Dict-NMT: Bilingual Dictionary based NMT for Extremely Low Resource Languages

Anonymous ACL submission

Abstract

Neural Machine Translation (NMT) 001 models have been effective on large bilingual datasets. However, the existing methods and techniques show that the model's performance is highly dependent on the number of examples in training data. For many languages, having such amount of corpora is a far fetched dream. Taking inspiration from monolingual speakers exploring new languages using bilingual dictio-011 naries, we investigate the applicability 012 of bilingual dictionaries for languages with extremely low, or no bilingual cor-014 pus. In this paper, we explore methods using bilingual dictionaries with an NMT model to improve translations for extremely low resource languages. We extend this work for multilingual systems, exhibiting zero-shot property. 021 We present a detailed analysis of effects of quality of dictionary, training dataset size, language family, etc., on the translation quality. Results on multiple low-resource test languages show a clear advantage of our bilin-026 gual dictionary-based method over the baselines.

1 Introduction

034

With the growing interest in improving automatic translation systems, deep learning-based models have played a significant role. They have a ubiquitous influence on such solutions. Neural Machine translation has been ruling the roost in recent times both in academia as well as in industries. It has outperformed other translation methods, and even human translators for some languages (Bojar et al., 2016) (Bentivogli et al., 2016) (Barrault et al., 2020). The encoder-decoder framework of NMT models allow them to transfer the semantic and syntactic information more precisely.

043

044

045

046

047

051

055

058

060

061

062

063

064

065

067

068

069

070

071

072

073

074

075

076

078

079

081

One of the major challenges for such languages is training corpora of sufficient size. Such models need for large bilingual or monolingual datasets, where usually it ranges between 1-50 million parallel sentences. For the extremely low resourced languages, dataset with size lesser than 20 thousand parallel sentences, NMT models have not been that successful (Östling and Tiedemann, 2017). The standard approach to this problem has mostly relied on techniques such as transfer learning (Zoph et al., 2016), and data augmentation approaches such as back-translation (Sennrich et al., 2015) (Przystupa and Abdul-Mageed, 2019) and data diversification (Nguyen et al., 2019).

The use of prior knowledge sources for translation of low-resource languages, such as bilingual dictionaries, is still under-explored. The work of (Duan et al., 2020) and (Nag et al., 2020) are closest to ours and both utilize bilingual dictionary, but they use additional large monolingual corpus while we use an extremely small test language's bilingual corpus or no bilingual corpora of the test language at all. The existing approaches mostly depend on the availability of some additional corpora like target monolingual corpus and target-to-source model for backtranslation, sister language for transfer learning or additional computations as in data diversification. One of the most common and widely available prior knowledge resources across low resourced languages is the bilingual dictionary which has shown potential in NMT in recent

<u>n</u>94

096

098

100

102

103

104

105

107

108

109

110

111

112

113

times.

In our work, we explore the use of bilingual dictionary for translation of extremely low languages. Any meaningful translation requires us to address the points as illustrated in Figure 1. In extremely low resource languages, NMT falls behind given the lack of enormous quantity of data required to train them properly. We study the potential of assisting the NMT models with the contextual dictionary transformation. Our proposed method involves the



Figure 1: Translation (a) word mapping task, which can be partially, or completely achieved with bilingual dictionary lookup. (b) is about the association a word has with its surrounding words, which in turn affects its alignment. (c) is the transformation of syntactic features of the source language to the source language.

use of bilingual dictionary for addressing the points (Figure 1 (a) and (b)) and an NMT model for (Figure 1 (c)), i.e., we use an NMT model to transform a distorted sentence into a meaningful sentence within the same language. Using this method, we propose two simple frameworks which can be extremely useful for languages with extremely less or no corpus available but having a bilingual dictionary. Summarizing the contributions of our paper as follows:

- We introduce a simple and effective method for incorporating a bilingual dictionary in a neural machine translation task.
- We propose a one-to-one bilingual dictionary based NMT model for extremely low-resource languages.
- We propose a many-to-one NMT model

capable of translating for languages it has never seen in the training sets.

115

116

117

118

119

120

121

122

123

124

125

126

127

128

129

130

131

132

133

134

135

136

We provide a brief description of our method in Section 2. We discuss the usage of the bilingual dictionary, tokenizer, and the NMT model. We explore the applicability of our proposed method in two settings, extremely less corpus and no corpus available for concerned language. In Section 3, we describe our one-to-one translation framework useful for translation in extremely low resource setting. We provide a detailed analysis with comparison among translation quality, dataset size and dictionary quality. In Section 4, we provide detailed information about our proposed many-to-one translation framework, which shows zero-shot property. We summarize and conclude our results and contributions in Section 5.

2 Dict-NMT: Assisting NMT model with bilingual dictionary

We propose a simple yet effective method 137 of translation, dict-NMT, using an NMT 138 model with the help of the respective 139 languages' bilingual dictionary. We use 140 a bilingual dictionary as word-to-word 141 translator to convert words from the 142 source language to an intermediate se-143 quence. This distorted sentence in the 144 target language is then fed to an NMT 145 model, here Transformer, to learn the re-146 lation between the intermediate sequence 147 and ground truth (Figure 2). This opens 148 up doors for various frameworks for trans-149 lation. One simple way is to apply this 150 method to a one-to-one translation system 151 (Section 3). Furthermore, one can also 152 devise a many-to-one translation frame-153 work (Section 4), where the NMT model 154 is trained on word-to-word translations 155 from various languages. This generalised 156 model can then be used even for languages 157 which were not used in the training data. 158 Other possible ways include fine-tuning 159 the generalised model on a specific lan-160 guage. Other data augmentation methods, 161 such as backtranslation and data diversifi-162 cation, are also applicable to our proposed 163 method. Another possible way of augment-164

ing data is by adding intra-shuffled (i.e., 165 words within a sentence *s* are shuffled), 166 noisy (replace tokens in s with random 167 tokens with some probability) sentences 168 from target language to the training data. We leave these methods for future work. 170

Bilingual Dictionary 2.1

171

172

174

175

176

177

178

181

182

183

186

187

189

190

192

193

194

195

201

211

A dictionary is a map of words from the source language to the target language, where the mapping can be one to many. Here, we consider mappings that are word to word and not word to phrase.

First, using the dictionary, we change the source language sentences into an intermediate sequence. This step would reduce the workload on our NMT model from learning the word meanings from the available small dataset. If a word in the source sentence is present in the dictionary then it is converted accordingly in intermediate sequence, otherwise, it remains unchanged in the intermediate sequence, i.e., we consider the word to be in the target language space. When using the dictionary, there might exist multiple target language words as meaning for a source language word. We settle this problem of polysemy by selecting the word most similar to the previous word's dictionary translation (using target language's pre-trained word embeddings). This would help us to preserve the contextual information.

More precisely, for any source language 197 S and target language T with a bilingual 198 dictionary $D_{S \to T}$, the first step is to trans-199 late the text in S word-to-word to T using $D_{S \to T}$. In case the mapping is not available for any word w in S, it is mapped to itself and is considered a random noise in T. For the case of polysemous words, we take the help of word embeddings of T. We select the word (in T) most similar 206 to the previous word's dictionary transla-207 tion (using target language's pre-trained word embeddings). For instance, given is a sentence $s = \{s_1, s_2, ..., s_n\}$ in S with 210 word-to-word translation $t = \{t_1, t_2, ..., t_n\}$ in T. For any $s_i(i > 1)$ having dictionary 212 translations $D_{S \to T}(s_i) = \{t_i^1, t_i^2, ..., t_i^m\}$, we 213



Figure 2: Our proposed method involving bilingual dictionary for NMT.

select its translation as

$$t_i = \underset{t \in D_{S \to T}(s_i)}{\operatorname{argmax similarity}(t_j, t)}$$
 219

214

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

where $j = \max\{j' < i | s_j \in Dom(D_{S \to T})\}$. We randomly select translation for the first occurring polysemous word. Here, for our experiments, we assume that the target language is a popular one, thus, decent word embeddings for T exist. This method would help us to preserve the contextual information. However, if the first randomly selected word is erroneous, the trailing polysemous words might have incorrect translations.

2.2 Tokenizer

Tokenization is the process of breaking the given text into smaller chunks. Since the model input and output are in the same language, we share the tokenizer for both of them. The intermediate representation might consist words from foreign language. Thus, instead of using the traditional whitespace tokenizer and giving all such words a <OOV> token, we use subspace tokenizer to handle the large amount of out-of-vocabulary words. This way, the noise created by the tokens of foreign language, would help the model being more robust.

2.3 NMT Model

Since both intermediate sequence and target belong to the same language, the NMT

Algorithm 1: Dict-NMT. $D_i = (S_i, T)$ is a set of parallel sentences from language S_i to T. Corresponding to the language pair (S_i, T) , we have a bilingual dictionary B_i where $B_i(s)$ is a word-to-word translation of s. We train the model M on $\{D_i\}_{i=1}^n$ using $\{B_i\}_{i=1}^n$.

1 **Procedure** Train $({D_i}_{i=1}^n, {B_i}_{i=1}^n, p)$: M randomly initialised NMT model. 2

Train M on create_dataset($\{D_i\}_{i=1}^n$, $\{B_i\}_n$, p) until it converges 3

return M 4

10

11

12

13

14 15

245

246

247

248

252

253

255

261

265

267

269

270

271

5 Function create_dataset $({D_i}_{i=1}^n, {B_i}_n, p)$:

```
D' = \phi
6
       for D \in \{D_i\}_{i=1}^n do
7
```

```
for (s, t) \in D do
8
```

```
9
```

```
if count_T(B_i(s)) \ge p then
```

```
D' = D' \cup (B_i(s), \mathbf{t})
```

```
/* count<sub>T</sub>(B<sub>s</sub>) = % of words in s having dictionary translations
shuffle D'
return D'
```

model is relieved from learning the word meanings. The model will now try to focus mostly on learning the grammar for the target language space. The NMT model learns the mappings from the source invariant representations coming from various languages to the target language and tries to generalise which would be beneficial for unknown languages.

Our proposed method can be applied to any NMT model. For our experiments, we use the state-of-the-art Transformer (Vaswani et al., 2017) model. Since, the intermediate sentence, i.e., the input for the Transformer, in itself does not make any sense, the attention mechanism helps to understand the dependencies of words through the whole sequence. The encoderdecoder framework allows us to find the meaning of the words not translated by the dictionary while preserving the context.

3 Dict-NMT for one-to-one translation

We propose a dictionary based one-to-one translation framework for extremely low resource settings. Given a language pair (S,T), we train an NMT model on the wordto-word dictionary translations of S, and T (i = 1 in Algo 1).

*/

272

273

274

275

277

278

279

281

284

285

286

290

292

293

294

295

3.1 Experimental Settings

We extensively check the effectiveness of the dictionary (by varying the dictionary percentage) across five European languages' translation tasks as well as the size of the bilingual corpora. We keep Transformer as our baseline model. We use 4 layer Transformer with 100 embedding/hidden units, and 400 feed-forward filter size. We tie source and target embeddings. We keep batch size 32, epochs 50, dropout 0.1 and optimizer Adam. In this work, we use the pre-trained BERT WordPiece tokenizer (Devlin et al., 2019), a subword tokenizer.

3.1.1 Bilingual dictionary

For our experiments, we use the publicly available Facebook MUSE's¹ bilingual dictionary, which consists of 110 large-scale ground-truth bilingual dictionaries (Conneau et al., 2017). For preserving the context while dictionary translation, we

¹https://github.com/facebookresearch/MUSE

297 298

308

310

311

314

315

318

320

323

326

327

332

338

340

342

344

use the Fasttext embeddings (Bojanowski et al., 2017).

3.1.2 Data

For our experiments, we consider Europarl v7 parallel corpus (Koehn, 2002) for Pt-En, Sv-En, Nl-En, Pt-En, and Fr-En language pairs. Here we selected English as our target language in all the cases. The intuition behind this is that any bilingual dictionary for an extremely low resource language would be created by taking a commonly used language, so that it can be of practical use. As English is one such language, we tried our experiments with it.

We filter each sentence such that it contains at most 80 tokens. We use these data in low resource setting, i.e. we use only 2K, 8K, 16K and 20K data size for each language. We create these datasets according to the percentage of words from each sentences available in the corresponding dictionary, precisely we did this for 50%-80% (Table 1). For each data size 0.05% is the test set and rest we use as training set.

3.2 Results and Analysis

We perform intensive experiments on the effectiveness of training data size and dictionary coverage on the performance of the translation system. Table 2 shows comparison between the baseline (bilingual dictionary based word-to-word translation) and our proposed method. The best scores for each language pair, along with the dictionary coverage are reported in the table. The best result is chosen over the dictionary coverage (50% - 80%), i.e. least percentage of words in each sentence available in the bilingual dictionary, and varied dataset size (2K - 20K). We report arithmetic mean of scores on 3 different datasets sampled from the same large data. The scores show a significant increase from simple word-to-word translations (4.8 - 8.6). We performed experiments for three language pairs with simple transformer as well. However, due to very less data, the model seemed to struggle considerably. For Pt-En, Sv-En, and Fr-En

Dict	Ro-Er	n Pt-	En l	Fr-En	It-En	Es-En			
%	(399.37)	K) (1.9	6 M)	(2 M)	(1.9 M)	(1.97 M)			
50	104	K 438	3.6K	$1.4\mathbf{M}$	941 K	$1.5\mathbf{M}$			
60	22.6	6K 77	7.7K 5	36.4K	202.6K	631.6K			
70	4	K 11	1.3 K 1	06.2K	26.8K	111.4K			
80	9'	75 2	2.5K	17.7K	4.8K	$15.9 \mathrm{K}$			
90	2	44	937	2.6K	$1.9 \mathrm{K}$	2.5K			
100	2	30	920	2K	$1.8 \mathrm{K}$	2K			
Italic									
Dict	Dict Da-I		Sv-En		e-En	Nl-En			
%	(1.97	7 M) (1.86 M	() (1	.92 M)	(2 M)			
50) 1.	3 M	1.7N	1	1.9 M	1.7 M			
60	506	.6K	$1.3\mathbf{M}$		$1.7\mathbf{M}$	$1\mathbf{M}$			
70 107.9		.9K	563.7K		$1\mathbf{M}$	$31.7 \mathrm{K}$			
80) 20	.2K	124.9 I	X 3	$17.5 \mathrm{K}$	55K			
90) 2	.6K	9.2K		$28.2 \mathrm{K}$	$4.7 \mathrm{K}$			
100) 1	.9K	$1.8\mathbf{I}$	X	3.3K	$1.9\mathbf{K}$			
Germanic									
Dict	Bg-En	Cz-En	Pl-	En	Sl-En	Sk-En			
%	(406.9K)	(646.6K) (632.	57 K) (623.49 M)	(640.72K)			
50	33.2K	136.8F	x 13	7.8K	64.6K	180.4K			
60	6.4K		v 3	b.5K	16 K	49.2K			

Slavic

8.6K

2.7K

1.1K

1K

3.9K

1.6K

1.1K

1.1K

11.3K

3.5K

1.4K

1.3K

70

80

90

100

 $1\mathbf{K}$

233K

58K

57K

9.5K

3.1K

1.4K

1.3K

Table 1: Dataset size after filtering sentences containing at least Dict% of dictionary words.

Language	Dict %	BLEU	
Language		Base	w D
Pt-En	55%	6.9	15.4
Sv-En	70%	8.4	17.0
Nl-En	65%	5.9	10.7
Fr-En	65%	8.6	15.4
Da-En	65%	7.7	16.1

Table 2: Results: Best BLEU score for each language. w D is "with dictionary", i.e. our proposed method, Base is "dictionary baseline" which is simply word-to-word translation . We get the best scores for maximum data size (i.e., 20K). Training set dictionary coverage (Dict %) is given for each corresponding score. The BLEU scores are calculated using SacreBLEU's corpus_bleu (Post, 2018)

language pairs, best scores on 20K dataset came to be 1.56, 1.39, and 1.37 respectively. This shows there is a clear advantage of using the proposed method for extremely low resource languages.

In Figure 3, we have heat-maps of BLEU

348 349

347



Figure 3: BLEU Scores for One-to-One translation method. We report the scores on test data with 80% dictionary coverage, as it was maximum in every case.

scores for different languages calculated
over different datasets. The x-axis shows
the size of the dataset and the y-axis
shows the dictionary percentage. We can
have the following observations from the
maps,

- **BLEU VS Dataset size:** The model clearly benefits from increasing the dataset size in an extremely low resource setting. There is a direct correlation between the score and the number of training examples.
- BLEU VS Training data dictionary 365 **coverage:** Dictionary coverage can be seen as inversely proportional to the amount of noise generated by the untranslated words from the dictionary. The best scores for each column of any map are always somewhere in the 371 middle (except 20K dataset of French). 372 We suspect this behavior is linked to finding the correct balance of noise and generalisation. With more noise 375 (less dictionary coverage), the method 376 seems to get more robust, however, 377 it underfits when it is exposed to too many of them.

4 Dict-NMT for many-to-one translation

A conventional idea for a many-to-one model would involve mapping the source text to a common representation space which would then further be used by the model to generate the translations. Fixing the target language, we can create a common representation for any given source language by translating the source text word-to-word into the target language using the bilingual dictionary. This is similar to how we humans translate any foreign language with the help of a bilingual dictionary. 381

382

383

384

385

386

387

388

389

390

391

392

393

394

395

396

397

398

399

400

401

402

403

404

405

406

407

408

We propose a many-to-one translation framework, which, just using a bilingual dictionary, can translate for languages which are not present in the training phase- absolute zero-shot translation. Given a test language pair (S,T), we train an NMT model on dictionary based wordto-word translations of language pairs $\{(S_i,T)\}_{i=1}^n$, where $S \neq S_i$ for i = 1,..,n. Our goal is to make the model invariant of source language. We achieve this via adding word-to-word dictionary translations from various languages coming from different families (Algo 1).

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

4.1 Experimental Settings

We perform a comprehensive study on the 410 effect of dataset size, no. of languages, 411 inclusion of test language family, and dic-412 tionary coverage in test set, on the transla-413 tion quality. We perform our experiments 414 on European languages with English as 415 416 the target language. We keep the tokenizer, NMT model and its hyperparameters simi-417 lar to that of previous experiment's setting. 418 We use Facebook MUSE's bilingual dictio-419 nary for this experiment as well. 420

4.1.1 Dataset

409

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

We perform our experiments on Europarl v7 parallel corpus, fixing English as our target language. We used languages from three families, namely, Italic (Romanian, Spanish, Portuguese, French, Italian), Slavic (Bulgarian, Czech, Polish, Slovene, Slovak) and Germanic (Danish, Swedish, German, Dutch). We use Romanian, Bulgarian, and Danish as our test languages. We analyse our results on training data, intra and inter combination of the language families with sizes 50K, 150K, 500K, and 1M. We test our experiments on 600 sentences. create_dataset()(Algo 1) shows how we created training data for our experiment. In our experiments, for a training set create_dataset($\{D_i\}_{i=1}^n$, $\{B_i\}_n$, p) (Algo 1), we take equal number of sentences from all n languages. The case of polysemy is handled the same way as it was in the previous experiment. We use the notation "All" for the combination of the above mentioned languages from all three language families (Italic, Germanic, and Slavic).

4.2 Results and Analysis

We present the best scores for three lan-448 guage pairs, Ro-En, Bg-En and Da-En in 449 Table 3. We further compare the scores 450 with word-to-word dictionary translations. 451 We choose the best score over varied train-452 ing data size (50K - 1M), Test set dictionary 453 coverage (0% - 80%), and combination of 454 language families. There is a significant 455 difference in scores of baseline and our 456 proposed method. Considering the fact 457

that the training sample has no examples from test languages, the resultant score demonstrates the zero-shot property of the proposed method.

Language	Dict %	BLEU		
Language		Base	w D	
Ro-En	50%	9.4	28.1	
Bg-En	50%	8.2	15.4	
Da-En	50%	7.7	13.4	

Table 3: Results: Best BLEU score for each language. w D is "with dictionary", i.e. our proposed method, Base is "dictionary base-line" which is simply word-to-word translation . We get the best scores for "All" dataset (Germanic Italic Slavic) with 500K parallel sentences. Training set dictionary coverage (Dict %) is given for each corresponding score. The test set of Ro and Bg has atleast 80% dictionary coverage, while Da has 70% (the raw dataset was too little to make a decent test dataset with similar number of samples).

We perform experiments to test effect of dataset size, inclusion of test language family, and test data dictionary coverage.

- BLEU VS Test data dictionary coverage: From figure 5, it is evident that the scores increase with the increase in dictionary coverage of test data, i.e., the NMT model gets better assisted with more word-to-word translations in a given sentence.
- **BLEU VS Training set data size:** With increase in data size, the scores increase as well (Table 4). However, it tends to converge on the data size between 500*K* and 1*M*.
- Effect of addition of language families: From Figure 6 it can be observed that the score stays the least for model trained just on Germanic family. There is a slight increase in score for Italic family. However, it increases significantly when we start combining language families together. We get the highest score for "All", which is a combination of all three language families. There is a slight decrease in score when we add sentences from two different languages.



Figure 4: BLEU VS Training Dataset Size



Figure 5: BLEU VS Test Data Dictionary Coverage. Here, the reported scores are for model trained on "All" dataset with size 500K

We suspect less number of parameters
of the model to be the reason behind
such behaviour. For better generalisation on more number of languages,
we believe larger NMT models would
be beneficial.



Figure 6: BLEU VS Effect of Test Family (Romanian) (TGS = Turkish + Greeek)

5 Conclusion

Using Europarl corpus, we showed that our method of incorporating bilingual dictionaries for NMT tasks can be quite effective. Given a dictionary, it not only works for languages with extremely low corpus, but also for languages with no parallel or monolingual corpus at all. We analyze the extent of improvement that can be done by varying dictionary percentage and with the range of size of datasets. This work can be extended by blending our method with other state-of-the-art approaches such as back translation, and transfer learning. We believe this work will motivate researchers to explore other possibilities of incorporating bilingual dictionaries for NMT in extremely low resource settings.

496

497

498

499

500

501

502

503

504

505

506

507

508

509

510

511

512

513

514

490

References

515

530

531

532

534

535

537

538

539

541

542

543

544

545

546

547

548

549

552

553

554

555

557 558

559

560

561

562

563

564

565

568

569

- Loïc Barrault, Magdalena Biesialska, Ondřej 516 Bojar, Marta R. Costa-jussà, Christian Fe-517 dermann, Yvette Graham, Roman Grundkiewicz, Barry Haddow, Matthias Huck, 519 Eric Joanis, Tom Kocmi, Philipp Koehn, 520 Chi-kiu Lo, Nikola Ljubešić, Christof Monz, 521 Makoto Morishita, Masaaki Nagata, Toshi-522 aki Nakazawa, Santanu Pal, Matt Post, and 523 Marcos Zampieri. 2020. Findings of the 524 2020 conference on machine translation 525 (WMT20). In Proceedings of the Fifth Con-526 ference on Machine Translation, pages 1–55, 527 Online. Association for Computational Lin-528 guistics.
 - Luisa Bentivogli, Arianna Bisazza, Mauro Cettolo, and Marcello Federico. 2016. Neural versus phrase-based machine translation quality: a case study. *arXiv preprint arXiv:1608.04631*.
 - Piotr Bojanowski, Edouard Grave, Armand Joulin, and Tomas Mikolov. 2017. Enriching word vectors with subword information. *Transactions of the Association for Computational Linguistics*, 5:135–146.
 - Ondřej Bojar, Rajen Chatterjee, Christian Federmann, Yvette Graham, Barry Haddow, Matthias Huck, Antonio Jimeno Yepes, Philipp Koehn, Varvara Logacheva, Christof Monz, Matteo Negri, Aurélie Névéol, Mariana Neves, Martin Popel, Matt Post, Raphael Rubino, Carolina Scarton, Lucia Specia, Marco Turchi, Karin Verspoor, and Marcos Zampieri. 2016. Findings of the 2016 conference on machine translation. In Proceedings of the First Conference on Machine Translation: Volume $\tilde{2}$, Shared Task Papers, pages 131-198, Berlin, Germany. Association for Computational Linguistics.
 - Alexis Conneau, Guillaume Lample, Marc'Aurelio Ranzato, Ludovic Denoyer, and Hervé Jégou. 2017. Word translation without parallel data. *arXiv preprint arXiv:1710.04087*.
 - Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. 2019. BERT: Pretraining of deep bidirectional transformers for language understanding. In Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long and Short Papers), pages 4171–4186, Minneapolis, Minnesota. Association for Computational Linguistics.
 - Xiangyu Duan, Baijun Ji, Hao Jia, Min Tan, Min Zhang, Boxing Chen, Weihua Luo, and

Yue Zhang. 2020. Bilingual dictionary based neural machine translation without using parallel sentences. In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, pages 1570– 1579. 573

574

575

576

577

578

579

580

582

583

584

585

586

587

588

589

590

591

592

593

594

595

596

597

598

599

600

601

602

603

604

605

606

607

609

610

611

612

613

614

615

616

617

618

619

620

621

622

- Philipp Koehn. 2002. Europarl: A multilingual corpus for evaluation of machine translation. Draft.
- Sreyashi Nag, Mihir Kale, Varun Lakshminarasimhan, and Swapnil Singhavi. 2020. Incorporating bilingual dictionaries for low resource semi-supervised neural machine translation. *arXiv preprint arXiv:2004.02071*.
- Xuan-Phi Nguyen, Shafiq Joty, Wu Kui, Ai Ti Aw, Jiuxiang Gu, Jason Kuen, Shafiq Joty, Jianfei Cai, Vlad Morariu, Handong Zhao, et al. 2019. Data diversification: An elegant strategy for neural machine translation.
- Robert Östling and Jörg Tiedemann. 2017. Neural machine translation for low-resource languages. *arXiv preprint arXiv:1708.05729*.
- Matt Post. 2018. A call for clarity in reporting BLEU scores. In Proceedings of the Third Conference on Machine Translation: Research Papers, pages 186–191, Belgium, Brussels. Association for Computational Linguistics.
- Michael Przystupa and Muhammad Abdul-Mageed. 2019. Neural machine translation of low-resource and similar languages with backtranslation. In Proceedings of the Fourth Conference on Machine Translation (Volume 3: Shared Task Papers, Day 2), pages 224–235.
- Rico Sennrich, Barry Haddow, and Alexandra Birch. 2015. Improving neural machine translation models with monolingual data. *arXiv preprint arXiv:1511.06709*.
- Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Lukasz Kaiser, and Illia Polosukhin.
 2017. Attention is all you need. arXiv preprint arXiv:1706.03762.
- Barret Zoph, Deniz Yuret, Jonathan May, and Kevin Knight. 2016. Transfer learning for low-resource neural machine translation. *arXiv preprint arXiv:1604.02201.*