

CycleGN: a Cycle Consistent approach for Neural Machine Translation training using the Transformer model in a shuffled dataset

Anonymous ACL submission

Abstract

CycleGN is a Transformer architecture using a Discriminator-less CycleGAN approach, specifically tailored for training Machine Translation models utilizing non-parallel datasets. Despite the widespread availability of large parallel corpora for numerous language pairs, the capacity to employ solely monolingual datasets would substantially expand the pool of training data. This approach is particularly beneficial for languages with scarce parallel text corpora.

The foundational concept of our research posits that in an ideal scenario, translations of translations should revert to the original source sentences. Consequently, we can simultaneously train a pair of models using a Cycle Consistency Loss framework. This method bears resemblance to the technique of back-translation, prevalently employed in Machine Translation, where a pre-trained translation model is used to generate new examples from a monolingual corpus, thereby artificially creating a parallel dataset for further training and refinement.

1 Introduction

The introduction of the Transformer architecture (Vaswani et al., 2017) marked a significant advancement in the field of Machine Translation, witnessing widespread adoption since its inception. Although self-attention mechanisms were not novel and had been investigated in prior studies (Bahdanau et al., 2016), the Transformer model demonstrated its formidable capabilities within Natural Language Processing (NLP). Characterized by its parallelized structure, the Transformer architecture facilitated computational efficiency, enabling the incorporation of a larger number of parameters. This enhancement has been exemplified in NLP systems like Charles University Block-Backtranslation-Improved Transformer Translation (cubitt) (Popel et al., 2020), which have surpassed the performance levels of human professionals in certain contexts.

Neural Machine Translation (NMT) datasets necessitate substantial text corpora, structured as aligned pairs. This alignment implies the requirement for sentences with equivalent meaning to be present in a minimum of two distinct languages, enabling the initiation of model training to forge linguistic linkages. Ongoing initiatives, including OPUS (Tiedemann and Thottingal, 2020) and Tatoeba (Tiedemann, 2012), are committed to facilitating public access to these datasets. Clearly, parallel datasets comprise a small subset of the volume of data in monolingual datasets.

Despite the remarkable efficacy exhibited by Large Language Models (LLMs) in MT (Machine Translation) without the necessity of exclusive training on parallel data (Zhu et al., 2023), their considerable magnitude renders them costly in terms of both training and operation. This economic burden consequently restricts their widespread availability.

Back-translation (Sennrich et al., 2016) is a technique leveraging a trained MT model to translate sentences from a monolingual dataset to produce corresponding pairs, thereby synthetically augmenting the training data. Our research is founded on the premise that the process of translating a sentence from a source language to a target language, followed by its retranslation from the target language back to the source language, allows for the measurement of the disparity between the original and the machine-translated sentences. This disparity serves as a metric to assess the efficacy of the models and facilitates the backpropagation of gradients within the networks. Notably, this methodology has been previously implemented in the realm of Image-to-Image Translation, as evidenced in the renowned CycleGAN study from Zhu et al. (2017).

2 Previous work

The TextCycleGAN model (Lorandi et al., 2023), while not utilizing the Transformer architecture nor

042
043
044
045
046
047
048
049
050
051
052
053
054
055
056
057
058
059
060
061
062
063
064
065
066
067
068
069
070
071
072
073
074
075
076
077
078
079
080
081

operating within the MT field, introduced an innovative strategy for text style transfer. This approach employed a CycleGAN on the Yelp dataset to facilitate the learning of mappings between positive and negative textual styles, notably in the absence of paired examples.

Shen et al. (2017) exemplified the feasibility of training two encoder-decoder networks in an unsupervised manner that enables the sharing of a latent space, thereby permitting style transfer. Lample et al. (2018), adopting a similar technique within the MT context, substantiated that the use of parallel datasets is not a prerequisite for effective translation.

3 Dataset

In the context of the current study, a “shuffled” dataset is defined as a parallel dataset wherein the sentences of one language have been systematically rearranged. Consequently, this results in a non-parallel corpus where it is guaranteed that each sentence has a corresponding translation located at an unspecified index within the dataset. The authors postulate that when employing sufficiently large monolingual datasets, which are not derived from shuffled parallel corpora, it is likely that most sentences will possess an accurate translation “somewhere” within the dataset.

For the purposes of this research, a shuffled dataset was utilized in lieu of a monolingual dataset. This choice was made to facilitate a direct comparison of our approach with conventional NMT training, employing an identical non-shuffled parallel dataset and the same model architecture.

The dataset employed in this study is the English-German language pair from the WMT23 challenge (Kocmi et al., 2023). Specifically, only the first half of this dataset was used for training, due to the current implementation’s high computational demands. This amounts to a total of approximately 27 million sentences. The data released for the WMT23 General MT task can be freely used for research purposes.

4 Training

For greater clarity, the mathematical notations from the original CycleGAN work will be employed in the present study. Given two languages \mathcal{X} and \mathcal{Y} with appropriate datasets, our objective is to obtain two NMT models $\mathcal{G} : \mathcal{X} \mapsto \mathcal{Y}$ and $\mathcal{F} : \mathcal{Y} \mapsto \mathcal{X}$ such that for $x \in \mathcal{X}$, $\mathcal{G}(x) = \hat{y}$, for $\hat{y} \in \mathcal{Y}$ and that

for $y \in \mathcal{Y}$, $\mathcal{F}(y) = \hat{x}$, for $\hat{x} \in \mathcal{X}$. If the translations are perfect, $\mathcal{G}(\mathcal{F}(y)) = y$ and $\mathcal{F}(\mathcal{G}(x)) = x$. By using the Cross-Entropy Loss (CEL) (Zhang and Sabuncu, 2018) in the role of the Cycle Consistency Loss (CCL), we can determine the distance between the original sentence and its double translation in order to compute the gradients.

As in the original CycleGAN work, our current study also implements an Identity Loss (IL), which relies on the CEL, to help with the training stability. As \mathcal{G} consists in a mapping $\mathcal{X} \mapsto \mathcal{Y}$, if given an input $y \in \mathcal{Y}$, we want to obtain an unchanged output such that $\mathcal{G}(y) = y$. The same is applied to \mathcal{F} , where we also compute the IL between $\mathcal{F}(x)$ and x . See Figure 1.

4.1 Obtaining labels

In the training process of a Transformer model, it is imperative to have prior knowledge of the labels, as the decoder predicts tokens sequentially. Each token prediction, barring the initial one, is contingent upon all preceding predictions. The act of selecting the most probable token constitutes a non-differentiable operation, thus precluding the possibility of backpropagation. By possessing prior knowledge of the reference translation, it becomes feasible to contrast each predicted token against the ground truth, enabling the calculation of loss at every step.

Teacher Forcing (Gers et al., 2002) is a technique that involves substituting the predicted token with the actual ground truth at each stage of the decoding process. This approach is designed to mitigate the cascading impact of early erroneous predictions in the sequence.

The CycleGN training process used here consists in a cooperation between \mathcal{G} and \mathcal{F} . The first step is to generate \hat{x} and \hat{y} , since labels are not required during inference, as backpropagation is unnecessary. Even though this step cannot be used to compute the gradients, it is crucial for the entire process. From $\mathcal{G}(\mathcal{F}(y)) = y$ and $\mathcal{F}(\mathcal{G}(x)) = x$, it follows that the label for \hat{y} is x and the label for \hat{x} is y . We can compute \hat{x} from $\mathcal{F}(\hat{y})$ with x as the label, and \hat{y} from $\mathcal{G}(\hat{x})$ with y as the label, and use the CCL between \hat{x} and x , and between \hat{y} and y to compute the gradients and backpropagate.

4.2 A Discriminator-less GAN

The CycleGAN methodology, as indicated by its nomenclature, is predicated on the Generative Adversarial Network (GAN) framework, initially

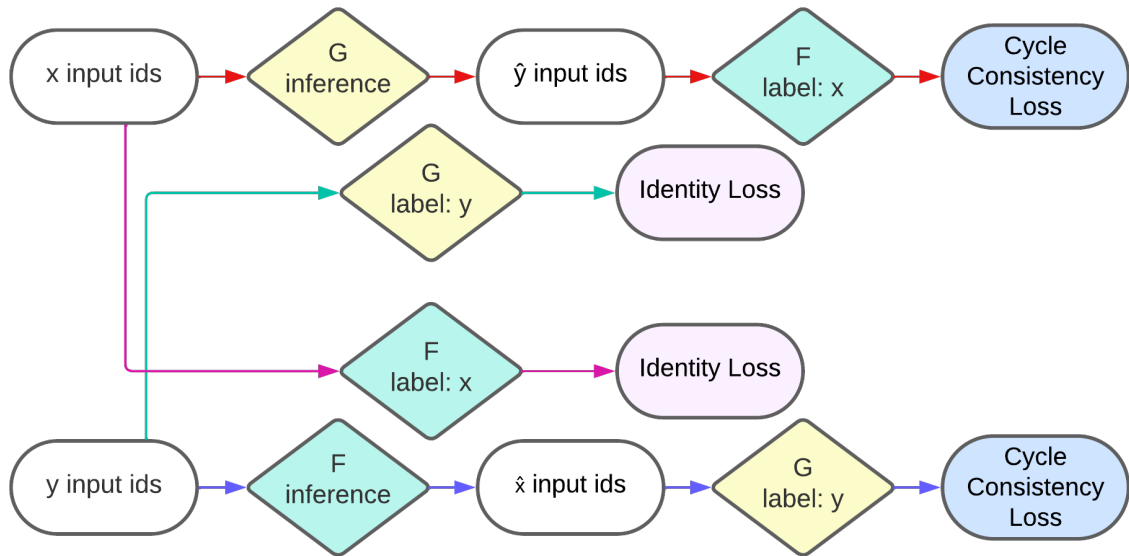


Figure 1: CycleGN training process.

introduced in Goodfellow et al. (2014). This paradigm involves the training of a Generator model in conjunction with another model, termed the Discriminator. The Discriminator is specifically trained to distinguish between authentic samples drawn from the dataset and synthetic samples produced by the Generator. In the CycleGAN training process, the Discriminators intervene after the generation of \hat{x} and \hat{y} , helping the training of the Generators. However, as mentioned in Section 4.1, there can be no gradient computation during the generation of \hat{x} and \hat{y} in a transformer model and as such, Discriminators cannot be used in the present work. This is why CycleGN is not an “Adversarial” approach, hence the name.

5 Model architecture

The architecture used for both models, \mathcal{G} and \mathcal{F} , is the Marian framework (Junczys-Dowmunt et al., 2018) implemented by Huggingface’s Transformers library (Wolf et al., 2020), which is licenced under the Apache Licence. While most parameters follow the default configuration, Table 1 references the changes that were made in order to reduce the computational cost of the architecture.

6 Vocabulary organization

Sequence2Sequence models employ either a unified tokenizer or two distinct tokenizers. In the case of a single tokenizer, it is trained using sentences

Parameter	Huggingface	Current work
Vocabulary size	58,101	32,000
Encoder layers	12	6
Decoder layers	12	6
Encoder attention heads	16	8
Decoder attention heads	16	8
Encoder feed-forward	4096	2048
Decoder feed-forward	4096	2048
Position embeddings	1024	128
Activation function	GELU	ReLU

Table 1: Non-default parameters in the configuration of Marian Transformer models

from both the source and target distributions, avoiding any duplicates. This approach facilitates the sharing of the encoder and decoder embedding layers, thereby diminishing computational demands and enhancing model accuracy (Press and Wolf, 2017).

Conversely, the alternative approach entails training one tokenizer on the source distribution and another one on the target distribution. While this method restricts the possibility of tying embeddings, it can potentially double the vocabulary size. The overall vocabulary size of the model in this scenario, is the cumulative total of the two individual vocabularies, barring shared tokens like punctuation symbols.

While contemporary Transformer models like Bidirectional Encoder Representations from Transformers (BERT) (Devlin et al., 2019) and Generative Pre-trained Transformers (GPT) (Radford et al.,

2018) typically utilize a single tokenizer, this study introduces a novel vocabulary methodology that amalgamates the aforementioned approaches. This method involves training two tokenizers, each for a respective language and with half the vocabulary size. Subsequently, the identifiers of one tokenizer are adjusted to prevent overlap, yielding a result analogous to a single tokenizer that includes duplicates across languages. It is important to note that special tokens such as $\langle eos \rangle$ (End of Sentence) and $\langle pad \rangle$ (Padding) are shared and not duplicated. This strategy is designed to simplify model analysis during development, albeit at the expense of a reduced vocabulary.

7 Pretraining

The CycleGN approach requires a pre-training step, as it will not converge at all without it. Indeed, as there is no Discriminator to ensure that \hat{x} belongs to \mathcal{X} and \hat{y} belongs to \mathcal{Y} , \mathcal{G} and \mathcal{F} can converge towards identity matrices. That is, if both \mathcal{G} and \mathcal{F} do not apply any change to their input, they can still achieve $\mathcal{G}(\mathcal{F}(y)) = y$ and $\mathcal{F}(\mathcal{G}(x)) = x$ without learning how to translate.

Masked Language Modeling (MLM) is a pre-training strategy first implemented in BERT, wherein a specified proportion of tokens within the input text are substituted with a unique $\langle mask \rangle$ token. The objective of the neural network under this paradigm is to accurately reconstruct the original sentence. This process enables the model to discern intricate relationships between words and to develop a profound representation of the language. This pre-training has revealed excellent performances in diverse NLP application such as sentiment analysis (Alaparthi and Mishra, 2021), text classification (Sun et al., 2020), Named Entity Recognition (NER) (Souza et al., 2020) (Chang et al., 2021) (Akhtyamova, 2020) and paraphrase detection (Khairova et al., 2022).

As MLM does not require any labeling, it is perfectly adapted to the CycleGN approach. A single model \mathcal{H} is trained on the entire dataset for a single epoch to reconstruct both languages, with 15% of the input tokens masked. When training the CycleGN, rather than randomly initializing \mathcal{G} and \mathcal{F} , the parameters from \mathcal{H} are directly copied to \mathcal{G} and \mathcal{F} . Indeed, as \mathcal{H} learns to reconstruct both language \mathcal{X} and \mathcal{Y} , it can be used to initialize both networks. Figure 2 shows the training process of \mathcal{H} .

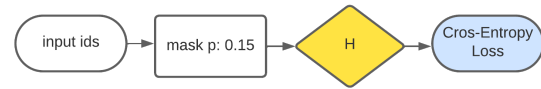


Figure 2: Masked Language Modeling training process.

8 Batch size

The original CycleGAN research mentions using a batch size of one, and while they did not state the reason in the research paper, one of the authors explained it in a GitHub issue (Junyanz, 2017) as a lack of GPU memory.

Rajput et al. (2021) examined the impact of batch size within the CycleGAN architecture, observing a significant decline in performance with its increase. This deterioration was evident both through the example images presented in that study and through the calculated cosine dissimilarity, indicating inferior model performance with larger batch sizes. However, quality was achieved at the expense of computational efficiency, as the training duration to achieve 200 epochs was 8 hours with a batch size of 1, but this was reduced to just 2 hours with a batch size of 64.

In the context of our research, however, the trade-off between quality improvement and computing resource, as observed in the aforementioned study, does not hold true. Utilizing a batch size of 1 in our experiments hindered any form of convergence. Consequently, a batch size of 16 was selected, as it represented the maximum capacity that could be accommodated within the available 24GB of GPU memory of the Nvidia 4090 used for this work.

9 Training stability

It is crucial for a CycleGAN architecture that the two models exhibit approximately equivalent levels of performance. Given the interdependent nature of these models, where the output of one serves as the input for the other, maintaining consistency between them during training is imperative. Without a strategy in place to prevent the performance of the models from diverging, it is possible for one model to gain the “upper hand” over the other.

9.1 Divergence between the Generators

Figure 3 presents the evolution of the CCL of an early prototype of CycleGN and it can clearly be seen that one of the two generators, \mathcal{F} , ends up performing much better than its counterpart \mathcal{G} , which

320 blocks any future training.

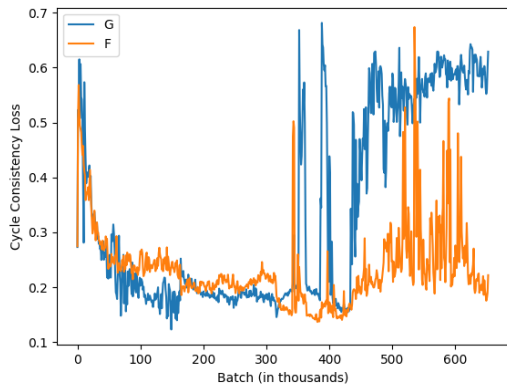


Figure 3: Evolution of the Cross-Entropy Loss during the training of an early prototype.

321 9.2 Gradient Clipping

322 Gradient clipping is a technique utilized in the training of Deep Learning (DL) models, to address the
323 problem of ‘exploding’ gradients. This issue occurs when gradients escalate to excessively high
324 values during training, leading to numerical instability and impeding the model’s convergence to an
325 optimal solution.

326 Gradient clipping can be implemented through two primary methods: norm clipping and value
327 clipping. Norm clipping involves establishing a threshold on the overall magnitude of the gradient
328 vector. On the other hand, value clipping involves individually adjusting elements of the gradient vector
329 that exceed the specified threshold.

330 By clipping the gradients by norm, with a threshold of 1.0, as advised by the Huggingface library,
331 the training stabilized and the divergence between \mathcal{G} and \mathcal{F} was observed to disappear.

340 10 One large epoch or multiple smaller ones?

341 The CycleGAN framework is recognized for its computational intensity due to several inherent factors.
342 Primarily, as CycleGAN operates on the principle of cycle consistency, it necessitates the training
343 of two GANs simultaneously – one for each direction of the transformation. This structure requires
344 substantial computational resources, as each GAN consists of both a Generator and a Discriminator.
345

346 The resource-intensiveness of the CycleGAN process, thus limits the size of the dataset that can

353 be used in a reasonable time. This necessitated a decision between training for a single epoch on a
354 large dataset, or training for multiple epochs on a smaller corpus arose.

355 We compared the CycleGAN model on the entire dataset under four different conditions:
356

- 357 1. One epoch containing 1% of the dataset 359
- 360 2. Five epochs containing 0.2% of the dataset 360
- 361 3. One epoch containing 2% of the dataset 361
- 362 4. Five epochs containing 0.4% of the dataset 362

363 We have selected the Crosslingual Optimised Metric for Evaluation of Translation (COMET)
364 score, as proposed by [Rei et al. \(2020\)](#), as our comparison criterion. This metric has proven to be one
365 of the most effective in recent WMT competitions, according to [Kocmi et al. \(2022\)](#), due to its strong
366 correlation with human judgment, aligning well with our goal of mirroring human evaluative standards.
367 Multiple COMET models have been made available and we chose the default “wmt22-comet-da” model.
368 The average scores obtained on 10,000 sentences that were not part of the model training set are presented in
369 Table 2.

Condition	English->German	German->English
1	0.2727	0.2715
2	0.2411	0.2635
3	0.2741	0.2665
4	0.2258	0.2658

Table 2: COMET scores of CycleGAN models depending on the dataset condition.

370 Models exposed to a larger portion of the total dataset demonstrate superior performance compared
371 to those limited to a smaller, repetitive subset, especially when the dataset encompasses over half
372 a million to a million sentences. We extrapolate this result to larger datasets and thus chose to train
373 our model for a single epoch on the largest dataset possible.

384 11 Results

385 To measure the performances of CycleGAN, every 1000th batch the CCL was averaged and 1,000 sentences
386 from the test set were translated to compute the COMET score.

387 Figure 4 demonstrates how the addition of gradient clipping helps with training stability.
388

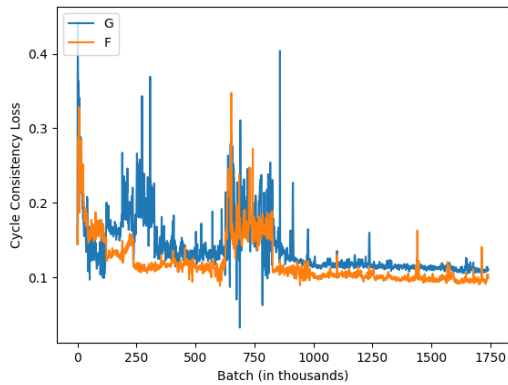


Figure 4: Evolution of the Cross-Entropy Loss during the training.

11.1 Translation quality

Even if tracking the CCL is an inexpensive manner to estimate the progress of the training of the CycleGN architecture, as mentioned in Section 7, it can also hide an absence of translation. Figure 5 demonstrates that the actual quality of translation, as measured by the COMET metric, increases with time. Note that the sudden drop is discussed in the next section.

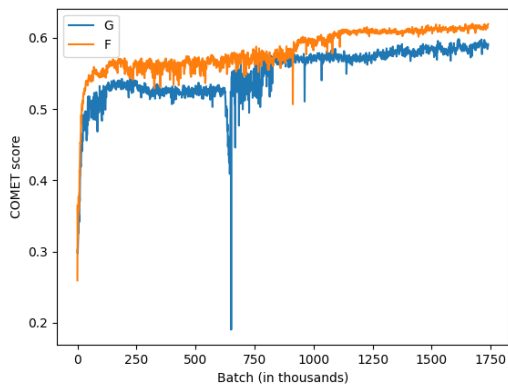


Figure 5: Evolution of the COMET score during the CycleGN training.

After the end of the training, a test set of 10,000 sentences per language were translated and the COMET scores are displayed in Table 3.

	English->German	German->English
Score	0.505	0.537

Table 3: COMET score of CycleGN models.

As mentioned in Section 3, in order to give a point of comparison, we trained a couple of

architecture-matched models using the parallel dataset. As in the case of the CycleGN training, these models were only trained for a single epoch on the first half of the WMT23 English-German language pair. Results are displayed in Table 4. We fully expected the COMET score of the CycleGN to be inferior to architectures using parallel corpora, but we believe the differences between the scores will reduce with larger datasets.

	English->German	German->English
Score	0.780	0.775

Table 4: COMET score of architecture-matched models.

11.2 The sudden drop

Upon examining Figure 5, there is an observable precipitous deterioration in the CEL of Generator \mathcal{G} post the 600,000th batch mark. Delving into the test set translations conducted at every 1,000th batch interval reveals substantial and abrupt modifications. Appendix A presents the evolution of the first three translations of the test set.

While these alterations, despite their detrimental effect on the translation’s quality, ostensibly do not exert a significant influence on the aggregate translation score at first, they are impressively accurate in predicting the drop in quality that ensues.

Examining the progression of alterations without delving into the translation quality, one can discern a clear pattern. Initially, an inverted comma is introduced at the onset of each sentence, subsequently appearing at the termination of most sentences as well. This is then substituted with a “(3)” at the start of each sentence, eventually being replaced by a letter “(b)”. This phase, primarily characterised by superficial quality degradation, gives way to a more pronounced collapse. Here, a significant portion of sentences is rendered as a parenthesis followed by a repeated letter “k”.

11.3 Recovery

Remarkably, this phase of decline vanishes in the subsequent batch, resulting in a minor, primarily cosmetic alteration in the output. This demonstrates that the training process is robust and can withstand even major disturbances to one of the two generators. This also shows the importance of accurately monitoring the accuracy achieved, to avoid stopping the training during such a drop.

448	12 Future Work		494
449	12.1 Activation function		495
450	The activation function in machine learning, especially in neural networks, plays a crucial role in determining the output of a node or neuron. It is a mathematical function that introduces non-linearity into the network, enabling it to learn and perform more complex tasks that linear functions cannot handle. The current CycleGN implementation relies on ReLU, but it seems GELU has now become the default activation function in Huggingface.		496
451			497
452			498
453			499
454			500
455			501
456			502
457			503
458			504
459	12.2 Longer dataset		505
460	Our current work has been trained on a small dataset compared to MT standards. Future work should try to see how convergence progresses with more iterations. Further computational optimizations are probably necessary to shorten the training time required.		506
461			507
462			508
463			509
464			510
465			511
466	12.3 Larger models		512
467	The current architecture relies on a total of 158,769,152 parameters, which is only about a third of the size of the default in the Huggingface library. Although Table 4 demonstrates that the current number of parameters is capable of producing better translations and an increase in both the number of epochs and the size of the dataset should be prioritized, larger models are common in NMT.		513
468			514
469			515
470			516
471			517
472			518
473			519
474			520
475	13 Source Code		521
476	The source code of CycleGN is available at [anonymized].		522
477			523
478			524
479			525
480			526
481			527
482			528
483			529
484			530
485			531
486			532
487			533
488			534
489			535
490			536
491			537
492			538
493			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
			801
			802
			803
			804
			805
			806
			807
			808
			809
			810
			811
			812
			813
			814
			815
			816
			817
			818
			819
			820
			821
			822
			823
			824
			825
			826
			827
			828
			829
			830
			831
			832
			833
			834
			835
			836
			837
			838
			839
			840
			841
			842
			843
			844
			845
			846
			847
			848
			849
			850
			851
			852
			853
			854
			855
			856
			857
			858
			859
			860
			861
			862
			863
			864
			865
			866
			867
			868
			869
			870
			871
			872
			873
			874
			875
			876
		</	

543	In our commitment to scientific integrity, we	Junyanz. 2017. Question: Batch size · issue 27 · junyanz/pytorch-cyclegan-and-pix2pix .	594
544	maintain transparency in our research methodolo-		595
545	gies, model development, and findings. We aim to		
546	make our results reproducible and accessible to the	Nina Khairova, Anastasiia Shapovalova, Orken Mamyrbayev, Nataliia Sharonova, and Kuralay. 2022. Using bert model to identify sentences paraphrase in the news corpus. In <i>CEUR Workshop Proceedings, volume 3171</i> , pages 38–48.	596
547	scientific community, contributing positively to the		597
548	field of machine translation.		598
			599
			600
549	Acknowledgements	Tom Kocmi, Eleftherios Avramidis, Rachel Bawden, Ondřej Bojar, Anton Dvorkovich, Christian Federmann, Mark Fishel, Markus Freitag, Thamme Gowda, Roman Grundkiewicz, Barry Haddow, Philipp Koehn, Benjamin Marie, Christof Monz, Makoto Morishita, Kenton Murray, Makoto Nagata, Toshiaki Nakazawa, Martin Popel, Maja Popović, and Mariya Shmatova. 2023. Findings of the 2023 conference on machine translation (WMT23): LLMs are here but not quite there yet . In <i>Proceedings of the Eighth Conference on Machine Translation</i> , pages 1–42, Singapore. Association for Computational Linguistics.	601
550	This publication has emanated from research con-		602
551	ducted with the financial support of Science Founda-		603
552	tion Ireland under Grant number 18/CRT/6183.		604
553	For the purpose of Open Access, the author has		605
554	applied a CC BY public copyright licence to any		606
555	Author Accepted Manuscript version arising from		607
556	this submission.		608
			609
			610
			611
			612
			613
557	References	Tom Kocmi, Rachel Bawden, Ondřej Bojar, Anton Dvorkovich, Christian Federmann, Mark Fishel, Thamme Gowda, Yvette Graham, Roman Grundkiewicz, Barry Haddow, Rebecca Knowles, Philipp Koehn, Christof Monz, Makoto Morishita, Masaaki Nagata, Toshiaki Nakazawa, Michal Novák, Martin Popel, and Maja Popović. 2022. Findings of the 2022 conference on machine translation (WMT22) . In <i>Proceedings of the Seventh Conference on Machine Translation (WMT)</i> , pages 1–45, Abu Dhabi, United Arab Emirates (Hybrid). Association for Computational Linguistics.	614
558	Liliya Akhtyamova. 2020. Named entity recognition in spanish biomedical literature: Short review and bert model . In <i>2020 26th Conference of Open Innovations Association (FRUCT)</i> , pages 1–7.		615
559			616
560			617
561			618
562	Shivaji Alaparathi and Manit Mishra. 2021. Bert: a sentiment analysis odyssey . <i>Journal of Marketing Analytics</i> , 9(2):118–126.		619
563			620
564			621
565	Dzmitry Bahdanau, Kyunghyun Cho, and Yoshua Ben-		622
566	gio. 2016. Neural machine translation by jointly learning to align and translate .		623
567			624
			625
568	Yuan Chang, Lei Kong, Kejia Jia, and Qinglei Meng.		626
569	2021. Chinese named entity recognition method based on bert . In <i>2021 IEEE International Conference on Data Science and Computer Application (ICDSCA)</i> , pages 294–299.		627
570			628
571			
572			
573	Jacob Devlin, Ming-Wei Chang, Kenton Lee, and	Guillaume Lample, Alexis Conneau, Ludovic Denoyer, and Marc’Aurelio Ranzato. 2018. Unsupervised machine translation using monolingual corpora only .	629
574	Kristina Toutanova. 2019. Bert: Pre-training of deep bidirectional transformers for language understanding .		630
575			631
576			632
577	Felix Gers, Nicol Schraudolph, and Jürgen Schmidhuber.	Michela Lorandi, Maram A.Mohamed, and Kevin McGuinness. 2023. Adapting the CycleGAN Architecture for Text Style Transfer . <i>Irish Machine Vision and Image Processing Conference</i> .	633
578	2002. Learning precise timing with lstm recurrent networks . <i>Journal of Machine Learning Research</i> , 3:115–143.		634
579			635
580			636
581	Ian J. Goodfellow, Jean Pouget-Abadie, Mehdi Mirza,	Martin Popel, Marketa Tomkova, Jakub Tomek, Łukasz Kaiser, Jakob Uszkoreit, Ondřej Bojar, and Zdeněk Žabokrtský. 2020. Transforming machine translation: a deep learning system reaches news translation quality comparable to human professionals . <i>Nature Communications</i> , 11(4381):1–15.	637
582	Bing Xu, David Warde-Farley, Sherjil Ozair, Aaron		638
583	Courville, and Yoshua Bengio. 2014. Generative adversarial networks .		639
584			640
585	Marcin Junczys-Dowmunt, Roman Grundkiewicz,	Alec Radford, Karthik Narasimhan, Tim Salimans, Ilya Sutskever, et al. 2018. Improving language understanding by generative pre-training .	641
586	Tomasz Dwojak, Hieu Hoang, Kenneth Heafield,		642
587	Tom Neckermann, Frank Seide, Ulrich Germann,		643
588	Alham Fikri Aji, Nikolay Bogoychev, André F. T.		
589	Martins, and Alexandra Birch. 2018. Marian: Fast neural machine translation in C++ . In <i>Proceedings of ACL 2018, System Demonstrations</i> , pages 116–121, Melbourne, Australia. Association for Computational Linguistics.		644
590			645
591			646
592			
593			
		Pranjali Singh Rajput, Kanya Satis, Sonnya Dellarosa, Wenxuan Huang, and Obinna Agba. 2021. cgans for cartoon to real-life images .	647
			648
		Ricardo Rei, Craig Stewart, Ana C Farinha, and Alon Lavi. 2020. COMET: A neural framework for MT	

649 evaluation. In *Proceedings of the 2020 Conference*
650 *on Empirical Methods in Natural Language Process-*
651 *ing (EMNLP)*, pages 2685–2702, Online. Association
652 for Computational Linguistics.

653 Rico Sennrich, Barry Haddow, and Alexandra Birch.
654 2016. [Improving neural machine translation models](#)
655 [with monolingual data.](#)

656 Tianxiao Shen, Tao Lei, Regina Barzilay, and Tommi
657 Jaakkola. 2017. [Style transfer from non-parallel text](#)
658 [by cross-alignment.](#)

659 Fábio Souza, Rodrigo Nogueira, and Roberto Lotufo.
660 2020. [Portuguese named entity recognition using](#)
661 [bert-crf.](#)

662 Chi Sun, Xipeng Qiu, Yige Xu, and Xuanjing Huang.
663 2020. [How to fine-tune bert for text classification?](#)

664 Jörg Tiedemann and Santhosh Thottingal. 2020. OPUS-
665 MT — Building open translation services for the
666 World. In *Proceedings of the 22nd Annual Confer-*
667 *ence of the European Association for Machine Trans-*
668 *lation (EAMT)*, Lisbon, Portugal.

669 Jörg Tiedemann. 2012. Parallel data, tools and inter-
670 faces in opus. In *Proceedings of the Eight Inter-*
671 *national Conference on Language Resources and*
672 *Evaluation (LREC'12)*, Istanbul, Turkey. European
673 Language Resources Association (ELRA).

674 Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob
675 Uszkoreit, Llion Jones, Aidan N. Gomez, Lukasz
676 Kaiser, and Illia Polosukhin. 2017. [Attention is all](#)
677 [you need.](#)

678 Thomas Wolf, Lysandre Debut, Victor Sanh, Julien
679 Chaumond, Clement Delangue, Anthony Moi, Pier-
680 ric Cistac, Tim Rault, Rémi Louf, Morgan Funtow-
681 icz, Joe Davison, Sam Shleifer, Patrick von Platen,
682 Clara Ma, Yacine Jernite, Julien Plu, Canwen Xu,
683 Teven Le Scao, Sylvain Gugger, Mariama Drame,
684 Quentin Lhoest, and Alexander M. Rush. 2020. [Hug-](#)
685 [gingface’s transformers: State-of-the-art natural lan-](#)
686 [guage processing.](#)

687 Zhilu Zhang and Mert R. Sabuncu. 2018. [Generalized](#)
688 [cross entropy loss for training deep neural networks](#)
689 [with noisy labels.](#)

690 Jun-Yan Zhu, Taesung Park, Phillip Isola, and Alexei A.
691 Efros. 2017. [Unpaired image-to-image translation](#)
692 [using cycle-consistent adversarial networks.](#)

693 Wenhao Zhu, Hongyi Liu, Qingxiu Dong, Jingjing Xu,
694 Shujian Huang, Lingpeng Kong, Jiajun Chen, and Lei
695 Li. 2023. [Multilingual machine translation with large](#)
696 [language models: Empirical results and analysis.](#)

