

# The Unreasonable Effectiveness of Random Target Embeddings for Continuous-Output Neural Machine Translation

Anonymous ACL submission

## Abstract

Continuous-output neural machine translation (CoNMT) replaces the discrete next-word prediction problem with an embedding prediction. The semantic structure of the target embedding space (*i.e.*, closeness of related words) is intuitively believed to be crucial. We challenge this assumption and show that completely random output embeddings can outperform laboriously pretrained ones, especially on larger datasets. Further investigation shows this surprising effect is strongest for rare words, due to the geometry of their embeddings. We shed further light on this finding by designing a mixed strategy that combines random and pre-trained embeddings for different tokens.

## 1 Introduction

Since text is naturally discrete, *i.e.*, each token in a target sentence is represented by an integer index in the vocabulary, neural machine translation (NMT), as many other language generation tasks, is trained mainly as a discrete-output model with softmax over the full vocabulary followed by the cross-entropy loss. Continuous-output neural machine translation (CoNMT) models, in contrast, are trained to predict the continuous representation based on the distances between vectors. It is an appealing line of study for computational and modeling related reasons (Kumar and Tsvetkov, 2019), as well as a reliable test bed for exploring the properties of continuous spaces that appear in modern deep generative models (Li et al., 2022). However, CoNMT introduces its own challenge, namely mapping to and from continuous space. During training, CoNMT model requires continuous targets, and while decoding, one needs to map back to the discrete text representation.

Text mapping to continuous space is widely explored in NLP and can be done using *embeddings* of tokens, words (Turian et al., 2010; Mikolov et al., 2013, 2018) and sentences (Reimers and Gurevych,

2019; Feng et al., 2022). Cosine similarity between word embeddings is well correlated with lexical similarity metrics, motivating the use of cosine distance against pretrained embeddings as the dominant training strategy for CoNMT. Nearest neighbor beam decoding would in this case include related words and, unlike discrete cross-entropy, the training strategy does not discourage synonyms.

Previous studies show that the quality of continuous-output models highly depends on the choice of embeddings (Li et al., 2022; Tokarchuk and Niculae, 2022; Kumar and Tsvetkov, 2019). In general, in CoNMT the embeddings are pre-trained and fixed, as otherwise making all embeddings equal leads to an unwanted global optimum. Obtaining pre-trained word embeddings can be computationally expensive, especially if one needs to train an embeddings model from scratch.

In this work we *randomly* initialize target embeddings for continuous-output models and keep them static during training. Arora et al. (2020) applied static random embeddings for text classification model’s input; however, to the best of our knowledge, the effect of untrained random target embeddings has not been previously studied in the literature, especially for text-generating tasks such as machine translation. Using random untrained embeddings as targets for training continuous-output models with distance measures confronts the idea of the semantic similarity importance. However, we show that random target embeddings perform close to their pre-trained counterpart, and even surpass them if there is enough data available. That means that meaningful semantic similarity is not the only factor contributing to the performance of the continuous-output models. We hypothesize and experimentally show that distances between embeddings play an important role for representation disentanglement. Our findings on three NMT tasks, namely WMT 2018 English→Turkish (en-tr), WMT 2016 English→Romanian (EnRo),

and WMT 2019 English  $\rightarrow$  German (en-de) indicate that random embeddings are more spread out and performing better on rare words for all language pairs. On the large-scale (en-de) CoNMT with random target embeddings are even substantially better overall. We propose simple, yet efficient combination of random and pre-trained embeddings and show that it helps improving models performance on both en-tr and ro-en

## 2 Continuous-Output NMT

The machine translation task involves learning to map sequences of input tokens  $\mathbf{x} = (x_1, \dots, x_m)$  to output tokens  $\mathbf{y} = (y_1, \dots, y_n)$ . In standard (discrete) NMT, each step is a multi-class next word prediction task, minimizing:

$$L_{\text{discrete}}(y_i = t; \mathbf{y}_{<i}, \mathbf{x}) = -\log p(y_i = t \mid \mathbf{y}_{<i}, \mathbf{x}) \\ = -\langle \mathbf{E}(t), \mathbf{h} \rangle + \log \sum_{t' \in V} \exp \langle \mathbf{E}(t'), \mathbf{h} \rangle, \quad (1)$$

where  $t$  is a token index,  $V$  is the vocabulary,  $\mathbf{E} : V \rightarrow \mathbb{R}^d$  is an embedding lookup, and  $\mathbf{h}$  is a transformer hidden state calculated in terms of  $\mathbf{x}$  and the output prefix  $\mathbf{y}_{<i}$ . The costly log-sum-exp and the penchant for continuous similarity metrics in NLP motivate a purely-continuous alternative:

$$L_{\text{cos}}(y_i = t; \mathbf{y}_{<i}, \mathbf{x}) = 1 - \cos(\mathbf{E}(t), \mathbf{h}). \quad (2)$$

Continuous NMT models were first studied by Kumar and Tsvetkov (2019), who also propose other probabilistic losses and later other margin-based objectives (Bhat et al., 2019), with limited gain and at the cost of additional hyperparameters; we therefore focus on the robust cosine objective. On the other hand, the choice of embeddings  $\mathbf{E}$  makes a much larger difference, especially due to the fact that all previous work keeps this parameter frozen: indeed, if it were trainable, Equation (2) would have trivial global optima by setting all  $\mathbf{E}(t)$  to the same vector for all  $t$ . With modern transformer architectures, the best performing embeddings overall tend to be the ‘‘oracle’’ output embeddings learned by a pretrained discrete MT system (Tokarchuk and Niculae, 2022). We highlight that the cosine loss is invariant to the norms of both the embeddings and of the decoder hidden state, and therefore we may restrict our modeling problem to the unit sphere.

Optimizing Equation (1) pushes the model  $\mathbf{h}$  away from all tokens different from the ‘‘gold’’ token, even if some other tokens (e.g., synonyms)

could otherwise be a good fit. Equation (2) has no such effect, leading to a promise of more diverse generations. An appealing intuition is that synonyms and related words being nearby in embedding space contributes to the performance of CoNMT and enables such diversity. However, this intuition is not consistent with practice. In fact, decoding is usually done by greedy nearest-neighbor lookup rather than beam search. Therefore, in this work, we challenge this conventional wisdom by considering completely random embeddings.

## 3 Random Embeddings Generation

We consider two different distributions from which to sample the  $|V|$  random embeddings.

**Spherical uniform.** We draw embeddings uniformly from the surface of the  $d$ -sphere, by drawing from a standard Gaussian and normalizing. (As the cosine loss is norm-invariant, uniform initialization is equivalent to the standard initialization of transformer embeddings.)

**Hypercube.** The corners of the hypercube  $\{-1, 1\}^d$  all have norm  $\sqrt{d}$  and thus form a discrete subset of a hypersphere. This motivates us to consider drawing embeddings from a scaled Rademacher distribution:

$$\mathbf{E}(y_i) = \mathbf{r}_i / \sqrt{d}; \quad \mathbf{r}_i \sim \text{Rademacher}(d).$$

Each coordinate of  $\mathbf{r}_i$  has 50% probability of being +1 and 50% of being -1. With this strategy, any two distinct embeddings have cosine distance at least  $2/d$ . Moreover, hypercubic embeddings can be stored as bit patterns and potentially allow for faster loss calculation with dedicated low-level implementations which we do not explore here.

## 4 Experimental Setup and Data

We train CoNMT systems against randomly-generated target embeddings and against pre-trained embeddings from discrete NMT systems.

Results are reported on three WMT translation tasks<sup>1</sup>: WMT 2016 Romanian  $\rightarrow$  English (ro-en), WMT 2018 English  $\rightarrow$  Turkish (en-tr) and WMT 2019 English  $\rightarrow$  German (en-de), the latter including back-translated data. Note that for en-tr we use only WMT 2018 training data with 207k training sentences to represent a challenging lower-resource and morphology-rich scenario. Data statistics are collected in Appendix A.

<sup>1</sup><https://www2.statmt.org/>

embeddings	en-tr		ro-en		en-de	
	BLEU $\uparrow$	BERTSc. $\uparrow$	BLEU $\uparrow$	BERTSc. $\uparrow$	BLEU $\uparrow$	BERTSc. $\uparrow$
discrete model	12.3	70.4	31.7	64.1	33.1	69.0
MTtransfer (beam=1)	10.1	67.1	29.0	58.5	31.3	66.2
MTtransfer	<b>10.4</b>	67.4	29.0	58.0	29.2	62.6
random uniform	8.9	65.1	28.8	58.8	31.8	<b>67.2</b>
random cube	8.7	64.6	28.7	58.8	31.4	66.9
combined	<b>10.4</b>	<b>68.3</b>	<b>29.5</b>	<b>60.4</b>	<b>32.0</b>	66.8

**Table 1:** BLEU and BertScore on ro-en newstest16, en-tr newstest2017 and newstest2016 en-de. We use a beam of 5 if not stated otherwise. In bold, we show the highest score among the continuous models in each column.

For subword tokenization we used the same SentencePiece (Kudo and Richardson, 2018) model for all language pairs, specifically the one used in the MBart multilingual model (Liu et al., 2020). This choice allows for unified preprocessing for all languages we cover. We validate that token-based models performs generally better than word-level models (Appendix C), even though subwords introduce an additional challenge of predicting subword continuation (Appendix C.1).

We used fairseq (Ott et al., 2019) framework for training our models. Baseline discrete models are trained with cross-entropy loss, label smoothing equal to 0.1 and effective batch size 65.5K tokens. Both discrete and continuous models are trained with learning rate  $5 \cdot 10^{-4}$ , 10k warm-up steps for ro-en and en-de, and 4k for the smaller en-tr dataset. All continuous models are trained with the cosine distance objective in Equation (2). Detailed description of training setup and parameters can be found in Appendix B.

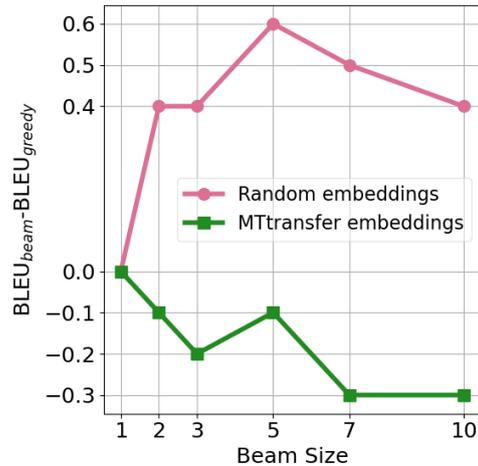
We measure translation accuracy using SacreBLEU<sup>2</sup> (Papineni et al., 2002; Post, 2018) and BertScore<sup>3</sup> (Zhang\* et al., 2020). Note that BertScore is scaled differently for each language, so the scores cannot be compared across languages.

## 5 Results and Discussion

**Scores.** Per Table 1, we find that random uniform embeddings outperform the MTtransfer baseline for en-de, match it closely for ro-en, and only underperform in the low-resource case for en-tr. We find that hypercube embeddings consistently perform worse than uniform embeddings; however, it is possible that their computational advantages can make up for this in some applications.

<sup>2</sup>nrefs:1|case:mixed|eff:no|tok:13a|smooth:exp|version:2.3.1

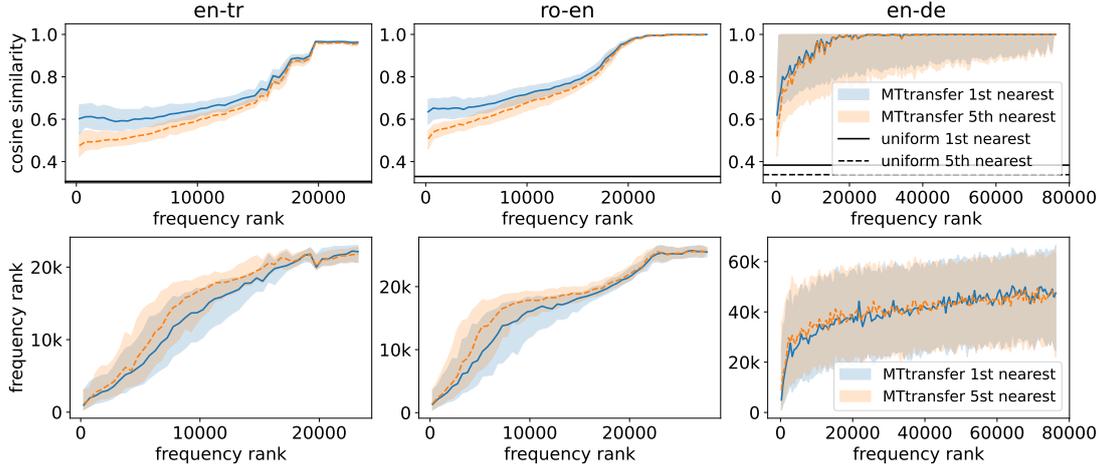
<sup>3</sup>implementation by [https://github.com/Tiiiiger/bert\\_score](https://github.com/Tiiiiger/bert_score)



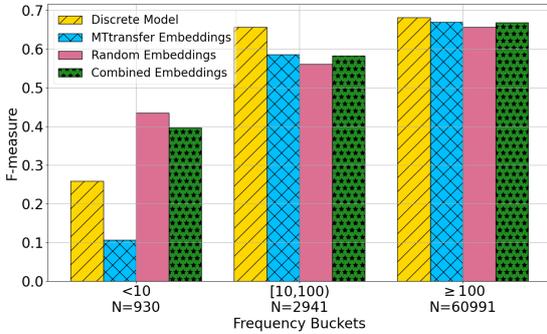
**Figure 1:** BLEU<sub>beam</sub> - BLEU<sub>greedy</sub> scores for the ro-en newsdev2016 for continuous output models with random and MTtransfer embeddings. Beam=1 BLEU score for the MTtransfer embeddings is equal to 30.0 and for uniform random embeddings 28.6

**Beam search.** Preliminary experiments with CoNMT models indicate little gain or even degradation from beam search, which is why we report results with greedy decoding for MTtransfer in Table 1. Further investigation in Figure 1 shows that the MTtransfer model degrades consistently, performing best in the greedy case, while the random embedding model benefits noticeably from a larger beam, in spite of neighboring words being random and not related. We discuss the details of the beam search in Appendix D.

**Frequency.** We perform a token-level evaluation using compare-mt (Neubig et al., 2019), computing the  $F_1$  score of matching a gold token (at its gold position), aggregated over bins defined by the token’s frequency in the training data. The result in Figure 3 reveals that random embeddings allow much better classification of rare tokens than even the discrete reference model. To understand this



**Figure 2:** Pre-trained embeddings demonstrate strong correlation between the frequency rank of each token and (top) the cosine similarity, and (bottom) the frequency rank of its nearby neighbors. Most rare words are identified with their nearest neighbor, which is also a rare word. Bin size 500; shaded area denotes 50% of values in each bin.



**Figure 3:** Token-level  $F_1$  test score grouped into three bins defined by training set frequency. The x label shows frequency boundaries and token counts per bucket.

effect, we study the geometry of the pre-trained embedding spaces in relation to frequency in Figure 2. The top row shows the relationship between the frequency rank (higher means rarer) and the similarity to its nearest- and fifth-nearest- neighbors. For all three language pairs we observe that most rare words become identical to their nearest neighbor. In contrast, for random embeddings this metric does not depend on rank and is always around 0.4. The bottom row of Figure 3 shows that the nearest neighbors of rare words tend also to be comparably rare. This geometry clarifies in part the surprising performance of random embeddings on rare tokens.

**Combined embeddings.** Our finding motivates combining pre-trained and random embeddings:

$$\mathbf{E}_{\text{cmb}}(y_i) = \frac{\alpha \mathbf{E}_{\text{MT}}(y_i) + (1 - \alpha) \mathbf{E}_{\text{rand}}(y_i)}{\|\alpha \mathbf{E}_{\text{MT}}(y_i) + (1 - \alpha) \mathbf{E}_{\text{rand}}(y_i)\|}.$$

To emphasize pre-trained distances more than the noise, we choose  $\alpha = 0.9$  for all language pairs.

This simple approach leads to overall improved performance, on almost all metrics and language pairs as shown in Table 1. Furthermore, Figure 3 confirms that combined embeddings preserve the performance of pre-trained embeddings on frequent tokens and increase  $F_1$  score on rare tokens. We further study the impact  $\alpha$  on ro-en in Appendix E and observe that for all considered  $\alpha \in [0.5, 0.9]$ , the combination outperforms random and pre-trained embeddings along both metrics; the specific value of  $\alpha$  in this range has only negligible impact.

## 6 Conclusion

Our experimental results show that randomly initialized target embeddings can achieve similar performance as pre-trained ones and even surpass them when a sufficiently large amount of data is available. The gap is most pronounced on very rare tokens. We also found that beam size  $> 1$  does not harm the performance of CoNMT with random target embeddings (compared to pre-trained target embeddings). We suggest combining random and pre-trained embeddings in attempt to maintain high accuracy on frequent tokens as well as rare tokens. This simple approach proved to be effective for en-tr and ro-en in terms of overall performance. However, more refined ways to combine random embeddings with semantically meaningful anchors may lead to more reliable improvements, and ideally hold the potential to remove the reliance on a pretrained model entirely. Finding the best ways to achieve this potential is an important avenue of future work for CoNMT and for continuous modeling of language representations more broadly.

## 278 Limitations

279 **Generalization.** Our experimental results show  
280 that semantic similarity of the targets embeddings  
281 does not play a major role for continuous-output  
282 NMT. However, this not necessarily holds for other  
283 text generation tasks like summarization or lan-  
284 guage modeling. To claim that random target  
285 embeddings can be successfully used for any text  
286 generation task yet has to be proved. In the future,  
287 we will conduct additional experiments on other  
288 text generation tasks, such as summarization and  
289 language modeling.

290 **Dataset Size.** Arora et al. (2020) argue that  
291 random embeddings can achieve comparable per-  
292 formance when the dataset size is big enough. In  
293 our work we report results on three language pairs  
294 with vast range of training samples. The gap be-  
295 tween pre-trained and random embeddings is much  
296 higher for en-tr with 207K training samples than  
297 for ro-en and en-de with 612K and 9.1M training  
298 samples. Moreover, on en-de random embeddings  
299 even outperform pre-trained ones. That hints that  
300 random embeddings indeed work only if there is  
301 sufficiently large amount of data available.

302 **Static Embeddings.** The formulation of the loss  
303 we use in our work, specifically cosine distance,  
304 leads to representation collapse when tuning target  
305 embeddings jointly with the model, That is why in  
306 our work the target embeddings are kept unchanged  
307 during training. Li et al. (2022) show that it  
308 is possible to design a loss that allows for joint  
309 training. However, we believe that fine-tuning of  
310 random embeddings is orthogonal to our study.

311 **Comparison with External Embeddings Mod-**  
312 **els.** In the scope of this work, we compared only  
313 embeddings extracted from the discrete NMT model  
314 (MTtransfer) and randomly generated embeddings.  
315 However, we do not compare random embeddings  
316 with external models like mBart (Liu et al., 2020)  
317 or fasttext (Bojanowski et al., 2016). That is inten-  
318 tional since Tokarchuk and Niculae (2022) showed  
319 that MTtransfer embeddings perform the best com-  
320 pared to the external models, and our goal was to  
321 compare to the best-performing baseline.

322 **Loss Function.** All our results are tied to the  
323 choice of the target objective function, precisely co-  
324 sine similarity. We chose cosine similarity to align  
325 our work with previous studies on CoNMT (Kumar  
326 and Tsvetkov, 2019; Tokarchuk and Niculae, 2022).  
327 We implicitly assumed that our embeddings lie on  
328 the sphere and have the norm equal to 1. In the

future, we would like to experiment with other geo- 329  
metrical spaces and verify if our findings are still 330  
valid. 331

## Risks 332

NMT as a technology is subject to dual-use con- 333  
cerns. We also want to stress that it is possible that 334  
random embedding models make different kinds 335  
of mistakes compared to other models, and they 336  
should be studied and treated with caution before 337  
deployment. CoNMT models are generally at an 338  
earlier stage of development and do not seem likely 339  
to replace the well-studied discrete models in de- 340  
ployed application in the very near future. 341

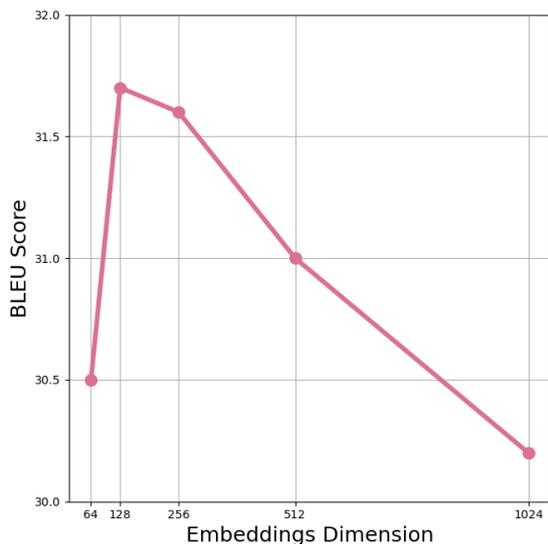
## References 342

- 343 Simran Arora, Avner May, Jian Zhang, and Christo- 343  
pher Ré. 2020. [Contextual embeddings: When are 344](#)  
[they worth it?](#) In *Proceedings of the 58th Annual 345*  
*Meeting of the Association for Computational Lin- 346*  
*guistics*, pages 2650–2663, Online. Association for 347  
Computational Linguistics. 348
- 349 Gayatri Bhat, Sachin Kumar, and Yulia Tsvetkov. 2019. 349  
[A margin-based loss with synthetic negative sam- 350](#)  
[ples for continuous-output machine translation](#). In 351  
*Proceedings of the 3rd Workshop on Neural Gener- 352*  
*ation and Translation*, pages 199–205, Hong Kong. 353  
Association for Computational Linguistics. 354
- 355 Piotr Bojanowski, Edouard Grave, Armand Joulin, 355  
and Tomas Mikolov. 2016. Enriching word vec- 356  
tors with subword information. *arXiv preprint 357*  
*arXiv:1607.04606*. 358
- 359 Fangxiaoyu Feng, Yinfei Yang, Daniel Cer, Naveen Ari- 359  
vazhagan, and Wei Wang. 2022. [Language-agnostic 360](#)  
[BERT sentence embedding](#). In *Proceedings of the 361*  
*60th Annual Meeting of the Association for Compu- 362*  
*tational Linguistics (Volume 1: Long Papers)*, pages 363  
878–891, Dublin, Ireland. Association for Computa- 364  
tional Linguistics. 365
- 366 Taku Kudo. 2018. [Subword regularization: Improv- 366](#)  
[ing neural network translation models with multiple 367](#)  
[subword candidates](#). In *Proceedings of the 56th An- 368*  
*nuual Meeting of the Association for Computational 369*  
*Linguistics (Volume 1: Long Papers)*, pages 66–75, 370  
Melbourne, Australia. Association for Computational 371  
Linguistics. 372
- 373 Taku Kudo and John Richardson. 2018. [SentencePiece: 373](#)  
[A simple and language independent subword tok- 374](#)  
[enizer and detokenizer for neural text processing](#). In 375  
*Proceedings of the 2018 Conference on Empirical 376*  
*Methods in Natural Language Processing: System 377*  
*Demonstrations*, pages 66–71, Brussels, Belgium. 378  
Association for Computational Linguistics. 379

380	Sachin Kumar and Yulia Tsvetkov. 2019. <a href="#">Von misefisher loss for training sequence to sequence models with continuous outputs</a> . In <i>International Conference on Learning Representations</i> .	436
381		437
382		438
383		439
384	Xiang Lisa Li, John Thickstun, Ishaan Gulrajani, Percy Liang, and Tatsunori Hashimoto. 2022. Diffusion-lm improves controllable text generation. <i>ArXiv</i> , abs/2205.14217.	440
385		441
386		442
387		443
388	Yinhan Liu, Jiatao Gu, Naman Goyal, Xian Li, Sergey Edunov, Marjan Ghazvininejad, Mike Lewis, and Luke Zettlemoyer. 2020. <a href="#">Multilingual denoising pre-training for neural machine translation</a> . <i>Transactions of the Association for Computational Linguistics</i> , 8:726–742.	444
389		445
390		446
391		447
392		448
393		449
394	Tomás Mikolov, Kai Chen, Greg Corrado, and Jeffrey Dean. 2013. <a href="#">Efficient estimation of word representations in vector space</a> . In <i>1st International Conference on Learning Representations, ICLR 2013, Scottsdale, Arizona, USA, May 2-4, 2013, Workshop Track Proceedings</i> .	450
395		451
396		452
397		453
398		454
399		455
400	Tomas Mikolov, Edouard Grave, Piotr Bojanowski, Christian Puhersch, and Armand Joulin. 2018. Advances in pre-training distributed word representations. In <i>Proceedings of the International Conference on Language Resources and Evaluation (LREC 2018)</i> .	456
401		457
402		458
403		459
404		460
405	Graham Neubig, Zi-Yi Dou, Junjie Hu, Paul Michel, Danish Pruthi, Xinyi Wang, and John Wieting. 2019. <a href="#">compare-mt: A tool for holistic comparison of language generation systems</a> . <i>CoRR</i> , abs/1903.07926.	461
406		462
407		463
408		464
409	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 <i>Proceedings of NAACL-HLT 2019: Demonstrations</i> .	465
410		466
411		467
412		468
413		469
414	Kishore Papineni, Salim Roukos, Todd Ward, and Wei-Jing Zhu. 2002. <a href="#">Bleu: a method for automatic evaluation of machine translation</a> . In <i>Proceedings of the 40th Annual Meeting of the Association for Computational Linguistics</i> , pages 311–318, Philadelphia, Pennsylvania, USA. Association for Computational Linguistics.	470
415		471
416		472
417		473
418		474
419		475
420		476
421	Matt Post. 2018. <a href="#">A call for clarity in reporting BLEU scores</a> . In <i>Proceedings of the Third Conference on Machine Translation: Research Papers</i> , pages 186–191, Brussels, Belgium. Association for Computational Linguistics.	477
422		478
423		479
424		480
425		481
426	Nils Reimers and Iryna Gurevych. 2019. <a href="#">Sentence-BERT: Sentence embeddings using Siamese BERT-networks</a> . In <i>Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing (EMNLP-IJCNLP)</i> , pages 3982–3992, Hong Kong, China. Association for Computational Linguistics.	482
427		483
428		484
429		485
430		486
431		487
432		488
433		489
434	Evgeniia Tokarchuk and Vlad Niculae. 2022. <a href="#">On target representation in continuous-output neural machine translation</a> . In <i>Proceedings of the 7th Workshop on Representation Learning for NLP</i> , pages 227–235, Dublin, Ireland. Association for Computational Linguistics.	490
435		491
		492
		493
		494
		495
		496
		497
		498
		499
		500
		501
		502
		503
		504
		505
		506
		507
		508
		509
		510
		511
		512
		513
		514
		515
		516
		517
		518
		519
		520
		521
		522
		523
		524
		525
		526
		527
		528
		529
		530
		531
		532
		533
		534
		535
		536
		537
		538
		539
		540
		541
		542
		543
		544
		545
		546
		547
		548
		549
		550
		551
		552
		553
		554
		555
		556
		557
		558
		559
		560
		561
		562
		563
		564
		565
		566
		567
		568
		569
		570
		571
		572
		573
		574
		575
		576
		577
		578
		579
		580
		581
		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
		608
		609
		610
		611
		612
		613
		614
		615
		616
		617
		618
		619
		620
		621
		622
		623
		624
		625
		626
		627
		628
		629
		630
		631
		632
		633
		634
		635
		636
		637
		638
		639
		640
		641
		642
		643
		644
		645
		646
		647
		648
		649
		650
		651
		652
		653
		654
		655
		656
		657
		658
		659
		660
		661
		662
		663
		664
		665
		666
		667
		668
		669
		670
		671
		672
		673
		674
		675
		676
		677
		678
		679
		680
		681
		682
		683
		684
		685
		686
		687
		688
		689
		690
		691
		692
		693
		694
		695
		696
		697
		698
		699
		700
		701
		702
		703
		704
		705
		706
		707
		708
		709
		710
		711
		712
		713
		714
		715
		716
		717
		718
		719
		720
		721
		722
		723
		724
		725
		726
		727
		728
		729
		730
		731
		732
		733
		734
		735
		736
		737
		738
		739
		740
		741
		742
		743
		744
		745
		746
		747
		748
		749
		750
		751
		752
		753
		754
		755
		756
		757
		758
		759
		760
		761
		762
		763
		764
		765
		766
		767
		768
		769
		770
		771
		772
		773
		774
		775
		776
		777
		778
		779
		780
		781
		782
		783
		784
		785
		786
		787
		788
		789
		790
		791
		792
		793
		794
		795
		796
		797
		798
		799
		800

	WMT ro-en			WMT en-tr				WMT en-de			
	train	dev16	test16	train	dev17	test17	test18	train	valid	test16	test18
sentences	612K	2K	2K	207K	1K	3K	3K	9.1M	2.2K	3K	3K
SPM vocabulary (tgt)	27.5K			23.3K				76K			
SPM % oov (tgt)	0.0	0.38	0.31	0.0	0.45	0.53	0.55	0.0	0.0	0.0	0.0

**Table 2:** Datasets Statistics



**Figure 4:** BLEU score of the discrete NMT models on newstest2016 ro-en.

we do all our experiments with dimension equal to 128.

## B.2 Training Parameters

We report fairseq yaml config in Listing 1. Language-pair-specific parameters are highlighted with a comment. Continuous transformer uses base Transformer architecture with 6 layers of encoder and decoder (Vaswani et al., 2017). Total number of training parameters is the following: ro-en discrete is 42M and ro-en continuous 74M; en-tr discrete is 40M and en-tr continuous 73M; en-de discrete is 132M and en-de continuous 123M.

We train our models using shared GPU cluster, which is equipped with GeForce GTX TITAN X as well as NVIDIA A100.

## C Word-level Embeddings

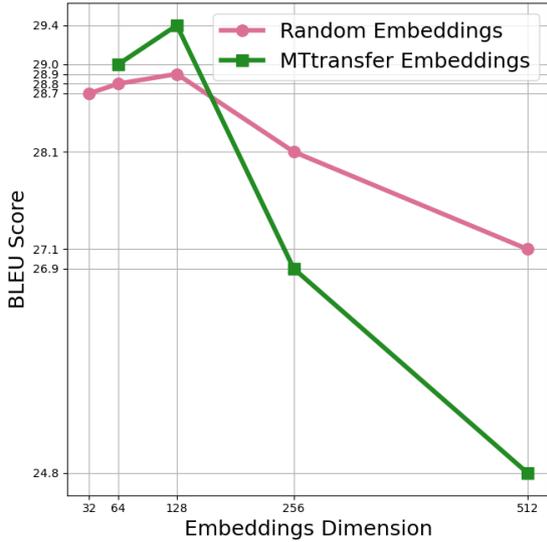
Since the continuous-output model struggles with subwords continuation and, at the same time, performs better on rare words, we conduct experiments on the word level. Word-level model tends to suffer from out-of-vocabulary issues (Table 2), so

**Listing 1** Training yaml config for CoNMT.

```

task:
  _name: translation
  data: language_specific_data
criterion:
  _name: cosine_ar_criterion
model:
  _name: continuous_transformer
  decoder:
    output_dim: 128
    learned_pos: true
  encoder:
    learned_pos: true
  dropout:
    0.3 # ro-en and en-tr
    0.1 # en-de
  target_embed_path: path_to_static_embeddings
  no_decoder_final_norm: false
optimizer:
  _name: adam
  adam_betas: (0.9,0.98)
lr_scheduler:
  _name: inverse_sqrt
warmup_updates:
  10000 # ro-en and en-de
  4000 # en-tr
warmup_init_lr: 1e-07
dataset:
  validate_after_updates: 10000
  max_tokens: 4096
  validate_interval_updates: 2000
optimization:
  lr: [0.0005]
  update_freq: [16]
  max_update: 50000
  stop_min_lr: 1e-09
checkpoint:
  no_epoch_checkpoints: true
  best_checkpoint_metric: bleu
  maximize_best_checkpoint_metric: true

```



**Figure 5:** BLEU score on ro-en newstest2016 of continuous-output model with various dimensionalities of random and pre-trained (MTtransfer) target embeddings.

discrete model performance drops respectively. Table 3 provides the comparison between the discrete word-level model and continuous-output model with random targets. Even though the continuous-output model struggles with subwords continuations, overall, using subwords allows us to have a stronger model both for discrete and continuous-output cases.

model	ro-en	en-tr
discrete words	28.5	8.9
continuous random words	27.6	5.6
discrete tokens	32.1	12.7
continuous random tokens	29.2	9.3

**Table 3:** BLEU scores for word level and tokens level models on validation set with greedy decoding.

### C.1 Subword Embeddings

We rely on the unigram language model for subword segmentation (Kudo, 2018) to train discrete and continuous-output NMT models as mentioned in Section §5. We hypothesize that it is harder for the continuous-output model to predict subwords than for the discrete model. Table 4 illustrates that the f1 macro average for the beginning of the spm tokens and continuation of the spm tokens differ a lot for discrete and continuous models. While the discrete model performs better on continuations, continuous models struggle with continuations of

subwords. However, overall scores for pre-trained and random targets are the same for continuation and random embeddings performs slightly better on the beginning of the subwords.

model	F1	
	SPM start	SPM cont.
discrete	0.12	0.14
pre-trained embeddings	0.10	0.09
random embeddings	0.11	0.09

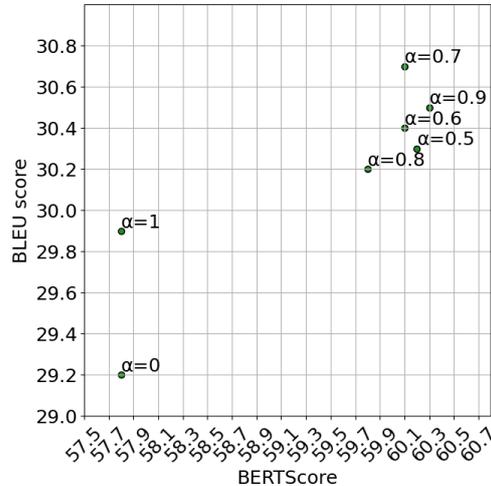
**Table 4:** F1 score on newstest2016 ro-en for beginning and continuation of the SentencePiece tokens.

## D Beam Search

In our work, we use implementation of the beam search provided by fairseq. However, instead of using log probabilities of the next token, we rely on the cosine similarity scores between output vector and all tokens in the vocabulary. We restrict maximum length of generated sentence to be length of the source sentence plus 200. For CoNMT, beam search may have a probabilistic interpretation by noticing that the cosine loss is equivalent to a Langevin (also known as vMF) log-likelihood with constant concentration parameter  $\kappa$ : in beam search we use this probabilistic interpretation and take

$$\log p(y_i = t \mid \mathbf{y}_{<i}, \mathbf{x}) = -\cos(\mathbf{E}(t), \mathbf{h}) + \log C_d(1),$$

*i.e.*, we apply the normalizing constant of the Langevin distribution for dimension  $d$  and fixed concentration  $\kappa = 1$ .



**Figure 6:** BLEU and BERTScores on ro-en newsdev2016 with different values of  $\alpha$ .

## 549 **E Combined Embeddings**

550 In Table 1 we report performance of combined  
551 embeddings with  $\alpha = 0.9$ . To study the effect of  
552  $\alpha$  on the models' performance, we conduct exper-  
553 iments on ro-en for  $\alpha \in [0.5, 0.9]$ . As shown  
554 in Figure 6, for all cases combined embeddings  
555 outperform pre-trained and random ones on both  
556 metrics.