

# Knowledge Based Template Machine Translation In Low-Resource Setting

Anonymous ACL submission

## Abstract

Incorporating tagging into neural machine translation (NMT) systems has shown promising results in helping translate rare words such as named entities (NE). However, translating NE in low-resource setting remains a challenge. In this work, we investigate the effect of using tags and NE hypernyms from knowledge graphs (KGs) in parallel corpus in different level of resource conditions. We find the tag-and-copy mechanism (tag the NEs in the source sentence and copy them to the target sentence) improves translation in high-resource settings only. Introducing copying also results in polarizing effects in translating different parts-of-speech (POS). Interestingly, we find that copy accuracy for hypernyms is consistently higher than that of entities. As a way of avoiding "hard" copying and utilizing hypernym in bootstrapping rare entities, we introduced a "soft" tagging mechanism and found consistent improvement in high and low-resource setting.

## 1 Introduction

NMT methods usually require significant training data. For low-resource languages, NMT models generally do not work as well, especially when translating NEs. With low occurrences and large variations, NEs often remain unseen until inference time. In this paper, we investigate the usefulness of using template tagging methods and hypernyms to generalize NMT under low-resource settings.

**Template Machine Translation** Template NMT usually involves tagging the input sentences such that the templates simplify the task for the model during translation. One of the first works addressing rare entities in translation uses multiple numbered unknown (*unks*) tokens to link up source and target sentences (Luong et al., 2015). With the introduction of such copy mechanism, models only need to copy (instead of translate) the unknown token from source to target sentence,

and (if needed) perform post-processing to replace the copied-over tags. Li et al. (2018a) replaces named entities with their type symbols (i.e. LOC, ORG) on both source and target side, and trains a character-level sequence to sequence model for NE translation. Crego et al. (2016) and Wang et al. (2017) use similar tagging mechanism, with the latter using a dictionary to translate tagged NE. Wang et al. (2019) and Li et al. (2018b) use a few tagging methods from code-switching, boundary tags (i.e. `<ORG>`, `<\ORG>`), to extra embedding to tag NE on both source and target side. Others have explored encouraging copying through constrained decoding (Hokamp and Liu, 2017, Post and Vilar, 2018), or modifying architecture or input format (Gu et al., 2018, Pham et al., 2018 Dinu et al., 2019).

**Knowledge Augmented Translation** In addition to tagging boundaries of NEs from previous section, a few methods also use POS and other linguistic features to improve NMT (Sennrich and Haddow, 2016, Modrzejewski et al., 2020, Hämäläinen and Alnajjar, 2019). Anwarus Salam et al. (2017) uses hypernyms in a statistical machine translation system for low-resource translation. Meanwhile, many have used KGs to improve NMT systems. Some use KGs for data augmentation (Zhao et al., 2021), while others combine NMT with knowledge graph embedding to improve translation quality (Lu et al., 2018, Zhao et al., 2020, Moussallem et al. (2019).

While our goal resembles similar efforts in template machine translation, we extend the tag types to a much wider range using hypernyms obtained through KGs. In addition, we perform extensive analysis to understand the pros and cons of copy mechanism under different resource conditions. Our paper provides 3 key insights:

- Copy mechanism improves translation only in high-resource setting.
- Copy models translate syntactic POS better

081 and semantic POS worse, yielding translation  
082 with similar sentence structures as the source.  
083 • Appending hypernyms to NEs can improve  
084 translation accuracy in low-resource settings.

## 085 2 Methods

086 We first use statistical alignment (FastAlign, Dyer  
087 et al., 2013) to build a phrase translation table. We  
088 then use DBpedia Spotlight entity linking system  
089 (Mendes et al., 2011)<sup>1</sup>) to find NEs within sen-  
090 tences that connects to English DBpedia<sup>2</sup>, as well  
091 as the translation of the NEs on target side through  
092 translation alignment. We substitute the NEs with  
093 corresponding templates. After model translation,  
094 we remove the tag<sup>3</sup>, either keep the translation al-  
095 ready in the tag or use the phrase translation table  
096 to translate copied entities. The system is modular  
097 and all code can be found in our repo<sup>4</sup>.

098 **Tagging Methods** We use the following tem-  
099 plates in our experiments (Table 1): **Tag** and **Trans**  
100 are similar to previous works shown to improve  
101 translation adequacy (Wang et al., 2019, Li et al.,  
102 2018b). In addition to the two methods, we also ex-  
103 periment with adding entity’s hypernym provided  
104 by DBpedia. Since hypernym is a more generalized  
105 term for the entity with higher term frequency, we  
106 expect the model to use it as context when trans-  
107 lating the sentence in addition to using it to copy.  
108 **Add** adds hypernym after entity tag, **TransA**  
109 adds hypernym after tag and translation, while **TransR**  
110 replaces original entity with hypernym and adds  
111 translation. For the target sentences, we replace  
112 the NE translations with the same templates as the  
113 source sentences.

114 In addition to enforcing a "hard" copying mech-  
115 anism using tagging templates, we also include a  
116 "soft" signal by simply adding the hypernym after  
117 the entity (**HypA**). On the target side, we append  
118 the translated hypernym if possible (from phrase  
119 translation table) otherwise we use the source lan-  
120 guage hypernym for **HypA**. Without direct signal  
121 for copying, we expect the model to rely on the  
122 hypernyms as context when translating NEs.

123 In our experiments, we ensure the same NEs are  
124 tagged across templates, with about 25% of all sen-  
125 tences tagged in each dataset (Appendix Table 6).

<sup>1</sup><https://www.dbpedia-spotlight.org/>

<sup>2</sup><https://www.dbpedia.org/>

<sup>3</sup>Our soft tagging approach, **HypA**, does not contain ex-  
plicit tag and requires no removal post translation

<sup>4</sup>Anonymized. Our code is included in a zip file as software  
component in the submission

## 2.1 NMT Model

126 For NMT model, we used XLM introduced by Con-  
127 neau et al. (2020)<sup>5</sup>. We use the same transformer  
128 architecture as Wang et al. (2019): 512 embedding  
129 size, 6 encoder and decoder layer, 8 multi-attention  
130 heads. Refer to Appendix Section A.5 for more de-  
131 tails. We train on both source  $\rightarrow$  target and target  
132  $\rightarrow$  source direction.  
133

## 3 Experiments

134 In order to evaluate our results in different resource  
135 amount settings, we test our methods in English-  
136 Chinese as well as English-Hausa. For English-  
137 Chinese, we randomly select 3 million pairs of  
138 sentences from MultiUN (Ziemski et al., 2016) as  
139 training dataset in high-resource setting. To eval-  
140 uate English-Chinese translation, we use WMT  
141 newstest datasets from 2017-2020. For English-  
142 Hausa, we combine available parallel corpus on  
143 WMT-21 website<sup>6</sup> including ParaCrawl (Bañón  
144 et al., 2020), Wikititles, Khamenei corpus, and  
145 English-Hausa Opus corpus (Tiedemann, 2012), in  
146 total of 740K parallel sentences. For simulated low-  
147 resource condition, we randomly sample 6K sen-  
148 tences from English-Hausa training set. We eval-  
149 uate English-Hausa translation on newsdev2021  
150 and newstest2021. We treat the WMT newstest as  
151 the out-of-domain datasets, and randomly select 5K  
152 valid and 5K test sentences as in-domain evaluation  
153 sets from each training dataset.  
154

155 Other than evaluating translation results with  
156 multi-BLEU metric, we also investigate the accu-  
157 racy of the copy mechanism. We report the copy  
158 accuracy for hypernym, entity, and entity transla-  
159 tion whenever possible. Additionally, we calculate  
160 the word translation accuracy by POS occurring  
161 before and after the tagged entity to observe the ef-  
162 fect of copying on the rest of the sentence. We use  
163 *en\_core\_web\_sm* and *zh\_core\_web\_sm* in SpaCy  
164 library for POS tagging. For Hausa, since there  
165 is not an available POS tagger, we use alignment  
166 file from FastAlign and project English POS to  
167 corresponding words in Hausa sentence, following  
168 Rasooli et al. (2021).

## 4 Results

### 4.1 English-Chinese (High-Resource)

170 **Tagging Improves Adequacy and Accuracy**  
171 We can see a clear improvement of around 1-4  
172

<sup>5</sup><https://github.com/facebookresearch/xlm>

<sup>6</sup><https://www.statmt.org/wmt21/translation-task.html>

<b>Base.</b>	myanmar was a highly civilized country.
<b>Tag</b>	<start> myanmar <end> was a highly civilized country.
<b>Add</b>	<start> myanmar <mid> state <end> was a highly civilized country.
<b>Trans</b>	<start> myanmar <mid> 缅甸 <end> was a highly civilized country.
<b>TransA</b>	<start>myanmar <mid1> 缅甸 <mid2>state <end> was a highly civilized country.
<b>TransR</b>	<start> state <mid> 缅甸 <end> was a highly civilized country.
<b>HypA</b>	myanmar state was a highly civilized country.

Table 1: Tagging Templates for English-Chinese source sentence. NE (in red) are replaced with templates (underlined), NE hypernyms are in blue and NE translations are in green. Best viewed in color.

Method	In-Domain	Out-of-Domain
Baseline (all)	33.30 ± 0.63	11.09 ± 0.78
- (tag-only)	34.64 ± 2.1	12.21 ± 0.81
Tag (all)	33.77 ± 0.24	11.26 ± 0.91
- (tag-only)	36.07 ± 0.28	12.89 ± 1.34
Add (all)	33.69 ± 0.21	11.29 ± 0.81
- (tag-only)	35.77 ± 0.36	12.89 ± 1.11
Trans (all)	33.77 ± 0.04	11.25 ± 0.90
- (tag-only)	35.80 ± 0.48	12.97 ± 1.00
TransA (all)	33.35 ± 0.28	11.32 ± 0.83
- (tag-only)	35.37 ± 0.65	13.03 ± 0.98
TransR (all)	33.84 ± 0.29	11.18 ± 0.87
- (tag-only)	35.73 ± 0.61	12.75 ± 0.88
HypA (all)	34.39 ± 0.14	11.48 ± 0.87
- (tag-only)	<b>37.54</b> ± 0.07	<b>13.69</b> ± 0.95

Table 2: Average and standard deviation of BLEU scores across evaluation sets for all tagging methods in English-Chinese. Evaluation is performed on whole dataset (**all**) and on tagged sentences only (**tag-only**). Best performances in tag-only subsets are in bold. Best performances in all datasets are underscored. See results for individual datasets in Appendix Table 7

BLEU point on average (Table 2). The improvements are much larger when we evaluate it on tag-only subsets. **HypA** outperforms other methods consistently. Similar trend is observed in Chinese-English Translation (see Appendix Table 8).

When looking at translation accuracy (Table 3) of the tagged NEs, we see about 35 points improvement in translation accuracy. This is expected because copying is much easier than translating. **HypA** method, while performing better in BLEU, does not improve NE translation accuracy as much because it does not enforce "hard" copying. **Tag** method performs best in translating NEs with 91.92% accuracy (assuming perfect phrase translation table).

Method	Entity	Translation	Hypernym
Baseline	-	55.38	-
Tag	91.92	-	-
Add	91.02	-	92.04
Trans	<b>92.12</b>	90.99	-
TransA	91.83	<b>91.27</b>	<b>92.97</b>
TransR	-	89.12	91.66
HypA	-	55.76	58.69

Table 3: Copy accuracy (percentage) mean for different parts of the tag in English-Chinese across evaluation sets. We equate correct NE translation in baseline to correct translation copy. The hypernym translation accuracy for **HypA** is approximated with the word translation accuracy after the entity.

**Effects of Copy Mechanism on Translation** As seen in Figure 1, copying provides benefits and downfalls. It improves translation accuracy for POSs which serve as structural syntactic signals in sentences such as conjunctions, particles, punctuation while decreasing accuracy for POSs containing more semantic information that require more context to translate (verb, adjective, adverb). Qualitatively, this is equivalent to producing translations with similar sentence structures to source sentence (Appendix Table 10). Since copying is a direct signal for models to ignore context and translate word by word for the entity, it is not surprising to see such polarizing effects on the rest of the sentences. Unexpectedly, despite being a "soft" copy signal, **HypA** shows similar effects. We suspect that the repeating semantic of appending hypernyms after NEs yields similar signal for models to follow word-by-word order sensitive translation.

Similarly in Table 2, we do not see significant BLEU improvement of tagging methods that contain hypernym (**Add**, **TransA**, **TransR**) over those that do not (**Tag**, **Trans**). We believe, by the same mechanism described above, the copy mechanism shifts models' priority from using the semantics of the hypernym to simply copying the word.

## 4.2 English-Hausa (Medium-Resource)

Full English-Hausa yields similar results as English-Chinese, except that the improvements in BLEU from tagged models over baseline become marginal (Appendix Table 11). **HypA** and **Tag** performs best in-domain while baseline performs best out-of-domain. Additionally, copy accuracy decreases from 90% to 80%, but remains 20% higher over baseline accuracy (Appendix Table 12). **Tag** still outperforms other methods in copy accuracy.

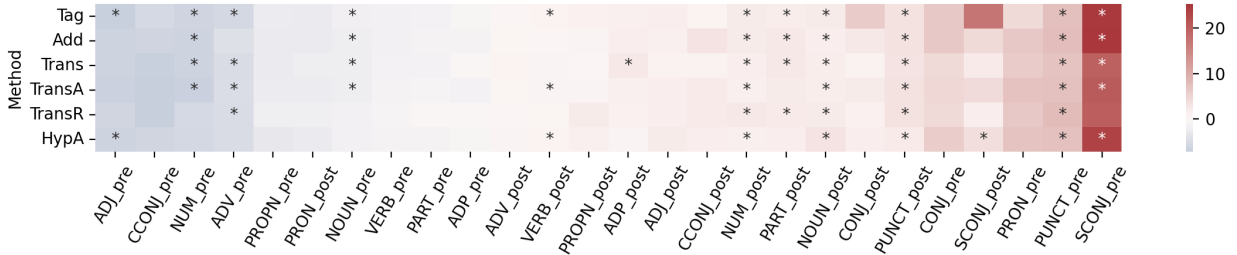


Figure 1: POS translation accuracy (percentage) difference against baseline before (**\_pre**) and after (**\_post**) the tagged entity in English-Chinese. \* indicates a statistical significant difference against baseline with p-value < 0.05

### 4.3 6K English-Hausa (Low-Resource)

Method	In-Domain	Out-of-Domain
Baseline	7.61 ± 0.21	3.80 ± 3.37
- (tag-only)	7.21 ± 0.85	3.40 ± 2.87
Tag (all)	7.39 ± 0.14	3.67 ± 3.12
- (tag-only)	6.69 ± 0.79	3.39 ± 3.13
Trans (all)	7.45 ± 0.08	3.91 ± 3.44
- (tag-only)	6.99 ± 0.92	<b>3.60</b> ± 3.44
HypA (all)	7.53 ± 0.25	3.52 ± 2.88
- (tag-only)	<b>7.82</b> ± 1.40	2.55 ± 1.89

Table 4: BLEU scores for 6K English-Hausa data. Only top performing methods are included.

Method	entity	translation	hypernym
Baseline	-	42.44	-
Tag	30.72↓	-	-
Add	34.48↓	-	55.66
Trans	37.81	35.69↓	-
TransA	<b>39.01</b>	37.53	<b>55.91</b>
TransR	-	30.61↓	55.39
HypA	-	<b>44.77</b> ↑	48.32

Table 5: Copy accuracy (mean) in 6K English-Hausa dataset models across evaluation sets. Arrows indicate statistical difference from baseline with p-value < 0.05.

In low-resource setting, tagging does not improve performance (Table 4). The NE copy accuracy drops below baseline. Interestingly, hypernyms are more consistently copied to the target side (Table 5). We believe this is due to hypernyms having higher term frequency in the training data. Compared to baseline, only **HypA** method is able to improve NE translation accuracy and obtain higher BLEU for tag-only subsets in-domain (Table 4). Despite not having as high of hypernym copy accuracy, the model is able to use hypernym as context, and improve NE translation.

## 5 Discussion

**Copy mechanism in low-resource?** As results show, copy mechanism is able to increase NE trans-

lation accuracy in both high and medium-resource but not in low-resource condition. Learning to copy requires significant amount of data. Once tags are recognized, the semantics of the content within are ignored. Translations become structurally similar to source sentence, while focusing less on semantics of words that depend on the context. Without enough data, "softer" methods of augmentation (**HypA** or extra embedding used by (Moussallem et al., 2019)) that incorporates hypernym in translation is a better choice. Work by (Currey et al., 2017), which copies target sentences to source side to create additional bitext, might be interesting alternatives to encourage copying.

**Low-Resource translation affected by term frequency.** As suggested by Table 5, before copy mechanism generalizes, models are more likely to copy words that occur more frequently (hypernyms). This points to potential directions in low-resource NLP in using hypernyms to bootstrap performance of other words or sentences. Data augmentation techniques like randomly inserting/replacing NEs with hypernyms could be potential ways of adding data points in low-resource settings and better generalize embedding space.

## 6 Conclusion

In our paper, we analyzed the tag-and-copy mechanism under different resource conditions. We found that learning to copy requires significant amount of resource which is often not achievable in low-resource languages. Additionally, we found that copying can induce polarizing effects on translating different POSs. It discouraged models from using contextual information, but provide "structural supervision". In low-resource setting, we found correlation between term frequency and copying accuracy. Our proposed method of appending hypernym after NEs was able to encourage better translation in both low and high-resource setting.

279	<b>Acknowledgements</b>		
280	<b>References</b>		
281	Khan Md Anwarus Salam, Setsuo Yamada, and Nishino	Chris Hokamp and Qun Liu. 2017. Lexically con-	335
282	Tetsuro. 2017. <a href="#">Improve example-based machine</a>	strained decoding for sequence generation using grid	336
283	<a href="#">translation quality for low-resource language using</a>	beam search. In <i>Proceedings of the 55th Annual</i>	337
284	<a href="#">ontology</a> . <i>International Journal of Networked and</i>	<i>Meeting of the Association for Computational Lin-</i>	338
285	<i>Distributed Computing</i> , 5:176.	<i>guistics (Volume 1: Long Papers)</i> , pages 1535–1546.	339
286	Marta Bañón, Pinzhen Chen, Barry Haddow, Kenneth	Belinda Z. Li, Sewon Min, Srinivasan Iyer, Yashar	340
287	Heafield, Hieu Hoang, Miquel Esplà-Gomis, Mikel L.	Mehdad, and Wen-tau Yih. 2020. Efficient one-pass	341
288	Forcada, Amir Kamran, Faheem Kirefu, Philipp	end-to-end entity linking for questions. In <i>EMNLP</i> .	342
289	Koehn, Sergio Ortiz Rojas, Leopoldo Pla Sempere,	Xiaoqing Li, Jinghui Yan, Jiajun Zhang, and Chengqing	343
290	Gema Ramírez-Sánchez, Elsa Sarrías, Marek Strelec,	Zong. 2018a. Neural name translation improves neu-	344
291	Brian Thompson, William Waites, Dion Wiggins, and	ral machine translation. In <i>China Workshop on Ma-</i>	345
292	Jaume Zaragoza. 2020. <a href="#">ParaCrawl: Web-scale acqui-</a>	<i>chine Translation</i> , pages 93–100. Springer.	346
293	<a href="#">sition of parallel corpora</a> . In <i>Proceedings of the 58th</i>	Zhongwei Li, Xuancong Wang, AiTi Aw, Eng Siong	347
294	<i>Annual Meeting of the Association for Computational</i>	Chng, and Haizhou Li. 2018b. Named-entity tagging	348
295	<i>Linguistics</i> , pages 4555–4567, Online. Association	and domain adaptation for better customized transla-	349
296	for Computational Linguistics.	tion. In <i>Proceedings of the Seventh Named Entities</i>	350
297	Alexis Conneau, Kartikay Khandelwal, Naman Goyal,	<i>Workshop</i> , pages 41–46.	351
298	Vishrav Chaudhary, Guillaume Wenzek, Francisco	Yu Lu, Jiajun Zhang, and Chengqing Zong. 2018. Ex-	352
299	Guzmán, Edouard Grave, Myle Ott, Luke Zettle-	ploiting knowledge graph in neural machine trans-	353
300	moyer, and Veselin Stoyanov. 2020. Unsupervised	lation. In <i>China Workshop on Machine Translation</i> ,	354
301	cross-lingual representation learning at scale. In	pages 27–38. Springer.	355
302	<i>ACL</i> .	Minh-Thang Luong, Ilya Sutskever, Quoc Le, Oriol	356
303	Josep Crego, Jungi Kim, Guillaume Klein, Anabel Re-	Vinyals, and Wojciech Zaremba. 2015. Addressing	357
304	bollo, Kathy Yang, Jean Senellart, Egor Akhanov,	the rare word problem in neural machine translation.	358
305	Patrice Brunelle, Aurelien Coquard, Yongchao Deng,	In <i>Proceedings of the 53rd Annual Meeting of the As-</i>	359
306	et al. 2016. Systran’s pure neural machine translation	<i>sociation for Computational Linguistics and the 7th</i>	360
307	systems. <i>arXiv preprint arXiv:1610.05540</i> .	<i>International Joint Conference on Natural Language</i>	361
308	Anna Currey, Antonio Valerio Miceli-Barone, and Ken-	<i>Processing (Volume 1: Long Papers)</i> , pages 11–19.	362
309	neth Heafield. 2017. Copied monolingual data im-	Pablo N. Mendes, Max Jakob, Andres Garcia-Silva, and	363
310	proves low-resource neural machine translation. In	Christian Bizer. 2011. Dbpedia spotlight: Shedding	364
311	<i>Proceedings of the Second Conference on Machine</i>	light on the web of documents. In <i>Proceedings of the</i>	365
312	<i>Translation</i> , pages 148–156.	<i>7th International Conference on Semantic Systems</i>	366
313	Georgiana Dinu, Prashant Mathur, Marcello Federico,	<i>(I-Semantics)</i> .	367
314	and Yaser Al-Onaizan. 2019. Training neural ma-	Maciej Modrzejewski, Miriam Exel, Bianka Buschbeck,	368
315	chine translation to apply terminology constraints. In	Thanh-Le Ha, and Alex Waibel. 2020. Incorporating	369
316	<i>Proceedings of the 57th Annual Meeting of the Asso-</i>	external annotation to improve named entity trans-	370
317	<i>ciation for Computational Linguistics</i> , pages 3063–	lation in nmt. In <i>Proceedings of the 22nd Annual</i>	371
318	3068.	<i>Conference of the European Association for Machine</i>	372
319	Chris Dyer, Victor Chahuneau, and Noah A Smith. 2013.	<i>Translation</i> , pages 45–51.	373
320	A simple, fast, and effective reparameterization of	Diego Moussallem, Mihael Arčan, Axel-Cyrille Ngonga	374
321	ibm model 2. In <i>Proceedings of the 2013 Conference</i>	Ngomo, and Paul Buitelaar. 2019. Augmenting	375
322	<i>of the North American Chapter of the Association</i>	neural machine translation with knowledge graphs.	376
323	<i>for Computational Linguistics: Human Language</i>	<i>arXiv preprint arXiv:1902.08816</i> .	377
324	<i>Technologies</i> , pages 644–648.	Ngoc-Quan Pham, Jan Niehues, and Alex Waibel. 2018.	378
325	Jiatao Gu, Yong Wang, Kyunghyun Cho, and Victor OK	Towards one-shot learning for rare-word translation	379
326	Li. 2018. Search engine guided neural machine trans-	with external experts. In <i>Proceedings of the 2nd</i>	380
327	lation. In <i>Proceedings of the AAAI Conference on</i>	<i>Workshop on Neural Machine Translation and Gen-</i>	381
328	<i>Artificial Intelligence</i> , volume 32.	<i>eration</i> , pages 100–109.	382
329	Mika Härmäläinen and Khalid Alnajjar. 2019. A	Matt Post and David Vilar. 2018. Fast lexically con-	383
330	template based approach for training nmt for low-	strained decoding with dynamic beam allocation for	384
331	resource uralic languages-a pilot with finnish. In	neural machine translation. In <i>Proceedings of the</i>	385
332	<i>Proceedings of the 2019 2nd International Conference</i>	<i>2018 Conference of the North American Chapter of</i>	386
333	<i>on Algorithms, Computing and Artificial Intelligence</i> ,	<i>the Association for Computational Linguistics: Hu-</i>	387
334	pages 520–525.	<i>man Language Technologies, Volume 1 (Long Pa-</i>	388
		<i>pers)</i> , pages 1314–1324.	389

390 Mohammad Sadegh Rasooli, Chris Callison-Burch, and  
391 Derry Tanti Wijaya. 2021. “wikily” supervised neu-  
392 ral translation tailored to cross-lingual tasks. In *Pro-*  
393 *ceedings of the 2021 Conference on Empirical Meth-*  
394 *ods in Natural Language Processing*, pages 1655–  
395 1670.

396 Rico Sennrich and Barry Haddow. 2016. Linguistic  
397 input features improve neural machine translation.  
398 In *Proceedings of the First Conference on Machine*  
399 *Translation: Volume 1, Research Papers*, pages 83–  
400 91.

401 Jörg Tiedemann. 2012. Parallel data, tools and inter-  
402 faces in opus. In *Proceedings of the Eight Inter-*  
403 *national Conference on Language Resources and*  
404 *Evaluation (LREC’12)*, Istanbul, Turkey. European  
405 Language Resources Association (ELRA).

406 Tao Wang, Shaohui Kuang, Deyi Xiong, and António  
407 Branco. 2019. Merging external bilingual pairs  
408 into neural machine translation. *arXiv preprint*  
409 *arXiv:1912.00567*.

410 Yuguang Wang, Shanbo Cheng, Liyang Jiang, Jiajun  
411 Yang, Wei Chen, Muze Li, Lin Shi, Yanfeng Wang,  
412 and Hongtao Yang. 2017. Sogou neural machine  
413 translation systems for wmt17. In *Proceedings of the*  
414 *Second Conference on Machine Translation*, pages  
415 410–415.

416 Yang Zhao, Lu Xiang, Junnan Zhu, Jiajun Zhang,  
417 Yu Zhou, and Chengqing Zong. 2020. Knowledge  
418 graph enhanced neural machine translation via multi-  
419 task learning on sub-entity granularity. In *Proceed-*  
420 *ings of the 28th International Conference on Compu-*  
421 *tational Linguistics*, pages 4495–4505.

422 Yang Zhao, Jiajun Zhang, Yu Zhou, and Chengqing  
423 Zong. 2021. Knowledge graphs enhanced neural ma-  
424 chine translation. In *Proceedings of the Twenty-Ninth*  
425 *International Conference on International Joint Con-*  
426 *ferences on Artificial Intelligence*, pages 4039–4045.

427 Michał Ziemski, Marcin Junczys-Dowmunt, and Bruno  
428 Pouliquen. 2016. [The United Nations parallel corpus v1.0](#). In *Proceedings of the Tenth International*  
429 *Conference on Language Resources and Evaluation*  
430 *(LREC’16)*, pages 3530–3534, Portorož, Slovenia.  
431 European Language Resources Association (ELRA).  
432

## 433 A Appendix

### 434 A.1 Text Preprocessing

435 We follow default preprocessing steps in XLM  
436 repo. For English and Hausa, we use Moses *tok-*  
437 *enizer.perl* script, after which we lower-case letters  
438 and remove accents. For Chinese, we use Moses  
439 *tokenizer\_PTB.perl* script.

### A.2 Special Tags in XLM Model

440 During tagging, in order to prevent creating addi-  
441 tional vocabulary, we use four of the special to-  
442 kens (i.e. <special2>, <special3>, <special4>,  
443 <special5>), that already exist in pretrained XLM-  
444 R model vocab, instead of actual <start>, <end>,  
445 etc.  
446

### A.3 Tagging Statistics

Language Pair	Train Size	Tag Size
English-Hausa	6 K	1.5 K (25.6%)
English-Hausa	746 K	191 K (25.6%)
English-Chinese	2,990 K	816 K (27.3%)

Table 6: Tagging Statistics in Training Sets

### A.4 Entity Linking

448 During experimentation, we have also tried more  
449 recent Entity linking systems such as BLINK (Li  
450 et al., 2020)<sup>7</sup>. In reality, we find BLINK tagging  
451 less entities as well as taking a longer time. We  
452 presume this is because BLINK expects normally-  
453 cased sentences while our entity linking occurs  
454 after input sentences are lower-cased.  
455

### A.5 Model Training Details

456 In all of our experiments, we use the pretrained  
457 XLM-R BPE vocab with 200,000 tokens, trained  
458 on 100 lanugages<sup>8</sup>. We use Adam optimizer, learn-  
459 ing rate 0.0001, epoch size 300000, dropout rate  
460 of 0.1. We fix number of tokens in a batch to be  
461 around 2000. To increase batch size with GPU  
462 memory constraint, we use gradient accumulation  
463 for every four batches to increase effective batch  
464 size. For low-resource condition with 6K train-  
465 ing sentences (see Section 3), we change epoch  
466 size to 120,000, dropout of 0.2, and enforce mini-  
467 mum sentence length to 10 words. All models are  
468 trained on NVIDIA V100 GPUs. Each English-  
469 Chinese model takes about 5 days to train (1 GPU  
470 time). Each English-Hausa model takes about 3  
471 days and each English-Hausa 6K model takes about  
472 15 hours.  
473

Method	subset	valid	test	nd2017	nt2017	nt2018	nt2019	nt2020	ntB2020
Baseline	all	32.85	33.75	11.23	10.77	11.02	10.20	12.54	10.78
Baseline	tag-only	33.15	36.12	13.22	12.69	12.13	11.30	12.69	11.20
Tag	all	33.59	33.94	11.20	11.38	<u>11.34</u>	10.14	12.85	10.66
Tag	tag-only	35.86	36.27	13.72	14.20	13.18	11.16	13.85	11.25
Add	all	33.53	33.84	11.15	<u>11.58</u>	11.19	10.36	12.71	10.72
Add	tag-only	35.51	36.03	13.25	14.48	12.88	12.17	13.33	11.20
Trans	all	33.74	33.80	11.23	11.10	10.72	<u>10.73</u>	13.04	10.68
Trans	tag-only	35.45	36.14	13.46	13.97	12.40	12.34	<b>14.04</b>	11.59
TransA	all	33.14	33.55	11.10	11.33	11.28	10.47	12.89	10.85
TransA	tag-only	34.90	35.83	13.50	13.72	<b>13.54</b>	12.02	13.84	11.53
TransR	all	33.63	34.05	11.10	11.08	11.18	10.31	12.82	10.61
TransR	tag-only	35.29	36.16	13.32	13.65	12.63	11.85	13.46	11.56
HypA	all	<u>34.29</u>	<u>34.39</u>	<u>11.31</u>	11.51	11.17	<u>10.73</u>	<u>13.18</u>	<u>10.99</u>
HypA	tag-only	<b>37.49</b>	<b>37.59</b>	<b>14.67</b>	<b>14.73</b>	13.49	<b>13.28</b>	13.76	<b>12.18</b>

Table 7: BLEU scores across evaluation sets for all tagging methods in English-Chinese. Evaluation is performed on whole dataset and on tagged sentences only. Best performances in tagged subset are in bold. Best performances in all datasets are underscored. Each point represents a single data point. (nd2017=newsdev2017, nt2017=newstest2017, etc)

## A.6 English-Chinese Full Results

## A.7 Chinese-English Translation Results

## A.8 Copy Efficiency In / Out of Domain

In English-Chinese translation results, we can observe that the copy accuracy for the tags is similar across different set regardless of the domain (Table 9), which is a good sign considering the drop in BLEU across the out-of-domain datasets. This indicate copy mechanism is a valuable method in translation avenues where entity translation accuracy is more valuable than adequacy (i.e. medical, scientific domain), confirming with results in Pham et al. (2018) and Dinu et al. (2019).

## A.9 English-Hausa POS Accuracy Qualitative Analysis

## A.10 English-Hausa 6K Translation Results

<sup>7</sup><https://github.com/facebookresearch/BLINK>

<sup>8</sup>See <https://github.com/facebookresearch/XLMthe-17-and-100-languages> for language details

Method	subset	valid	test	nd2017	nt2017	nt2018	nt2019	nt2020	ntB2020
Baseline	all	38.46	42.33	12.06	12.74	13	10.37	12.13	11.65
Baseline	tag-only	43.28	44.87	13.01	13.81	14.16	11	12.88	12.47
Tag	all	41.47	42.56	12.53	12.76	13.06	10.55	12.48	<u>11.84</u>
Tag	tag-only	44.01	45.13	14.51	13.87	14.57	11.94	13.43	13.17
Add	all	41.42	42.37	12.76	13.14	12.74	10.38	12.46	11.83
Add	tag-only	43.82	44.86	14.73	14.11	14.33	11.54	<b>13.67</b>	<b>13.26</b>
Trans	all	41.31	42.42	12.35	13	13.17	10.42	12.21	11.61
Trans	tag-only	43.4	44.8	13.84	14.26	14.96	12.14	13.26	13.02
TransA	all	41.1	42.17	12.76	<u>13.21</u>	13.13	10.66	12.07	11.52
TransA	tag-only	42.99	44.39	14.3	<b>14.69</b>	14.84	12.24	13.12	12.72
TransR	all	41.21	42.28	<u>12.8</u>	13.03	12.88	<u>10.75</u>	<u>12.52</u>	11.81
TransR	tag-only	43.49	44.75	<b>15.03</b>	14.26	14.69	12.26	13.39	12.82
HypA	all	<u>41.84</u>	<u>42.99</u>	12.47	12.98	<u>13.29</u>	10.48	12.2	11.68
HypA	tag-only	<b>45.32</b>	<b>46.08</b>	14.76	14.55	<b>15.07</b>	<b>12.62</b>	13.18	13.23

Table 8: BLEU scores across evaluation sets for all tagging methods in Chinese-English. There is a consistent 0.5-2 point improvement with tagged methods over baseline. Each point represents a single data point.

	Valid	Test	nd2017	nt2017
H	91.98	90.92	94.88	97.19
E	91.84	90.5	92.79	91.8
T	91.91	90.15	93.17	93.91
	nt2018	nt2019	nt2020	ntB2020
H	94.45	94.16	88.48	91.73
E	93.76	93.67	88.02	92.27
T	92.37	92.94	86.41	89.33

Table 9: Copy Accuracy of TransA model across different in and out-of-domain evaluation datasets. Each point represents a single data point. H=Hypernym, E=Entity, T=Entity translation



<b>Label</b>	in the <b>gambia</b> 's interim paper , it was noted that major factors in poverty among rural women include their predominance in subsistence agriculture , where they have less access than men to mechanized technologies , and the fact that , in addition to farming , they work longer hours than men carrying out household tasks .
<b>Baseline</b>	the interim document of the <b>gambia</b> indicated that rural women 's poverty was mainly due to their livelihood agriculture , which was less skilled than men ; and that they were more time spent than men to run their household than men , in addition to their work .
<b>Tag</b>	the <u>&lt;special2&gt; <b>gambia</b> &lt;special5&gt;</u> interim paper indicated that the main cause of poverty among rural women was their main livelihood agriculture , less access to mechanized technologies than men ; and that in addition to farming , they were more time-consuming than men .
<b>Add</b>	the <u>&lt;special2&gt; <b>gambia</b> &lt;special3&gt; <b>country</b> &lt;special5&gt;</u> 's interim paper noted that the main causes of poverty among rural women were their primary work in subsistence agriculture , more than men 's access to mechanical techniques , and that they would have more time than men to take their household roles in addition to their farm .
<b>Trans</b>	the <u>&lt;special2&gt; <b>gambia</b> &lt;special3&gt; <b>冈比亚</b> &lt;special5&gt;</u> 's provisional document noted that the main causes of poverty among rural women are their primary subsistence agriculture , less than men 's access to mechanized technologies , and that in addition to their farm , they are more time than men to operate household .
<b>TransA</b>	in the <u>&lt;special2&gt; <b>gambia</b> &lt;special3&gt; <b>冈比亚</b> &lt;special4&gt; <b>country</b> &lt;special5&gt;</u> 's interim paper , it was noted that major factors in poverty among rural women include their predominance in subsistence agriculture , where they have less access than men to mechanized technologies , and the fact that , in addition to farming , they work longer hours than men carrying out household tasks .
<b>TransR</b>	the provisional document of the <u>&lt;special2&gt; <b>country</b> &lt;special3&gt; <b>冈比亚</b> &lt;special5&gt;</u> indicates that the main causes of poverty among rural women are their predominance in livelihood agriculture , less access to mechanized technologies than men , and that they are more time than men to take up their housework in addition to their agricultural work .
<b>HypA</b>	the interim document of the <b>gambia country</b> indicated that the main reason for poverty among rural women was their predominant livelihood farming , less than the mechanized technique of access to men ; and that they were also taking more time than men to operate their household tasks .

Table 10: Translation example before post-translation tag removal. In Chinese-English translation setting, we compare all model translation results with ground truth English sentence. In all tagging methods, models tend to produce more similar sentence structures due to similar syntactic word choices. Given fixed sentence structures, there is less emphasis on translating the rest of the words that contain more semantic variations (verbs, adjectives, adverbs, etc.). NE (in red) are replaced with templates (underlined), NE hypernyms are in blue and NE translations are in green. Best viewed in color.

Method	valid	test	nd2021	nt2021
Base(all)	32.94	32.89	<u>11.31</u>	21.62
- (tag-only)	35.35	37.12	11.50	<b>23.18</b>
Tag(all)	<u>33.17</u>	32.99	10.77	<u>21.84</u>
- (tag)	<b>35.91</b>	37.28	11.86	23.13
Add (all)	32.25	32.62	11.16	21.42
- (tag-only)	34.58	36.44	12.07	22.54
Trans(all)	32.27	32.29	10.85	21.56
- (tag-only)	35.45	36.14	<b>12.01</b>	22.71
TransA	32.22	32.3	10.58	21.38
- (tag-only)	33.88	35.94	11.33	22.56
TransR	32.65	32.77	11.18	21.74
- (tag-only)	34.74	36.73	12.38	22.71
HypA(all)	33.02	<u>33.00</u>	9.59	20.24
- (tag-only)	35.89	<b>37.39</b>	8.12	15.42

Table 11: BLEU scores with English-Hausa full data. Each point represents a single data point.

Method	entity	translation	hypernym
Tag	81.93	-	-
Add	79.16	-	79.34
Trans	<b>82.10</b>	<b>81.30</b>	-
TransR	-	80.99	<b>81.86</b>
TransA	80.87	80.23	80.90
HypA	-	61.00	64.29
Baseline	-	59.56	-

Table 12: Copy accuracy mean with English-Hausa full data. Aggregated across all evaluation datasets.

Method	valid	test	nd2021	nt2021
Base (all)	<u>7.75</u>	<u>7.46</u>	1.41	6.18
- (tag-only)	6.61	7.81	<b>1.37</b>	5.43
Tag (all)	7.49	7.29	1.46	5.87
- (tag-only)	6.13	7.25	1.18	5.6
Add (all)	7.59	7.52	1.38	6.29
- (tag-only)	6.19	7.61	1.25	5.48
Trans (all)	7.51	7.39	<u>1.48</u>	6.34
- (tag-only)	6.34	7.64	1.16	6.03
TransA (all)	7.14	7.12	1.35	6.13
- (tag-only)	5.86	7.3	1.19	5.56
TransR (all)	7.35	7.32	1.4	<u>6.5</u>
- (tag-only)	5.73	7.22	1.03	5.74
HypA (all)	7.71	7.35	1.48	5.55
- (tag-only)	<b>6.83</b>	<b>8.18</b>	1.21	3.88

Table 13: BLEU scores in 6K English-Hausa data for all models across individual evaluation sets. Each point represents a single data point. nd2021=newsdev2021, nt2021=newstest2021