Segmenting Beyond Defaults: Asymmetrical Byte Pair Encoding for **Optimal Machine Translation Performance**

Anonymous ACL submission

Abstract

In Machine Translation (MT) research, we often come across recipes recommending a set of fixed hyperparameters to train segmenta-005 tion models to segment words regardless of the amount of text or language pair involved. Although a fixed hyperparameter for the word segmentation model can reduce training resource overhead, we find that using the same number of merge operations (NMO) on both source and target languages - symmetric Byte Pair Encoding (BPE), for different language pairs and text sizes does not guarantee optimal Machine Translation system performance. In 015 this work, we explore and identify BPE segmentation recipes across various data sizes and language pairs to obtain optimal performance. We find that using *asymmetric* BPE improves results compared to symmetric BPE, particularly in low-resource scenarios (50K, 100K, 500K) by (5.32,4.46,0.7) CHRF++ scores (with p < 0.05) on average for English-Hindi. We further validate our findings on the other six pairs, English↔Telugu, Shona, Norwegian, Kyrgyz, Hausa and Inuktitut, to show the consistency of this work. A statistically significant improvement is observed using asymmetric BPE config-028 urations in 10 of 12 systems when comparing symmetric BPE configurations. Our findings indicate that using a high NMO for the source (4K to 32K) and a low NMO (0.5K to 2K)provides optimal results, particularly in lowresource contexts.

1 Introduction

004

007

012

017

027

034

There have been commendable and considerable efforts in the recent past to bring low resource language pairs into the fold of Neural Machine Translation paradigm (e.g. Workshop on Technologies for MT of Low Resource Languages). However, in this process, it is common to apply successful configurations or methodologies, such as hyperparameters for data preprocessing pipelines 042

and neural networks from past work, without suffi-043 cient consideration of their suitability for the specific language pair in question. For example, if 045 we take a preprocessing step, such as word segmentation, it plays a vital role by dividing words 047 into "subwords" to improve learning. In machine translation, word segmentation/subword tokenisa-049 tion reduces vocabulary size, manages rare and unknown words, and improves MT performance. No-051 table methods include BPE (Sennrich et al., 2016), word piece (Devlin et al., 2019), sentence piece (Kudo and Richardson, 2018), and morfessor (Smit 054 et al., 2014). BPE compresses the data by merging 055 frequent character pairs into new symbols (Gage, 1994) and thus the resultant subword tokenisation 057 model has number of merge operations (NMO) as a hyperparameter. The NMO determines the degree of word segmentation: a lower NMO, such 060 as 500 (see Table 1), leads to smaller vocabularies 061 and more segmentation, whereas a higher NMO, 062 such as 32K, creates larger vocabularies with mini-063 mal segmentation (segmenting mostly rare or un-064 known words). More often than not, same num-065 ber of merge operations are performed on both 066 the source and the target languages. Most of the 067 work using BPE as preprocessing step has taken 068 the same NMO on source and target for various lan-069 guage pairs without taking the language pair into 070 account or the amount of data involved. Recent 071 research highlights the importance of examining 072 BPE hyperparameters in MT systems, particularly 073 in low-resource scenarios (Ding et al., 2019; Abid, 074 2020). However, there is still uniformity in terms 075 of having the same NMO for both the source and the target (symmetrical BPE) (Huck et al., 2017; 077 Ortega et al., 2020; Lankford et al., 2021; Domingo 078 et al., 2023; Lee et al., 2024) and very little explo-079 ration has been done on asymmetrical BPE for MT. Work like Ngo Ho and Yvon (2021) explored asymmetric BPE for alignment between languages, but not for MT.

In the context of MT, we define the "BPE config-084 uration" as $m_1 m_2$ where m_1 and m_2 are the number of merge operations (NMO) for the source and 086 target languages, respectively. Our work explores the effects of symmetric and asymmetric BPE configurations for English-Hindi language pair in different dataset availability scenarios and shows that 090 we achieve best performance using asymmetric configuration. We further validate our findings on other six language pairs English \leftrightarrow {Telugu, Shona, Norwegian, Kyrgyz, Hausa, and Inuktitut } across different language families to show the consistency of this work. These languages were chosen to represent different typology and morphological structures. We find that, depending on data availability and the language pair involved, the optimal BPE configuration is likely to be asymmetric. Specifi-100 cally, for low-resource scenarios, we find that, irre-101 spective of the language pair, the optimal configu-102 ration should have 4K to 32K NMO on the source 103 side and 500 to 2K on the target. 104

> Section 2 summarises previous efforts to use symmetric BPE merge operations to improve MT performance. Section 3 explains our motivation for finding optimal BPE configurations by exploring asymmetric BPE. Section 4 outlines our experimental setup and presents the performance of the English-Hindi MT system on FLORES. Section 5 evaluates the setup for other language pairs in low resource context, concluding our observations in Section 6.

105

106

107

109

110

111

112

113

114

115

2 Related Work - Symmetrical BPE

116 In quite a lot of work on BPE to find the optimal segmentation of sub-words for low to high 117 resource settings, the recipe/BPE configuration of 118 keeping NMO the same for both source and target 119 is prevalent. Ding et al. (2019) showed that in low-120 resource settings with transformer architectures, 121 using 0-4000 (4K) merge operations is optimal, 122 as traditional 32K operations can reduce BLEU 123 scores by up to 4 points. Abid (2020) found lower 124 merge operations are better for English, Egyptian, 125 and Levantine in low-resource scenarios. Domingo 126 et al. (2023) revealed that the performance of the 127 system varies with tokenizers in the target lan-128 129 guage. There is a class of works which modify the segmentation strategy based on the typologi-130 cal and morphological properties of one or both 131 of the languages in the pair. Ortega et al. (2020) improved BPE by restricting segmentation based 133



Figure 1: CHRF++ Scores for Symmetrical BPE (32K,4K) vs Asymmetrical BPE $(m1 \neq m2)$

134

135

136

137

138

139

140

141

142

143

144

145

146

147

148

149

150

151

152

153

154

155

156

157

on token suffixes in polysynthetic languages. Lankford et al. (2021) highlighted the importance of NMO in low-resource settings for MT performance affecting English-to-Irish systems. Lee et al. (2024) addressed BPE's over-segmentation issue in morphologically rich languages like Korean by incorporating longer words from monolingual corpora, linking difficulties to morphological complexity. However, there also has been works which have shown potential directions to find different subword tokenising strategies. Work like Huck et al. (2017) and Poncelas et al. (2020) used cascading segmentation strategies. Xu et al. (2021) applied Marginal Utility to select appropriate vocabularies (VOLT), which refined segmentation methods and further improved results. Ngo Ho and Yvon (2021) used various BPE NMO in source and target languages independently to improve word alignment between language pairs, but this was not extended to check the performance in the MT system. Unfortunately, we did not find any study that explored asymmetric BPE configurations in diverse resource availability scenario as presented here.

3 Exploring Asymmetrical BPE

In practice, for a BPE configuration $m_1_m_2$, the 158 values of m_1 and m_2 are the same and the NMO 159 for both the source and target range from 8K to 40K 160 (Wu, 2016; Denkowski and Neubig, 2017; Cherry 161 et al., 2018; Renduchintala et al., 2019), but Ding 162 et al. (2019); Dewangan et al. (2021) found that 163 these are not ideal for low-resource language pairs. 164 Ding et al. (2019) observed $m_1 = m_2 <=4K$ 165 NMO is better for low-resource settings than 32K. 166 This finding aligns with our experiments on 0.1 mil-167 lion sentence pairs (English \leftrightarrow {Hindi, Telugu}), 168 showing better performance with 4K NMO (Fig-169

Sentence	bosusco, 54, runs an adventure tourism bureau.						
500 NMO	bo@@ su@@ sc@@ o, 5@@ 4, r@@ un@@ s an						
	ad@@ v@@ en@@ ture t@@ our@@ is@@ m bu@@ re@@ a@@ u .						
32K NMO	bo@@ su@@ sco, 54, runs an adventure tourism bureau.						

Table 1: Different NMO Effect

ure 1). Dewangan et al. (2021) found that identical 170 BPE configurations do not guarantee similar perfor-171 mance across different language pairs, as demon-172 strated by the English-Hindi vs. English-Telugu 173 comparisons (Figure 1) for 4K NMO. 174

175

176

177

178

179

181

183

186

189

190

191

192

194

To choose NMO, one should consider datasize and language pair as nuanced BPE strategies have benefitted morphologically complex languages (Ortega et al., 2020; Mujadia and Sharma, 2021). We examine symmetrical BPE configurations that employ the same NMO for both source and target languages, and investigate alternative BPE configurations by independently altering the NMO m_1 and m_2 for the source and target languages in English-Hindi, examining datasets ranging from low resource (50K) to high resource (8 million). We find 185 that this approach delivered better results (Figure 1) in low-resource contexts. Our extensive experiments with English-Hindi, evaluated on FLORES (Goyal et al., 2022), show that this kind of atypical BPE configuration yields better performance. We further validate the exploration of BPE configurations for optimal results by extending these experiments to English \leftrightarrow {Telugu, Shona, Norwe-193 gian, Kyrgyz, Hausa, and Inuktitut }.

Our work yields compelling results in favour 195 of optimising NMO based on training data size 196 and language pair. Figure 2 offers a conceptual overview of the optimal ranges of the BPE config-198 uration derived from our experiments in English-199 Hindi, as we move from low to high-resource settings. In this context, "ranges" denote the spectrum 201 of NMO values that serve as hyperparameters for the source and target subword tokenization mod-203 els in word segmentation. The performance gap between the best and symmetrical BPE systems is depicted by varying shades of green, with the 207 most significant improvement seen in low-resource settings (darker shade of green). As the dataset size 208 grows (for medium to high resource contexts), the performance variance between different configurations decreases (lighter shade of green). 211

Evaluation on English \iff Hindi 4

212

213

214

215

216

217

218

219

220

221

222

223

224

225

226

227

228

229

230

231

232

233

234

235

236

237

238

239

240

241

242

243

244

245

246

247

248

249

250

251

252

We explore BPE configurations with the Samanantar dataset (Ramesh et al., 2022) for English-Hindi¹, with 8 million parallel sentences. English text is tokenised, normalised and lowercase using scripts from Moses, while for the preprocessing of the Hindi text we use the Indic NLP library (Kunchukuttan, 2020). We simulate various training set size scenarios by grouping sentences based on English sentence length (Table 2) and randomly sample datasets of sizes 0.05 million (M), 0.1M, 0.5M, 1M, 4M and 8M, maintaining sentence length propotion. For each language and dataset size, the BPE tokenizer is trained on eight different NMOs - 0.5K, 1K, 2K, 4K, 8K, 16K, 25K and 32K. After subword segmentation, all possible BPE configurations (src_{500} _tgt₅₀₀, src_{500} _tgt₁₀₀₀, src₅₀₀_tgt₂₀₀₀, etc.) are trained with Transformer architecture (Vaswani et al., 2017) using Fairseq $(Ott et al., 2019)^2$. Training a single BPE configuration $m_1_m_2$ for all dataset sizes required an average of 1040 GPU hours using a 1080 Ti. This results in 64 system configurations trained per language direction for 6 datasets, totalling 768 systems ($64 \times 6 \times 2$), and resultant systems are evaluated on FLORES (Goyal et al., 2022) Dataset using CHRF++ (Popović, 2015) to assess the impacts of the BPE configurations.

Figures 3 present the performance of all configurations for English \iff Hindi systems in a low resource scenario (for data set sizes of 0.05M, 0.1M and 0.5M). And Figures 4 show the performance of all configurations on 1M, 4M and 8M dataset sizes. Each subgraph represents performance on a particular dataset size, with the x-axis being the source NMO. The black stepped dotted lines indicate the maximum CHRF++ score for each dataset size considering for each source NMOs.

In figure 3 for low-resource environments (0.05M, 0.1M and 0.5M) systems, as noted by (Ding et al., 2019), the use of symmetric BPE

¹English Hindi bitext under the created directory of version 2 of Samanantar is used.

²Transformer Base Architecture

		Source NMO							
		0.5K	1K	2K	4K	8K	16K	25K	32K
Target NMO	0.5K								
	1K					Optimal Fo			
	2K								
	4K			Optimal Fo	or Medium I	Resource			
	8K								
	16K								
	25K								
	32K				Optimal For High Resource			urce	

Figure 2: Changes in Optimal BPE Configuration from Low- to High-Resource Settings

Length bin	1 to 10	11 to 15	16 to 20	21 to 25	26 to 30	31 to 35	35 to 40	>=41	Total
No. of sentences	2792334	1655162	1150396	854091	617318	420583	275774	414926	8180584
Percentage	34.13	20.23	14.06	10.44	7.55	5.14	3.37	5.07	100

Table 2: Distribution of sentences in groups based on token length for full data

configuration with lower NMOs improves performance over high NMOs. However, the best results 254 are achieved using asymmetric BPE configurations 255 when the source has a higher NMO than the target. 256 We see a maximum performance gain when the source NMO is very high and the target NMO very low (we see consistent performance with the target 259 260 NMO = 500). Conversely, when the target's NMO is greater than that of the source, performance de-261 clines, like for the Hindi to English 0.1M dataset, 262 performance of 500_25K and 500_32K was worse than symmetric BPE configurations.

In medium resource settings (1M), optimal con-265 figuration for both the source and target NMO moves towards medium (2K - 8K) range of values, with relatively small score variations across 269 different configurations \approx 3 CHRF++ between extremes. In contrast, among configurations with 270 low resources, the best system outperforms the 271 weakest system by a margin of \approx 15 CHRF++ 272 scores and the best symmetric BPE configuration by \approx 5 CHRF++ scores. For the high-resource scenario, the score disparity between the most effective and least effective configuration is minimal (< 2 CHRF++), with the most effective system em-277 ploying a 32K NMO for the target. This pattern 278 of shifting optimal BPE configurations across the 279 size of the dataset is also evident if we consider only the last point (32K on x-axis) in all subplots in Figures 3 and 4. This highlights that modifying the NMO on the target side, especially in a low resource scenario, plays a more vital role in determining the optimal BPE configuration. For example, for English to Hindi MT systems considering all the configurations $32K_x$ for all dataset sizes, we see a gradual change in the configurations of the best performing systems starting with 500 NMO for 0.05M to 32K NMO for 8M.

287

289

290

291

292

293

295

296

297

298

300

301

302

303

304

305

306

307

308

309

310

311

312

313

314

315

We also see that as far as symmetric BPE configurations are concerned, for low resource MT systems it under performs compared to the asymmetric configuration. But as we move to larger dataset size, symmetrical configurations start to perform comparable to asymmetrical configurations. However, using BPE configurations in an asymmetric fashion has yielded a statistically significant improvement in a low-resource environment.

We also compare our systems using optimal BPE configurations with VOLT (Xu et al., 2021)³. Figure 5 shows the comparison of the CHRF++ scores between the systems using VOLT tokenisation, optimal BPE, and baseline configurations⁴. Systems using optimal BPE settings surpass those tokenised with VOLT in all dataset sizes, especially for low resource setting with statistical significant (p<0.05).

5 Exploring Asymmetrical BPE Configurations for other language pairs

To confirm the effectiveness of asymmetric BPE configurations to improve machine translation (MT), we extended our work with atypical BPE configurations using 0.1 million sentence pairs in English \leftrightarrow {Telugu, Shona, Norwegian, Kyrgyz, Hausa and Inuktitut}⁵. The English-Telugu

³using hyperparameters specified in their work

⁴We selected the top performing BPE configuration as baseline configuration where source NMO = target NMO

⁵Sampled in a manner similar to English-Hindi sampling



(e) 0.5 Million English to Hindi



Figure 3: Evaluation of English \leftrightarrow Hindi MT Systems for 0.05M, 0.1M and 0.5M dataset sizes on FLORES, x-axis is source NMO and y-axis is CHRF++ scores



(e) 8 Million English to Hindi

(f) 8 Million Hindi to English

Figure 4: Evaluation of English \leftrightarrow Hindi MT Systems for 1M, 4M and 8M dataset sizes on FLORES, x-axis is source NMO and y-axis is CHRF++ scores



Figure 5: CHRF++ score comparison of Asymmetric BPE with VOLT



Figure 6: CHRF++ scores improvement with asymmetrical over symmetrical BPE

pair is taken from Ramesh et al. (2022), English-316 {Shona, Norwegian, Kyrgyz, Hausa} from Gowda 317 et al. (2021) and English-Inuktitut from Joanis et al. (2020). These language pairs are intended to ex-319 amine the impact of both symmetric and asymmet-320 ric BPE configurations in a low resource scenario 321 across languages belonging to various language 322 families and exhibiting diverse levels of morphological complexity. Using symmetric configurations 324 $(4K_4K, 32K_32K)$ as baselines and asymmetric configurations (8K_500, 16K_500) based on insights gained from optimal BPE settings in English-Hindi. We perform evaluations on the FLORES 328 test set, with the exception of Inuktitut, which is 330 evaluated using the test set from Joanis et al. (2020). Figure 6 compares the performance of asymmetric and symmetric BPE. Asymmetric BPE improves 332 translation in four of six X to English translation systems and for all systems while translating from 334 English to X languages. The improvements are statistically significant (p<0.05, languages indicated 336 by * in the figure) compared to symmetric BPE, underscoring the importance of investigating BPE configurations beyond standard choices, particu-339 larly for low-resource language pairs.

6 Conclusion

In-depth examination of BPE configurations across diverse language pairs and differing dataset sizes reveals that typical configurations (n_n) do not always produce optimal results. As referenced in Section 2, in low-resource settings, systems benefit from using symmetric *n* NMO configurations when *n* is significantly smaller than 32K; our experiments with asymmetric BPE n_m show that further improvement in translation performance is possible, under low-resource conditions, when $n \gg$ *m* and *n*, *m* represent NMOs for source and target respectively.

In conclusion, this study underscores the importance of moving beyond the default segmentation options typically used in machine translation, especially when working with languages that have limited data available. In cases with medium- to high-availability of training data, symmetric Byte Pair Encoding (BPE) configurations might be adequate. However, their effectiveness tends to diminish when applied to low-resource scenarios. By utilising configurations with larger number of merger operations (NMO) for the source language and a 342

343

344

345

346

347

348

349

350

351

352

353

354

355

356

357

358

360

361

362

363

364

smaller NMO for the target language, significant
improvements in translation quality can be realised.
These asymmetric BPE configurations prove to be
consistently effective across various language families, regardless of their differing morphological
complexities, thereby emphasising the necessity
for specifically tailored configurations to optimise
translation results in a low resource scenario.

373 Limitation

This study is restricted by the computationally expensive resources required to analyse all BPE configurations for each language pair. However, the findings demonstrate that particular configuration ranges can consistently elevate translation quality, especially in low-resource setups, thus reducing the search grid by significant margin.

References

390

394

395

400

401

402 403

404

405

406

407

408

409

410

411

412

413

414

- Wael Abid. 2020. The SADID evaluation datasets for low-resource spoken language machine translation of Arabic dialects. In *Proceedings of the 28th International Conference on Computational Linguistics*, pages 6030–6043, Barcelona, Spain (Online). International Committee on Computational Linguistics.
- Colin Cherry, George Foster, Ankur Bapna, Orhan Firat, and Wolfgang Macherey. 2018. Revisiting characterbased neural machine translation with capacity and compression. In Proceedings of the 2018 Conference on Empirical Methods in Natural Language Processing, pages 4295–4305, Brussels, Belgium. Association for Computational Linguistics.
- Michael Denkowski and Graham Neubig. 2017. Stronger baselines for trustable results in neural machine translation. In *Proceedings of the First Workshop on Neural Machine Translation*, pages 18–27, Vancouver. Association for Computational Linguistics.
- Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. 2019. Bert: Pre-training of deep bidirectional transformers for language understanding. *Preprint*, arXiv:1810.04805.
- Shubham Dewangan, Shreya Alva, Nitish Joshi, and Pushpak Bhattacharyya. 2021. Experience of neural machine translation between indian languages. *Machine Translation*, 35(1):71–99.
- Shuoyang Ding, Adithya Renduchintala, and Kevin Duh. 2019. A call for prudent choice of subword merge operations in neural machine translation. In Proceedings of Machine Translation Summit XVII: Research Track, pages 204–213, Dublin, Ireland. European Association for Machine Translation.

Miguel Domingo, Mercedes García-Martínez, Alexandre Helle, Francisco Casacuberta, and Manuel Herranz. 2023. How much does tokenization affect neural machine translation? In *Computational Linguistics and Intelligent Text Processing*, pages 545–554, Cham. Springer Nature Switzerland. 415

416

417

418

419

420

421

422

423

424

425

426

427

428

429

430

431

432

433

434

435

436

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

- Philip Gage. 1994. A new algorithm for data compression. *C Users Journal*, 12(2):23–38.
- Thamme Gowda, Zhao Zhang, Chris Mattmann, and Jonathan May. 2021. Many-to-English machine translation tools, data, and pretrained models. In *Proceedings of the 59th Annual Meeting of the Association for Computational Linguistics and the 11th International Joint Conference on Natural Language Processing: System Demonstrations*, pages 306–316, Online. Association for Computational Linguistics.
- Naman Goyal, Cynthia Gao, Vishrav Chaudhary, Peng-Jen Chen, Guillaume Wenzek, Da Ju, Sanjana Krishnan, Marc'Aurelio Ranzato, Francisco Guzmán, and Angela Fan. 2022. The Flores-101 evaluation benchmark for low-resource and multilingual machine translation. *Transactions of the Association for Computational Linguistics*, 10:522–538.
- Matthias Huck, Simon Riess, and Alexander Fraser. 2017. Target-side word segmentation strategies for neural machine translation. In *Proceedings of the Second Conference on Machine Translation*, pages 56–67, Copenhagen, Denmark. Association for Computational Linguistics.
- Eric Joanis, Rebecca Knowles, Roland Kuhn, Samuel Larkin, Patrick Littell, Chi-kiu Lo, Darlene Stewart, and Jeffrey Micher. 2020. The Nunavut Hansard Inuktitut–English parallel corpus 3.0 with preliminary machine translation results. In *Proceedings* of the Twelfth Language Resources and Evaluation Conference, pages 2562–2572, Marseille, France. European Language Resources Association.
- Taku Kudo and John Richardson. 2018. SentencePiece: A simple and language independent subword tokenizer and detokenizer for neural text processing. In Proceedings of the 2018 Conference on Empirical Methods in Natural Language Processing: System Demonstrations, pages 66–71, Brussels, Belgium. Association for Computational Linguistics.
- Anoop Kunchukuttan. 2020. The IndicNLP Library. https://github.com/anoopkunchukuttan/ indic_nlp_library/blob/master/docs/ indicnlp.pdf.
- Seamus Lankford, Haithem Alfi, and Andy Way. 2021. Transformers for low-resource languages: Is féidir linn! In *Proceedings of Machine Translation Summit XVIII: Research Track*, pages 48–60, Virtual. Association for Machine Translation in the Americas.
- Jungseob Lee, Hyeonseok Moon, Seungjun Lee, Chanjun Park, Sugyeong Eo, Hyunwoong Ko, Jaehyung Seo, Seungyoon Lee, and Heuiseok Lim. 2024.

472 473

471

474 475

476

477 478

- 479
- 480 481
- 482 483
- 484 485

486

- 487 488
- 489
- 490
- 491 492
- 493 494
- 495 496
- 497 498

499

- 501
- 502
- 503 504

505

506 507 508

509 510

511

- 512 513
- 514 515
- 516 517
- 518 519

520

523

524 525

526

Research Track, pages 244-255, Dublin, Ireland. European Association for Machine Translation. 527

- Length-aware byte pair encoding for mitigating oversegmentation in Korean machine translation. In Findings of the Association for Computational Linguistics: ACL 2024, pages 2287-2303, Bangkok, Thailand. Association for Computational Linguistics.
- Vandan Mujadia and Dipti Misra Sharma. 2021. English-Marathi neural machine translation for LoResMT 2021. In Proceedings of the 4th Workshop on Technologies for MT of Low Resource Languages (LoResMT2021), pages 151-157, Virtual. Association for Machine Translation in the Americas.
- Anh Khoa Ngo Ho and François Yvon. 2021. Optimizing word alignments with better subword tokenization. In Proceedings of Machine Translation Summit XVIII: Research Track, pages 256–269, Virtual. Association for Machine Translation in the Americas.
- John E Ortega, Richard Castro Mamani, and Kyunghyun Cho. 2020. Neural machine translation with a polysynthetic low resource language. Machine Translation, 34(4):325-346.
- Myle Ott, Sergey Edunov, Alexei Baevski, Angela Fan, Sam Gross, Nathan Ng, David Grangier, and Michael Auli. 2019. fairseq: A fast, extensible toolkit for sequence modeling. In Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics (Demonstrations), pages 48-53, Minneapolis, Minnesota. Association for Computational Linguistics.
- Alberto Poncelas, Jan Buts, James Hadley, and Andy Way. 2020. Using multiple subwords to improve English-Esperanto automated literary translation quality. In Proceedings of the 3rd Workshop on Technologies for MT of Low Resource Languages, pages 108-117, Suzhou, China. Association for Computational Linguistics.
- Maja Popović. 2015. chrF: character n-gram F-score for automatic MT evaluation. In Proceedings of the Tenth Workshop on Statistical Machine Translation, pages 392-395, Lisbon, Portugal. Association for Computational Linguistics.
- Gowtham Ramesh, Sumanth Doddapaneni, Aravinth Bheemaraj, Mayank Jobanputra, Raghavan AK, Ajitesh Sharma, Sujit Sahoo, Harshita Diddee, Mahalakshmi J, Divyanshu Kakwani, Navneet Kumar, Aswin Pradeep, Srihari Nagaraj, Kumar Deepak, Vivek Raghavan, Anoop Kunchukuttan, Pratyush Kumar, and Mitesh Shantadevi Khapra. 2022. Samanantar: The Largest Publicly Available Parallel Corpora Collection for 11 Indic Languages. Transactions of the Association for Computational Linguistics, 10:145-162.

Adithya Renduchintala, Pamela Shapiro, Kevin Duh,

and Philipp Koehn. 2019. Character-aware decoder for translation into morphologically rich languages.

In Proceedings of Machine Translation Summit XVII:

Rico Sennrich, Barry Haddow, and Alexandra Birch. 2016. Neural machine translation of rare words with subword units. In Proceedings of the 54th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers), pages 1715–1725, Berlin, Germany. Association for Computational Linguistics.

528

529

530

531

532

533

535

536

537

538

539

540

541

542

543

544

545

546

547

548

549

550

551

552

553

554

555

556

557

- Peter Smit, Sami Virpioja, Stig-Arne Grönroos, and Mikko Kurimo. 2014. Morfessor 2.0: Toolkit for statistical morphological segmentation. In Proceedings of the Demonstrations at the 14th Conference of the European Chapter of the Association for Computational Linguistics, pages 21-24, Gothenburg, Sweden. Association for Computational Linguistics.
- Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz Kaiser, and Illia Polosukhin. 2017. Attention is all you need. Advances in neural information processing systems, 30.
- Yonghui Wu. 2016. Google's neural machine translation system: Bridging the gap between human and machine translation. arXiv preprint arXiv:1609.08144.
- Jingjing Xu, Hao Zhou, Chun Gan, Zaixiang Zheng, and Lei Li. 2021. Vocabulary learning via optimal transport for neural machine translation. In 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), pages 7361-7373, Online. Association for Computational Linguistics.