemann. 2020. The Tatoeba Translation Challenge - Realistic data sets for low resource and multilingual MT. In Proceedings of the Fifth Conference on Machine Translation, pages 1174–1182, Online. Association for Computational Linguistics.
- Aohan Zeng, Xiao Liu, Zhengxiao Du, Zihan Wang, Hanyu Lai, Ming Ding, Zhuoyi Yang, Yifan Xu, Wendi Zheng, Xiao Xia, Weng Lam Tam, Zixuan Ma, Yufei Xue, Jidong Zhai, Wenguang Chen, Peng Zhang, Yuxiao Dong, and Jie Tang. 2023a. GLM-130B: An open bilingual pre-trained model.
- Jiali Zeng, Fandong Meng, Yongjing Yin, and Jie Zhou. 2023b. TIM: Teaching large language models to translate with comparison.