

# Generative Pretraining for Paraphrase Evaluation

Anonymous ACL submission

## Abstract

We introduce ParaBLEU, a paraphrase representation learning model and evaluation metric for text generation. Unlike previous approaches, ParaBLEU learns to understand paraphrasis using generative conditioning as a pre-training objective. ParaBLEU correlates more strongly with human judgements than existing metrics, obtaining new state-of-the-art results on the 2017 WMT Metrics Shared Task. We show that our model is robust to data scarcity, exceeding previous state-of-the-art performance using only 50% of the available training data and surpassing BLEU, ROUGE and METEOR with only 40 labelled examples. Finally, we demonstrate that ParaBLEU can be used to conditionally generate novel paraphrases from a single demonstration, which we use to confirm our hypothesis that it learns abstract, generalized paraphrase representations.

## 1 Introduction

Representing the relationship between two pieces of text, be it through a simple algorithm or a deep neural network, has a long history and diverse use-cases that include the evaluation of text generation models (Wiseman et al., 2017; Van Der Lee et al., 2019) and the clinical evaluation of human speech (Johnson et al., 2003; Weintraub et al., 2018). One of the earliest examples of such a representation is the Levenshtein distance (Levenshtein, 1966), which describes the number of character-level edits required to transform one piece of text into another. This metric now forms part of a wider family of edit-distance-based metrics that includes the word error rate (WER) and the translation error rate (TER) (Och, 2003). Other algorithms, such as ROUGE (Lin, 2004), METEOR (Banerjee and Lavie, 2005) and the widely used BLEU metric (Papineni et al., 2002), perform exact or approximate  $n$ -gram matching between the two texts.

These low-level approaches bear little resemblance to the human process of comparing two

texts, which benefits from a deep prior understanding of the semantic and syntactic symmetries of language (Novikova et al., 2017). For example, pairs like “she was no ordinary burglar” and “she was an ordinary burglar” are close in edit-distance-space but semantically disparate. The goal of an automatic text evaluation metric is typically to be a good proxy for human judgements, which is clearly task-dependent. More recently, neural approaches have begun to close the gap between automatic and human judgements of semantic text similarity using Transformer-based language models such as BERT (Zhang et al., 2019a; Sellam et al., 2020). They aim to leverage the transferable knowledge gained by the model during pretraining on large text corpora. The relationship between two texts is similarly modelled, albeit implicitly, by sequence-to-sequence models such as BART (Lewis et al., 2019) and T5 (Raffel et al., 2019). We consider paraphrase evaluation and paraphrase generation to be two instances of *paraphrase representation learning*.

Linguistically, a paraphrase is a restatement that preserves essential meaning, with arbitrary levels of literality, fidelity and completeness. In practice, what qualifies as a good paraphrase is context-specific. One motivation for considering paraphrase evaluation as a representation learning problem is the varied nature of paraphrase evaluation tasks, which may have an emphasis on semantic equivalence (e.g. PAWS (Zhang et al., 2019b) and MRPC (Dolan and Brockett, 2005)), logical entailment versus contradiction (e.g. MultiNLI (Williams et al., 2017) and SNLI (Bowman et al., 2015)), and the acceptability of the generated text (e.g. the WMT Metrics Shared Task (Bojar et al., 2017)). Considering even broader applications such as clinical speech analysis further motivates learning generalized paraphrase representations.

In this paper, we introduce ParaBLEU, a paraphrase representation learning model that predicts a

conditioning factor for sequence-to-sequence paraphrase generation as one of its pretraining objectives, inspired by style transfer in text-to-speech (Skerry-Ryan et al., 2018) and text generation systems (Yang et al., 2018; Lample et al., 2018). ParaBLEU addresses the primary issue with neural paraphrase evaluation models to date: the selection of a sufficiently generalized pretraining objective that primes the model for strong performance on downstream paraphrase evaluation tasks when data is scarce. Previous state-of-the-art neural models have either used a broad multi-task learning approach or eschewed additional pretraining altogether. The former case may encourage the model to learn the biases of inferior or inappropriate metrics, while the latter leaves room for optimization. Non-neural models, such as BLEU, TER, ROUGE and BERTScore (Zhang et al., 2019a), benefit from requiring no training data not being subject to domain shift but cannot, however, learn to exploit task-specific nuances of what defines ‘good’ paraphrasing.

We evaluate ParaBLEU’s ability to predict human judgements of paraphrases using the English subset of the 2017 WMT Metrics Shared Task. A useful neural text similarity metric should be robust to data scarcity, so we assess performance as a function of the fine-tuning dataset size. Finally, using the ParaBLEU pretraining model as a paraphrase generation system, we explore our hypothesis that the model reasons in high-level paraphrastic concepts rather than low-level edits through an explainability study, and demonstrate that ParaBLEU can operate as a conditional paraphrase generation model.

## 2 Approach

In this section, we describe and justify the set of inductive biases we build into ParaBLEU, along with a description of the model architecture and pretraining/fine-tuning strategy. We consider a reference text  $x$  and a candidate text  $\hat{x}$ . We wish to learn a function  $f : f(x, \hat{x}) \rightarrow y$ , where  $y \in \mathbb{R}^N$  is a single- or multi-dimensional paraphrase representation, which could be a scalar score.

### 2.1 Inductive biases

Our approach begins by decomposing paraphrase representation learning into three overlapping factors:

1. **Edit-space representation learning:** Build-

ing a representation of high-level syntactic and semantic differences between  $x$  and  $\hat{x}$ , contrasted with the low-level pseudo-syntactic/-semantic operations considered by edit-distance-based and  $n$ -gram based metrics.

2. **Candidate acceptability judgement:** Evaluating the grammaticality, coherence and naturalness of  $\hat{x}$  in isolation. Perplexity (Jelinek et al., 1977) with respect to a given language model is one proxy for this.
3. **Semantic equivalence:** Assessing whether  $x$  and  $\hat{x}$  convey the same essential meaning precisely, as opposed to merely being semantically similar. This is related to entailment classification tasks and, more broadly, the interaction between language and formal logic.

Exploiting this factorization, we hypothesize that the following inductive biases are beneficial to a paraphrase representation learning model:

- **Using pretrained language models:** All three factors require a general understanding of the semantic and syntactic structures of language, making transfer learning from powerful pretrained language models such as BERT (Devlin et al., 2018) appealing.
- **Non-local attention as bitext alignment:** Factors (1) and (3) require performing context-aware ‘matching’ between  $x$  and  $\hat{x}$ . This is similar to the statistical method of bitext alignment (Tiedemann, 2011). Attention mechanisms within a Transformer (Vaswani et al., 2017) are an obvious candidate for learnable context-aware matching, which has precedent in paraphrasing tasks and the next-sentence-prediction objective of the original BERT pre-training. We note that  $x$  and  $\hat{x}$  side-by-side violates attention locality, meaning local attention mechanisms, such as those used in T5, may be suboptimal for longer text-pairs.
- **Bottlenecked conditional generation objective:** A key insight is that a strong factor (1) representation  $z \in \mathbb{R}^M$  where  $h : h(x, \hat{x}) \rightarrow z$  is one that can condition the sampling of  $\hat{x}$  from  $x$  through some generative model  $g : g(x | z) \rightarrow \hat{x}$ . One trivial solution to this is  $h(x, \hat{x}) = \hat{x}$ . To avoid this case, we introduce an information bottleneck on  $z$  such

180 that it is advantageous for the model to learn  
181 to represent high-level abstractions, which are  
182 cheaper than copying  $\hat{x}$  through the bottle-  
183 neck if they are sufficiently abstract compared  
184 with  $\hat{x}$ , the bottleneck is sufficiently tight, and  
185 the decoder can jointly learn the same ab-  
186 stractions. It is likely advantageous to use  
187 a pretrained sequence-to-sequence language  
188 model, which can already reason in linguistic  
189 concepts.

- 190 • **Masked language modelling objective:** Factor  
191 (2) can be addressed by an MLM objective,  
192 which alone is sufficient for a neural network  
193 to learn a language model (Devlin et al., 2018).  
194 Performing masked language modelling on a  
195 reference-candidate pair also encourages the  
196 network to use  $x$  to help unmask  $\hat{x}$  and vice  
197 versa, strengthening the alignment bias useful  
198 for factors (1) and (2).
- 199 • **Entailment classification objective:** Factor  
200 (3) is similar to the classification of whether  
201  $x$  logically entails  $\hat{x}$ . There are a number of  
202 sentence-pair datasets with entailment labels  
203 that could be used to construct this loss; see  
204 Table 5.

## 205 2.2 ParaBLEU

206 Inspired by style transfer in text-to-speech (Skerry-  
207 Ryan et al., 2018) and text generation systems  
208 (Yang et al., 2018; Lample et al., 2018), we propose  
209 the architecture shown in Figure 1. The grey box  
210 indicates the Transformer encoder we wish to pre-  
211 train, which we refer to as the ‘edit encoder’. Factor-  
212 ization of the task leads to three complementary  
213 objectives: a cross-entropy masked language mod-  
214 elling loss  $\mathcal{L}_{MLM}$ , a binary cross-entropy entail-  
215 ment classification loss  $\mathcal{L}_{CLS}$  and a cross-entropy  
216 autoregressive causal language modelling loss  $\mathcal{L}_{AR}$ .  
217 An additional sequence-to-sequence Transformer  
218 model is used during pretraining to provide a learn-  
219 ing signal. The proposed bottleneck lies within the  
220 feedforward network (FFN). The full pretraining  
221 loss is given by:

$$222 \mathcal{L}_{pre} := \mathcal{L}_{AR} + \alpha \cdot \mathcal{L}_{MLM} + \beta \cdot \mathcal{L}_{CLS}, \quad (1)$$

223 where  $\alpha$  and  $\beta$  are tunable hyperparameters. We  
224 probe the importance of each objective in the abla-  
225 tion studies in Section 4.2. At fine-tuning time, the  
226 sequence-to-sequence model is discarded and the  
227 edit encoder is fine-tuned using a linear projection

228 on top of the pooled output. Throughout this work,  
229 our pooling layers simply take the beginning-of-  
230 sequence token.

231 Our architecture places restrictions on valid com-  
232 binations of pretrained models. We found in prac-  
233 tice that using an encoder-only pretrained language  
234 model to initialize the edit encoder, and a sequence-  
235 to-sequence pretrained language model to initialize  
236 the sequence-to-sequence model, works best. This  
237 is likely because encoder-only models are encour-  
238 aged to encode strong representations at the final  
239 layer, and these representations have already been  
240 directly pretrained with an MLM objective. For  
241 technical ease we require that the models have  
242 a consistent tokenizer and vocabulary, and that  
243 the pretrained checkpoints are available through  
244 the HuggingFace transformers library (Wolf  
245 et al., 2019). In this paper, we consider the com-  
246 bination RoBERTa (Liu et al., 2019) + BART, but  
247 we note that both multilingual (XLM-R (Conneau  
248 et al., 2019) + mBART (Liu et al., 2020)) and  
249 long (Longformer + Longformer-Encoder-Decoder  
250 (LED) (Beltagy et al., 2020)) combinations exist.  
251 We consider both base and large variants, which  
252 correspond to RoBERTa<sub>base</sub> and RoBERTa<sub>large</sub>. In  
253 both cases, we use a BART<sub>base</sub> checkpoint.

## 254 2.3 Related work

255 To contextualize this work, we provide a summary  
256 of related architectures and describe the ways in  
257 which they are similar/dissimilar to our proposed  
258 model. BLEURT (Sellam et al., 2020) and FSET  
259 (Kazemnejad et al., 2020) are the most relevant.

260 BLEURT is a neural automatic evaluation met-  
261 rican for text generation. Starting from a pretrained  
262 BERT model, it is further pretrained to predict a  
263 number of pre-existing metrics, such as BLEU,  
264 ROUGE and BERTScore. ParaBLEU, by contrast,  
265 does not use pre-existing metrics as training ob-  
266 jectives, instead using generative conditioning as a  
267 more general signal for paraphrase representation  
268 learning. FSET is a retrieval-based paraphrase gen-  
269 eration system in which a sentence  $z$  is paraphrased  
270 by first locating a similar reference sentence from  
271 a large bank of reference/candidate pairs, then ex-  
272 tracting and replaying similar low-level edits on  $z$ .  
273 Common to ParaBLEU and FSET is the use of a  
274 Transformer for paraphrase style transfer, with dif-  
275 fering architectural details. However, FSET is de-  
276 signed to transpose low-level edits and so requires  
277 lexically similar examples; whereas ParaBLEU is

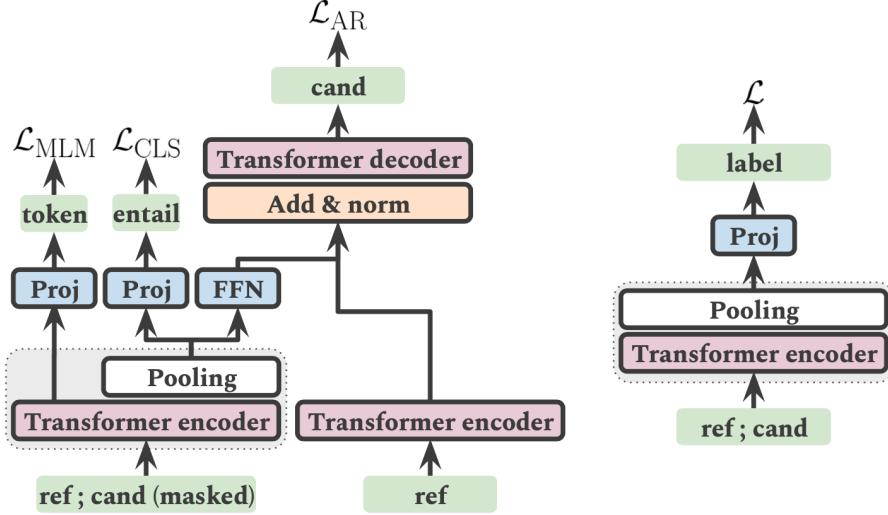


Figure 1: The pretraining (left) and fine-tuning (right) setups for ParaBLEU. ‘ref ; cand’ indicates the canonical method for combining a reference and candidate sentence for a given language model.  $\mathcal{L}_{\text{AR}}$  is an autoregressive causal language modelling loss,  $\mathcal{L}_{\text{MLM}}$  a masked language modelling loss, and  $\mathcal{L}_{\text{CLS}}$  an entailment classification loss. The feedforward network (FFN) includes two affine layers, the middle dimension of which can be used to create a bottleneck (see Section 2.1). Dropout layers and activations are omitted for brevity.

278 explicitly designed to learn high-level, reference-  
 279 invariant paraphrase representations using a factorized  
 280 objective. The **musical style Transformer**  
 281 **autoencoder** (Choi et al., 2020) uses a similar  
 282 Transformer-based style transfer architecture to  
 283 conditionally generate new music in controllable  
 284 styles. Other examples in text-to-speech systems  
 285 perform style transfer by encoding the prosody of  
 286 a source sentence into a bottlenecked reference em-  
 287 bedding (Skerry-Ryan et al., 2018) or disentangled  
 288 style tokens (Wang et al., 2018b).

289 There is a wealth of recent literature on con-  
 290 trollable paraphrase generation and linguistic style  
 291 transfer, some of which we highlight here. **T5**  
 292 leverages a huge text corpus as pretraining for con-  
 293 ditional generation using ‘commands’ encoded as  
 294 text, which includes paraphrastic tasks such as sum-  
 295 marization. **Linguistic style transfer** (Yang et al.,  
 296 2018; Zhao et al., 2018; Jin et al., 2020) work aims  
 297 to extract the style of a piece of text and map it  
 298 onto another piece of text without changing its se-  
 299 mantic meaning. **STRAP** (Krishna et al., 2020)  
 300 generates paraphrases in controllable styles by mix-  
 301 ing and matching multiple style-specific fine-tuned  
 302 GPT-2 models. **REAP** (Goyal and Durrett, 2020)  
 303 uses a Transformer to syntactically diverse generate  
 304 paraphrases by including an additional position em-  
 305 bedding representing the syntactic tree. **DNPG** (Li  
 306 et al., 2019) is a paraphrase generation system that  
 307 uses a cascade of Transformer encoders/decoders to

control whether paraphrasing is sentential/phrasal.

### 3 Data

In this section, we describe the pretraining and  
 310 fine-tuning datasets we use in our studies.  
 311

#### 3.1 WMT Metrics Shared Task

The WMT Metrics Shared Task is an annual bench-  
 313 mark for automated evaluation metrics for transla-  
 314 tion systems, where the goal is to predict average  
 315 human ratings comparing the machine-translated  
 316 candidate  $\hat{x}$  with human-translated reference  $x$ ,  
 317 both of which have been translated from the same  
 318 source sentence.  
 319

We use an identical setup to (Sellam et al., 2020)  
 320 and (Zhang et al., 2019a), where we use the subset  
 321 of data for which the candidate and reference are  
 322 in English, which we will refer to as the to-English  
 323 subset. The source, which is unused, can be in any  
 324 non-English language, the set of which varies from  
 325 year-to-year. We produce results for the WMT Met-  
 326 rics Shared Task 2017 (WMT17) using the official  
 327 test, and train on the to-English subsets of WMT15  
 328 and WMT16. The training sets contains 5,360 ex-  
 329 amples. The distributions of example length in  
 330 tokens is shown in Figure 2.  
 331

We report the agreement between the metric and  
 332 the human scores using two related correlation co-  
 333 efficients: absolute Kendall  $|\tau|$  and absolute Pear-  
 334 son  $|r|$ , the latter of which was the official metric  
 335

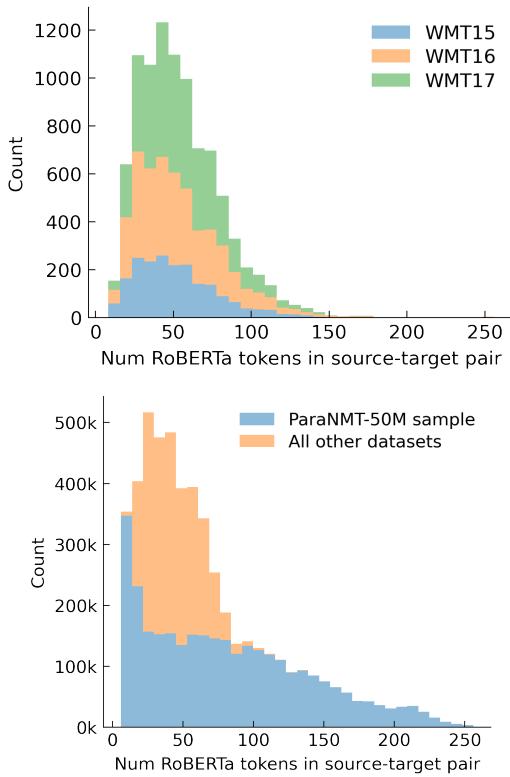


Figure 2: Stacked histograms showing the distribution of the number of tokens in the WMT Metrics Shared Task data (left) and ParaCorpus (right).

of the 2017 task. In our summary results in the main paper, we average these metrics across all source languages but not over reference/candidate language. Full results are provided in Appendix D.

### 3.2 ParaCorpus

In addition to our design choices, we also encourage a robust and generalizable pretraining by using a dataset covers a variety of styles and lengths. We collate a number of paraphrase datasets to create a single pretraining dataset we call ParaCorpus. The composition of the dataset is shown in Table 5, with a total of  $\sim 5.1$ m examples. All examples have reference and candidate texts and around one third additionally have binary entailment labels. Where the source dataset included three-way labels ‘entailment’/‘contradiction’/‘neutral’, ‘entailment’ was mapped to 1 and the others to 0. A subset of ParaNMT-50M (Wieting and Gimpel, 2017), which includes noisier, speech-like examples, was included to add additional stylistic diversity to the dataset, and to increase the population of the dataset with combined token lengths above 128, which we hypothesize will make the model more robust to the longer examples seen in the WMT

Table 1: Summary results for WMT17. The metrics reported are absolute Kendall  $|\tau|$  and Pearson  $|r|$  averaged across each source language. Full results can be found in Appendix D.

Model	$ \tau $	$ r $
BLEU	0.292	0.423
TER	0.352	0.475
ROUGE	0.354	0.518
METEOR	0.301	0.443
chrF++	0.396	0.578
BLEURT-large	0.625	0.818
BERTScore-RoBERTa <sub>large</sub>	0.567	0.759
BERTScore-T5 <sub>large</sub>	0.536	0.738
BERTScore-DeBERTa <sub>large</sub>	0.580	0.773
MoverScore	0.322	0.454
ParaBLEU <sub>large</sub>	<b>0.653</b>	<b>0.843</b>
ParaBLEU <sub>base</sub>	0.589	0.785

datasets. Token lengths are shown in Figure 2.

## 4 Experiments

In this section, we present results on WMT17, benchmarked against the current state-of-the-art approach, along with widely used neural,  $n$ -gram and edit-distance-based metrics. We study ParaBLEU performance as a function of number of pretraining steps and the size of the fine-tuning dataset. Finally, we perform ablations to test the impact of the inductive biases and resultant architectural decisions described in Section 2.

We report results for both ParaBLEU<sub>base</sub>, based on RoBERTa<sub>base</sub> (12 layers, 768 hidden units, 12 heads), and our default model ParaBLEU<sub>large</sub>, based on RoBERTa<sub>large</sub> (24 layers, 1,024 hidden units, 16 heads). Both models are trained near-identically for 4 epochs on ParaCorpus. Further pretraining details can be found in Appendix A. For fine-tuning, we use a batch size of 32, a learning rate of 1e-5 and train for 40k steps, with a validation set size of 10% (unless otherwise stated). No reference texts are shared between the train and validation sets, following (Sellam et al., 2020). Pre-training ParaBLEU<sub>large</sub> takes  $\sim 10$ h on a 16 A100 GPU machine. Fine-tuning takes  $\sim 8$ h on a single A100 GPU machine.

### 4.1 Results

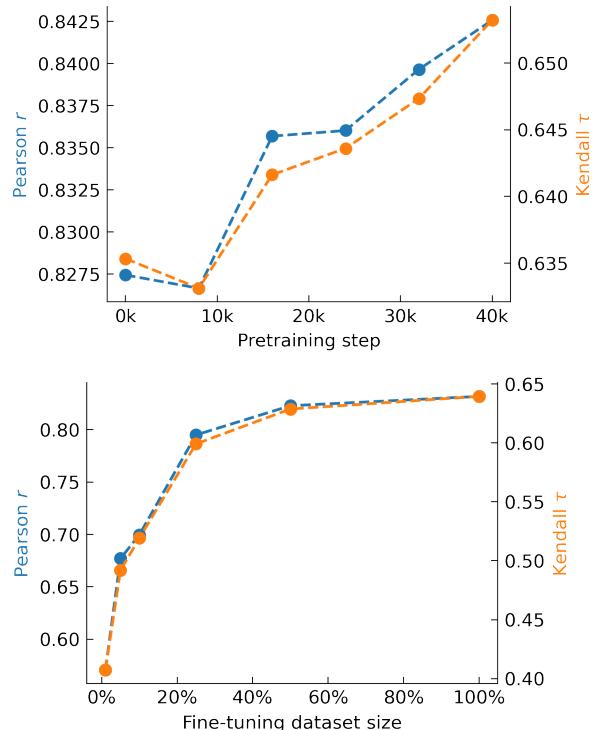
ParaBLEU results on WMT17 are given in Table 1, along with a number of baselines. Base-

389 lines include BLEURT, described in Section 2.3),  
 390 along with BERTScore, a non-learned neural metric  
 391 that uses a matching algorithm on top of neural  
 392 word embeddings, similar to  $n$ -gram matching ap-  
 393 proaches. MoverScore (Zhao et al., 2019) is similar  
 394 to BERTScore, but uses an optimal transport algo-  
 395 rithm. BLEU, ROUGE, METEOR and chrF++ are  
 396 widely used  $n$ -gram-based methods, working at  
 397 the word, subword or character level. TER is an  
 398 edit-distance-based metric, similar to WER.

399 ParaBLEU<sub>large</sub> achieves new state-of-the-art re-  
 400 sults on WMT17, exceeding the previous state-of-  
 401 the-art approach, BLEURT, on both correlation  
 402 metrics. We note that non-neural metrics perform  
 403 the worst, of which the character-level  $n$ -gram-  
 404 matching algorithm chrF++ performs the best. Non-  
 405 learned neural metrics (BERTScore and Mover-  
 406 Score) tend to perform better, and learned neural  
 407 metrics (BLEURT and ParaBLEU) perform the  
 408 best. BLEU, the most widely used metric, has  
 409 the poorest correlation with human judgements.  
 410 This is consistent with results seen previously in  
 411 the literature (Zhang et al., 2019a; Sellam et al.,  
 412 2020). The significant drop in performance from  
 413 ParaBLEU<sub>large</sub> to ParaBLEU<sub>base</sub> highlights the  
 414 benefit of larger, more expressive pretrained lan-  
 415 guage models.

416 Figure 3 probes performance as a function of  
 417 number of pretraining steps and the size of the  
 418 fine-tuning dataset for ParaBLEU<sub>large</sub>. As expected,  
 419 pretraining for longer increases downstream task  
 420 performance. However, we note that 40k steps,  
 421 approximately 4 epochs of ParaCorpus, does not  
 422 yet reach diminishing returns on WMT17 perfor-  
 423 mance. We therefore recommend pretraining for  
 424 significantly longer. Both BERT and RoBERTa  
 425 are pretrained for 40 epochs (Liu et al., 2019; Lan  
 426 et al., 2019); the T5 authors ablate their dataset size  
 427 at a fixed number of steps and conclude that per-  
 428 formance does not significantly degrade up to and  
 429 including 64 epochs (Raffel et al., 2019); conversely,  
 430 the BLEURT authors see diminishing returns on  
 431 downstream task performance after 2 pretraining  
 432 epochs (Sellam et al., 2020).

433 For the fine-tuning dataset size study, we con-  
 434 sistently use a validation set size of 25% to facilitate  
 435 the small-data results. Despite the training set (the  
 436 English subsets of WMT15 and WMT16) forming  
 437 a relatively small dataset, ParaBLEU<sub>large</sub> trained  
 438 on 50% of the available data (2,010 training ex-  
 439 amples, 670 validation examples) still beats the



440 Figure 3: Performance of ParaBLEU<sub>large</sub> on WMT17  
 441 as a function of number of pretraining steps (left) and  
 442 the fine-tuning dataset size (right). Note that the Pearson  
 443  $r$  results (blue) use the left  $y$ -axis, whereas Kendall  $\tau$   
 444 (orange) uses the right  $y$ -axis.

445 previous state-of-the-art, BLEURT, yielding a Pear-  
 446 son correlation of 0.823. The impact of reducing  
 447 the train size from 100% (4,020 training examples,  
 448 1,340 validation examples) to 25% (1,005 training  
 449 examples, 335 validation examples) has a relatively  
 450 small effect on performance, reducing Pearson  $r$   
 451 from 0.832 to 0.795. With a dataset size of only  
 452 1% (40 training examples, 14 validation examples),  
 453 ParaBLEU<sub>large</sub> achieves a Pearson  $r$  of 0.571, still  
 454 correlating significantly more strongly with human  
 455 judgements than BLEU, TER, ROUGE, METEOR  
 456 and MoverScore. We attribute this to the suitability  
 457 of the generalized pretraining objective for priming  
 458 the model for paraphrase evaluation tasks.

## 4.2 Ablations

459 To more directly test the hypotheses in Section 2.1,  
 460 we perform ablations in which we remove each  
 461 component of the factorized objective in turn. The  
 462 results of this are shown in Table 2. Each part  
 463 of the objective is associated with an increase in  
 464 downstream task performance. The most signifi-  
 465 cant degradation comes from removing the MLM  
 466 loss. Possible reasons for this include: the MLM

Table 2: Ablation results on WMT17. The metrics reported are the absolute Kendall  $|\tau|$  and Pearson  $|r|$  correlation coefficients averaged across each reference language.

Model	$ \tau $	$ r $
Baseline (ParaBLEU <sub>large</sub> )	<b>0.653</b>	<b>0.843</b>
No MLM loss ( $\mathcal{L}_{MLM}$ )	0.633	0.826
No autoregressive loss ( $\mathcal{L}_{AR}$ )	0.642	0.834
No entailment classification loss ( $\mathcal{L}_{CLS}$ )	0.644	0.837

loss’ contribution to candidate acceptability judgement are crucial; the MLM loss acts as a regularizer, encouraging the edit encoder to represent paraphrases in linguistic concepts rather than low-level edits; and the MLM loss further encourages bitext alignment behaviour, as described in Section 2.1.

## 5 One-shot paraphrase generation

As our final study, we exploit the generative nature of the pretraining architecture to test our claim that the edit encoder reasons in high-level paraphrastic concepts rather than low-level edits. To do this, we diverge from the pretraining setup, in which the same reference text is passed to both the edit encoder and the sequence-to-sequence model, by passing a different, unseen reference to the sequence-to-sequence model. Akin to (Brown et al., 2020; Gao et al., 2020), the hope is that the ‘demonstration paraphrase’ acts as a conditioning factor for paraphrasing the unseen sentence in a similar way.

If the model is reasoning in low-level edits or otherwise ‘cheating’, we expect to see:

- Thematic/word leakage from the encoder candidate to the generated candidate, caused by the candidate being autoencoded. This is the undesirable behaviour we sought to address using a bottleneck.
- Ungrammatical or otherwise unacceptable output with made-up words and/or bad word order, caused by the encoding of low-level edits scrambling the generator reference tokens.

If the model is reasoning in high-level paraphrastic concepts, we expect to see:

- Consistently grammatical, acceptable output.

- The flavour of the paraphrase mirroring the conditioning, e.g. the altering of a linguistic style, mood or tense.

We generate text using beam-search (Medress et al., 1977). We sample references at random from the MRPC dataset. The demonstration candidate is a hand-crafted paraphrase of the demonstration reference that embodies a pre-specified paraphrase type. We report the predicted entailment score of the demonstration reference and candidate, along with the candidate generated by the model. The examples in Table 3 and in Appendix E are selected at random.

The generation results are shown in Table 3. We include two sets of results for each paraphrastic type (e.g. ‘negative’): one where the demonstration reference/candidate differ in this concept, and one where both embody the concept. Since we wish to encode the *difference* between the demonstration reference/candidate texts, the desired behaviour when the demonstration pair is identical is no change. If this is not the case, it is likely that the edit encoder is just autoencoding the candidate using high-level linguistic concepts, similar to linguistic style transfer. Further examples are given in Appendix E.

The results present a strong case that the encoder is representing high-level paraphrastic concepts. It is able to successfully identify changes in mood, style and tense between the demonstration reference and candidate, and transpose them onto the unseen reference to make a largely grammatical and appropriately paraphrased sentence. We do not see significant leakage of concepts, words or styles between the demonstration candidate and the generated candidate, instead the expected transfer of paraphrase style.

## 6 Conclusions

In this paper, we introduced ParaBLEU, a paraphrase representation learning model and associated paraphrase evaluation metric. We demonstrated that the metric yields state-of-the-art correlation with human paraphrase judgements and is robust to data scarcity. We motivated its pre-training strategy through a set of inductive biases, which we tested through ablation studies. Finally, we reframed the pretraining as a one-shot paraphrase generation model and gathered evidence that ParaBLEU represents meaningful paraphrastic information.

Table 3: One-shot paraphrase generation results. In every case, the reference is “*Congo’s war began in 1998 when Uganda and Rwanda invaded to back rebels fighting to topple the central government.*” The difference between the identity and non-identity version is shown in boldface. ‘Ent.’ is the predicted demonstration reference/candidate entailment probability from the classification loss.

Type	Demo reference	Demo candidate	Ent.	Generated candidate
negation	Both bidders agreed to assume about \$90 million in debt owed on the planes	The bidders did not agree to assume about \$90 million in debt owed on the planes	0.0%	The Congo War began in 1998 when Uganda and Rwanda <b>did not invade</b> to back rebels fighting to topple the central government.
identity (negation)	The bidders did not agree to assume about \$90 million in debt owed on the planes	The bidders did not agree to assume about \$90 million in debt owed on the planes	92.6%	The Congo War began in 1998 when Uganda and Rwanda invaded to back rebels fighting to topple the central government.
summary	Both bidders agreed to assume about \$90 million in debt owed on the planes	They agreed to assume the debt	99.0%	<b>They invaded to back</b> the rebels fighting to topple the central government.
identity (summary)	They agreed to assume the debt	They agreed to assume the debt	95.6%	Congo began in 1998 when Uganda and Rwanda invaded to support the rebels fighting to topple the central government.
passive	Both bidders agreed to assume about \$90 million in debt owed on the planes	\$90 million in debt owed on the planes was agreed to be assumed by both bidders	99.8%	<b>The war in Congo</b> began in 1998 when Uganda and Rwanda <b>were invaded to back</b> rebels fighting to topple the central government.
identity (passive)	\$90 million in debt owed on the planes was agreed to be assumed by both bidders	\$90 million in debt owed on the planes was agreed to be assumed by both bidders	92.7%	Congo War began in 1998 when Uganda and Rwanda invaded to back rebels fighting to topple the central government.
conditional perfect	Both bidders agreed to assume about \$90 million in debt owed on the planes	Both bidders would have agreed to assume about \$90 million in debt owed on the planes	69.4%	<b>Congo’s war</b> began in 1998, when Uganda and Rwanda <b>would have invaded to back</b> rebels fighting to topple the central government.
identity (conditional perfect)	Both bidders would have agreed to assume about \$90 million in debt owed on the planes	Both bidders would have agreed to assume about \$90 million in debt owed on the planes	90.7%	Congo War began in 1998 when Uganda and Rwanda invaded to back rebels fighting to topple the central government.
personal	Both bidders agreed to assume about \$90 million in debt owed on the planes	We agreed to assume about \$90 million in debt owed on the planes	95.4%	<b>We started the Congo war</b> in 1998 when Uganda and Rwanda invaded to support rebels fighting to topple the central government.
identity (personal)	We agreed to assume about \$90 million in debt owed on the planes	We agreed to assume about \$90 million in debt owed on the planes	93.2%	The Congo War began in 1998 when Uganda and Rwanda invaded to back rebels fighting to topple the central government.
informal	Both bidders agreed to assume about \$90 million in debt owed on the planes	Both guys agreed to take on like \$90 million that they owed for some planes	98.8%	<b>Congo’s war started</b> in 1998 when Uganda and Rwanda invaded to <b>get some</b> rebels fighting to topple the central government.
identity (informal)	Both guys agreed to take on like \$90 million that they owed for some planes	Both guys agreed to take on like \$90 million that they owed for some planes	94.5%	Congo War began in 1998 when Uganda and Rwanda invaded to back rebels fighting to topple the central government.

## References

- Satanjeev Banerjee and Alon Lavie. 2005. Meteor: An automatic metric for mt evaluation with improved correlation with human judgments. In *Proceedings of the acl workshop on intrinsic and extrinsic evaluation measures for machine translation and/or summarization*, pages 65–72.
- Iz Beltagy, Matthew E Peters, and Arman Cohan. 2020. Longformer: The long-document transformer. *arXiv preprint arXiv:2004.05150*.
- Ondřej Bojar, Yvette Graham, and Amir Kamran. 2017. Results of the WMT17 metrics shared task. In *Proceedings of the Second Conference on Machine Translation*, pages 489–513, Copenhagen, Denmark. Association for Computational Linguistics.
- Samuel R Bowman, Gabor Angeli, Christopher Potts, and Christopher D Manning. 2015. A large annotated corpus for learning natural language inference. *arXiv preprint arXiv:1508.05326*.
- Tom B Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared Kaplan, Prafulla Dhariwal, Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, et al. 2020. Language models are few-shot learners. *arXiv preprint arXiv:2005.14165*.
- Kristy Choi, Curtis Hawthorne, Ian Simon, Monica Dinulescu, and Jesse Engel. 2020. Encoding musical style with transformer autoencoders. In *International Conference on Machine Learning*, pages 1899–1908. PMLR.
- Alexis Conneau, Kartikay Khandelwal, Naman Goyal, Vishrav Chaudhary, Guillaume Wenzek, Francisco Guzmán, Edouard Grave, Myle Ott, Luke Zettlemoyer, and Veselin Stoyanov. 2019. Unsupervised cross-lingual representation learning at scale. *arXiv preprint arXiv:1911.02116*.
- Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. 2018. Bert: Pre-training of deep bidirectional transformers for language understanding. *arXiv preprint arXiv:1810.04805*.
- William B Dolan and Chris Brockett. 2005. Automatically constructing a corpus of sentential paraphrases. In *Proceedings of the Third International Workshop on Paraphrasing (IWP2005)*.
- Qingxiu Dong, Xiaojun Wan, and Yue Cao. 2021. Parasci: A large scientific paraphrase dataset for longer paraphrase generation. *arXiv preprint arXiv:2101.08382*.
- Tianyu Gao, Adam Fisch, and Danqi Chen. 2020. Making pre-trained language models better few-shot learners. *arXiv preprint arXiv:2012.15723*.
- Tanya Goyal and Greg Durrett. 2020. Neural syntactic preordering for controlled paraphrase generation. *arXiv preprint arXiv:2005.02013*.
- Fred Jelinek, Robert L Mercer, Lalit R Bahl, and James K Baker. 1977. Perplexity—a measure of the difficulty of speech recognition tasks. *The Journal of the Acoustical Society of America*, 62(S1):S63–S63.
- Di Jin, Zhijing Jin, Zhitong Hu, Olga Vechtomova, and Rada Mihalcea. 2020. Deep learning for text style transfer: A survey. *CoRR, abs/2011.00416*.
- David K Johnson, Martha Storandt, and David A Balota. 2003. Discourse analysis of logical memory recall in normal aging and in dementia of the alzheimer type. *Neuropsychology*, 17(1):82.
- Amirhossein Kazemnejad, Mohammadreza Salehi, and Mahdieh Soleymani Baghshah. 2020. Paraphrase generation by learning how to edit from samples. In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, pages 6010–6021.
- Kalpesh Krishna, John Wieting, and Mohit Iyyer. 2020. Reformulating unsupervised style transfer as paraphrase generation. *arXiv preprint arXiv:2010.05700*.
- Guillaume Lample, Sandeep Subramanian, Eric Smith, Ludovic Denoyer, Marc’Aurelio Ranzato, and Y-Lan Boureau. 2018. Multiple-attribute text rewriting. In *International Conference on Learning Representations*.
- Zhenzhong Lan, Mingda Chen, Sebastian Goodman, Kevin Gimpel, Piyush Sharma, and Radu Soricut. 2019. Albert: A lite bert for self-supervised learning of language representations. *arXiv preprint arXiv:1909.11942*.
- Vladimir I Levenshtein. 1966. Binary codes capable of correcting deletions, insertions, and reversals. In *Soviet physics doklady*, pages 707–710. Soviet Union.
- Mike Lewis, Yinhan Liu, Naman Goyal, Marjan Ghazvininejad, Abdelrahman Mohamed, Omer Levy, Ves Stoyanov, and Luke Zettlemoyer. 2019. Bart: Denoising sequence-to-sequence pre-training for natural language generation, translation, and comprehension. *arXiv preprint arXiv:1910.13461*.
- Zichao Li, Xin Jiang, Lifeng Shang, and Qun Liu. 2019. Decomposable neural paraphrase generation. *arXiv preprint arXiv:1906.09741*.
- Chin-Yew Lin. 2004. Rouge: A package for automatic evaluation of summaries. In *Text summarization branches out*, pages 74–81.
- Yinhan Liu, Jiatao Gu, Naman Goyal, Xian Li, Sergey Edunov, Marjan Ghazvininejad, Mike Lewis, and Luke Zettlemoyer. 2020. Multilingual denoising pre-training for neural machine translation. *Transactions of the Association for Computational Linguistics*, 8:726–742.
- Yinhan Liu, Myle Ott, Naman Goyal, Jingfei Du, Mandar Joshi, Danqi Chen, Omer Levy, Mike Lewis, Luke Zettlemoyer, and Veselin Stoyanov. 2019.

655	Roberta: A robustly optimized bert pretraining approach. <i>arXiv preprint arXiv:1907.11692</i> .	708
656		709
657		710
658		711
659		712
660		713
661		714
662		
663	Mark F. Medress, Franklin S Cooper, Jim W. Forgie, CC Green, Dennis H. Klatt, Michael H. O’Malley, Edward P Neuburg, Allen Newell, DR Reddy, B Ritea, et al. 1977. Speech understanding systems: Report of a steering committee. <i>Artificial Intelligence</i> , 9(3):307–316.	715
664		716
665		717
666		718
667	Jekaterina Novikova, Ondřej Dušek, Amanda Cercas Curry, and Verena Rieser. 2017. Why we need new evaluation metrics for nlg. <i>arXiv preprint arXiv:1707.06875</i> .	719
668		720
669		721
670		
671	Franz Josef Och. 2003. Minimum error rate training in statistical machine translation. In <i>Proceedings of the 41st annual meeting of the Association for Computational Linguistics</i> , pages 160–167.	722
672		723
673		724
674		725
675		
676	Kishore Papineni, Salim Roukos, Todd Ward, and Wei Jing Zhu. 2002. Bleu: a method for automatic evaluation of machine translation. In <i>Proceedings of the 40th annual meeting of the Association for Computational Linguistics</i> , pages 311–318.	726
677		727
678		728
679		729
680		
681	Colin Raffel, Noam Shazeer, Adam Roberts, Katherine Lee, Sharan Narang, Michael Matena, Yanqi Zhou, Wei Li, and Peter J Liu. 2019. Exploring the limits of transfer learning with a unified text-to-text transformer. <i>arXiv preprint arXiv:1910.10683</i> .	730
682		731
683		732
684	Thibault Sellam, Dipanjan Das, and Ankur P Parikh. 2020. Bleurt: Learning robust metrics for text generation. <i>arXiv preprint arXiv:2004.04696</i> .	733
685		734
686		735
687		736
688		737
689	RJ Skerry-Ryan, Eric Battenberg, Ying Xiao, Yuxuan Wang, Daisy Stanton, Joel Shor, Ron Weiss, Rob Clark, and Rif A Saurous. 2018. Towards end-to-end prosody transfer for expressive speech synthesis with tacotron. In <i>international conference on machine learning</i> , pages 4693–4702. PMLR.	738
690		
691	Jörg Tiedemann. 2011. Bitext alignment. <i>Synthesis Lectures on Human Language Technologies</i> , 4(2):1–165.	739
692		740
693		741
694		742
695		
696	Chris Van Der Lee, Albert Gatt, Emiel Van Miltenburg, Sander Wubben, and Emiel Krahmer. 2019. Best practices for the human evaluation of automatically generated text. In <i>Proceedings of the 12th International Conference on Natural Language Generation</i> , pages 355–368.	743
697		744
698		745
699		746
700		
701	Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz Kaiser, and Illia Polosukhin. 2017. Attention is all you need. <i>arXiv preprint arXiv:1706.03762</i> .	747
702		748
703	Alex Wang, Amanpreet Singh, Julian Michael, Felix Hill, Omer Levy, and Samuel R Bowman. 2018a. Glue: A multi-task benchmark and analysis platform for natural language understanding. <i>arXiv preprint arXiv:1804.07461</i> .	749
704		
705		
706		
707		
708	Yuxuan Wang, Daisy Stanton, Yu Zhang, RJ-Skerry Ryan, Eric Battenberg, Joel Shor, Ying Xiao, Ye Jia, Fei Ren, and Rif A Saurous. 2018b. Style tokens: Unsupervised style modeling, control and transfer in end-to-end speech synthesis. In <i>International Conference on Machine Learning</i> , pages 5180–5189. PMLR.	750
709		751
710		752
711		753
712		754
713		
714		
715	Sandra Weintraub, Lilah Besser, Hiroko H Dodge, Merilee Teylan, Steven Ferris, Felicia C Goldstein, Bruno Giordani, Joel Kramer, David Loewenstein, Dan Marson, et al. 2018. Version 3 of the alzheimer disease centers’ neuropsychological test battery in the uniform data set (uds). <i>Alzheimer disease and associated disorders</i> , 32(1):10.	755
716		756
717		757
718		758
719		
720		
721		
722	John Wieting and Kevin Gimpel. 2017. Paranmt-50m: Pushing the limits of paraphrastic sentence embeddings with millions of machine translations. <i>arXiv preprint arXiv:1711.05732</i> .	759
723		760
724		761
725		
726	Adina Williams, Nikita Nangia, and Samuel R Bowman. 2017. A broad-coverage challenge corpus for sentence understanding through inference. <i>arXiv preprint arXiv:1704.05426</i> .	762
727		763
728		764
729		
730	Sam Wiseman, Stuart M Shieber, and Alexander M Rush. 2017. Challenges in data-to-document generation. <i>arXiv preprint arXiv:1707.08052</i> .	765
731		766
732		
733	Thomas Wolf, Lysandre Debut, Victor Sanh, Julien Chaumond, Clement Delangue, Anthony Moi, Pierrette Cistac, Tim Rault, Rémi Louf, Morgan Funtowicz, et al. 2019. Huggingface’s transformers: State-of-the-art natural language processing. <i>arXiv preprint arXiv:1910.03771</i> .	767
734		768
735		769
736		770
737		771
738		
739	Zichao Yang, Zhiting Hu, Chris Dyer, Eric P Xing, and Taylor Berg-Kirkpatrick. 2018. Unsupervised text style transfer using language models as discriminators. <i>arXiv preprint arXiv:1805.11749</i> .	772
740		773
741		774
742		775
743	Tianyi Zhang, Varsha Kishore, Felix Wu, Kilian Q Weinberger, and Yoav Artzi. 2019a. Bertscore: Evaluating text generation with bert. <i>arXiv preprint arXiv:1904.09675</i> .	776
744		777
745		778
746		779
747	Yuan Zhang, Jason Baldridge, and Luheng He. 2019b. Paws: Paraphrase adversaries from word scrambling. <i>arXiv preprint arXiv:1904.01130</i> .	780
748		781
749		
750	Wei Zhao, Maxime Peyrard, Fei Liu, Yang Gao, Christian M Meyer, and Steffen Eger. 2019. Moverscore: Text generation evaluating with contextualized embeddings and earth mover distance. <i>arXiv preprint arXiv:1909.02622</i> .	782
751		783
752		784
753		785
754		786
755	Yanpeng Zhao, Wei Bi, Deng Cai, Xiaojiang Liu, Kewei Tu, and Shuming Shi. 2018. Language style transfer from sentences with arbitrary unknown styles. <i>arXiv preprint arXiv:1808.04071</i> .	787
756		788
757		789
758		

759            **A Pretraining hyperparameters**

760            Table 4 shows the hyperparamaters used for the  
761            ParaBLEU<sub>base</sub> and ParaBLEU<sub>large</sub> models during  
762            pretraining.  $\alpha$  and  $\beta$  are the loss weights from  
763            Equation 1.

764            **B ParaCorpus**

765            Table 5 provides a description of the composition  
766            of the pretraining dataset.

767            **C Microsoft Research Paraphrase  
768            Corpus results**

769            We additionally ran a study on the Microsoft Re-  
770            search Paraphrase Corpus (MRPC) (Dolan and  
771            Brockett, 2005), a constituent of the GLUE bench-  
772            mark (Wang et al., 2018a). MRPC contains 5,801  
773            sentence pairs each accompanied by hand-labelled  
774            binary judgement of whether the pair constitutes a  
775            paraphrase. The data is split into a train set (4,076  
776            sentence pairs of which 2,753 are paraphrases) and  
777            a test set (1,725 sentence pairs of which 1,147 are  
778            paraphrases).

779            We fine-tune our ParaBLEU models on the  
780            MRPC train set using the fine-tuning procedure  
781            detailed in (Liu et al., 2019) and predict on the held-  
782            out test set. For baselines we use the ALBERT<sub>large</sub>  
783            (Lan et al., 2019) and the RoBERTa<sub>large</sub> (Liu et al.,  
784            2019) models fine-tuned using their respective hy-  
785            perparameters.

Model	Accuracy	F1 score
ALBERT <sub>large</sub>	88.2	91.3
RoBERTa <sub>large</sub>	89.5	92.2
ParaBLEU <sub>large</sub>	88.8	91.5
ParaBLEU <sub>base</sub>	85.2	88.9

Table 6: The results from the Microsoft Research Paraphrase Corpus (MRPC).

786            From Table 6 we observe that our default model  
787            ParaBLEU<sub>large</sub> underperforms compared to the  
788            model it is based on, RoBERTa<sub>large</sub>. This could  
789            be because the hyperparameter sweep we used for  
790            our ParaBLEU models (the same sweep as recom-  
791            mended by the authors of RoBERTa<sub>large</sub>) is subop-  
792            timal and a broader hyperparameter sweep may be  
793            required.

Hyperparameter	ParaBLEU <sub>base</sub>	ParaBLEU <sub>large</sub>
Edit encoder base model	RoBERTa <sub>base</sub>	RoBERTa <sub>large</sub>
Sequence-to-sequence base model	BART <sub>base</sub>	BART <sub>base</sub>
Batch size (per GPU; examples)	64	32
Batch size (per GPU; max tokens)	16,384	8,192
Learning rate (per GPU)	4e-4	1e-4
Warmup steps	1,200	2,400
Train length (updates)	20k	40k
Train length (epochs)	4	4
Gradient accumulation steps	1	2
$\alpha$	2.0	2.0
$\beta$	10.0	10.0

Table 4: Pretraining hyperparameters for the ParaBLEU<sub>base</sub> and ParaBLEU<sub>large</sub> models used in this paper. These were adapted for a larger architecture from the RoBERTa paper (Liu et al., 2019) and not subject to tuning.

Table 5: ParaCorpus composition.

Dataset	Subsets included	Nature	Size	Ent. labels	Ref
PAWS	Wiki-train; train	QQP- Sentence pairs with high semantic over- lap	740k	✓	(Zhang et al., 2019b)
SNLI	Train	Human-written entailment sentence pairs	550k	✓†	(Bowman et al., 2015)
MultiNLI	Train	Multi-genre entail- ment sentence pairs	390k	✓†	(Williams et al., 2017)
ParaSCI	ACL-train; arXiv- train	Human-written academic paraphrase pairs	340k	✗	(Dong et al., 2021)
ParaNMT-50M	Random sample (see main text)	Varied paraphrase pairs from machine translation	3.1m	✗	(Wieting and Gimpel, 2017)
ParaCorpus	-	-	5.1m	Partial	-

## D Full to-English WMT results

Table 7 shows the full WMT17 results, which are summarized in main paper Table 1. See Section 4.1 for more details.

Table 7: Full to-English results for WMT17. The metrics reported are absolute Kendall  $|\tau|$  and Pearson  $|r|$ . Models are fine-tuned on the English subset of WMT15 and WMT16. For a language pair ‘x-y’, the original reference was in language ‘x’, and both human and machine translations are in language ‘y’. For results averaged across all source languages, see the main paper.

Model	lv-en $ \tau  /  r $	tr-en $ \tau  /  r $	zh-en $ \tau  /  r $	ru-en $ \tau  /  r $	de-en $ \tau  /  r $	cs-en $ \tau  /  r $	fi-en $ \tau  /  r $
BLEU	0.215 / 0.334	0.313 / 0.461	0.344 / 0.488	0.313 / 0.431	0.259 / 0.372	0.255 / 0.373	0.342 / 0.503
TER	0.329 / 0.439	0.393 / 0.472	0.365 / 0.493	0.358 / 0.509	0.295 / 0.403	0.315 / 0.458	0.411 / 0.548
ROUGE	0.303 / 0.459	0.395 / 0.56	0.366 / 0.542	0.343 / 0.488	0.336 / 0.488	0.302 / 0.462	0.434 / 0.628
METEOR	0.258 / 0.403	0.375 / 0.554	0.352 / 0.521	0.353 / 0.491	0.307 / 0.445	0.287 / 0.448	0.402 / 0.597
MoverScore	0.252 / 0.350	0.314 / 0.493	0.345 / 0.485	0.375 / 0.493	0.296 / 0.401	0.317 / 0.433	0.356 / 0.521
chRF++	0.333 / 0.520	0.432 / 0.614	0.405 / 0.593	0.415 / 0.588	0.365 / 0.534	0.35 / 0.523	0.475 / 0.678
BLEURT	0.644 / 0.835	0.629 / 0.824	0.602 / 0.814	0.613 / 0.811	0.599 / 0.792	0.593 / 0.773	0.695 / 0.878
BERTScore-RoBERT <sub>large</sub>	0.555 / 0.756	0.569 / 0.751	0.568 / 0.775	0.555 / 0.746	0.554 / 0.745	0.522 / 0.71	0.646 / 0.833
BERTScore-T5 <sub>large</sub>	0.529 / 0.74	0.53 / 0.721	0.532 / 0.749	0.531 / 0.74	0.5 / 0.699	0.485 / 0.69	0.643 / 0.831
BERTScore-DeBERT <sub>large</sub>	0.581 / 0.785	0.579 / 0.755	0.584 / 0.795	0.576 / 0.771	0.561 / 0.751	0.537 / 0.729	0.642 / 0.825
ParABLEU <sub>large</sub>	0.641 / 0.832	0.643 / 0.846	0.586 / 0.791	0.628 / 0.824	0.612 / 0.796	0.607 / 0.797	0.695 / 0.881
ParABLEU <sub>base</sub>	0.603 / 0.805	0.627 / 0.824	0.565 / 0.777	0.580 / 0.780	0.568 / 0.764	0.530 / 0.704	0.649 / 0.838

## E More generation examples

This section includes additional examples of one-shot paraphrase generation sampled from the MRPC dataset. See Section 5 for more information.

Table 8: Additional one-shot paraphrase generation results. ‘Ent. score’ is the predicted demonstration reference/candidate entailment probability from the classification loss. See Section 5 for more information.

#	Type	Demo reference	Demo candidate	Ent. score	Reference	Generated candidate
1	identity	Both bidders agreed to assume about \$ 90 million in debt owed on the planes	Both bidders agreed to assume about \$ 90 million in debt owed on the planes	90.0%	Robert Stewart, a spokesman for Park Place , the parent company of Caesars Palace , said he was surprised by the court 's decision .	Robert Stewart, a spokesman for Park Place, the parent company of Caesars Palace, said he was surprised by the court decision.
2	identity	Both bidders agreed to assume about \$ 90 million in debt owed on the planes	Both bidders agreed to assume about \$ 90 million in debt owed on the planes	90.0%	Democrats now hope to increase the value of awards proposed by Hatch and to create a mechanism to ensure the fund remains solvent .	Democrats now hope to increase the value of awards proposed by Hatch and create a mechanism to ensure that the fund remains solvent.
3	identity	Both bidders agreed to assume about \$ 90 million in debt owed on the planes	Both bidders agreed to assume about \$ 90 million in debt owed on the planes	90.0%	Indonesia 's army has often been accused of human rights abuses during GAM 's battle for independence , charges it has generally denied while accusing the separatists of committing rights violations .	Indonesia's army has often been accused of human rights abuses during GAM's battle for independence, charges it generally denied while accusing the separatists of committing rights violations.
4	identity	Both bidders agreed to assume about \$ 90 million in debt owed on the planes	Both bidders agreed to assume about \$ 90 million in debt owed on the planes	90.0%	Progress Software plans to acquire privately held DataDirect Technologies for about \$ 88 million in cash , the companies said Friday .	Progress Software plans to acquire privately held DataDirect Technologies for about \$ 88 million in cash, the companies said Friday.
5	identity	Both bidders agreed to assume about \$ 90 million in debt owed on the planes	Both bidders agreed to assume about \$ 90 million in debt owed on the planes	90.0%	A positive PSA test has to be followed up with a biopsy or other procedures before cancer can be confirmed .	A positive PSA test must be followed up with a biopsy or other procedures before cancer can be confirmed.

Table 9: Additional one-shot paraphrase generation results. ‘Ent. score’ is the predicted demonstration reference/candidate entailment probability from the classification loss.

#	Type	Demo reference	Demo candidate	Ent. score	Reference	Generated candidate
6	negation	Both bidders agreed to assume about \$ 90 million in debt owed on the planes	The bidders did not agree to assume about \$ 90 million in debt owed on the planes	0.0%	Robert Stewart , a spokesman for Park Place , the parent company of Caesars Palace , said he was surprised by the court's decision .	Robert Stewart, a spokesman for Park Place, the parent company of Caesars Palace, did not say he was surprised by the court decision.
7	negation	Both bidders agreed to assume about \$ 90 million in debt owed on the planes	The bidders did not agree to assume about \$ 90 million in debt owed on the planes	0.0%	Democrats now hope to increase the value of awards proposed by Hatch and to create a mechanism to ensure the fund remains solvent .	Democrats now hope to increase the value of awards proposed by Hatch and to create a mechanism to ensure the fund does not remain solvent.
8	negation	Both bidders agreed to assume about \$ 90 million in debt owed on the planes	The bidders did not agree to assume about \$ 90 million in debt owed on the planes	0.0%	Indonesia 's army has often been accused of human rights abuses during GAM 's battle for independence , charges it has generally denied while accusing the separatists of committing rights violations .	The Indonesian army did not often be accused of human rights abuses during the GAM struggle for independence, charges it generally denied while accusing the separatists of committing rights violations.
9	negation	Both bidders agreed to assume about \$ 90 million in debt owed on the planes	The bidders did not agree to assume about \$ 90 million in debt owed on the planes	0.0%	Progress Software plans to acquire privately held DataDirect Technologies for about \$ 88 million in cash , the companies said Friday .	Progress Software plans to acquire the privately held DataDirect Technologies for about \$ 88 million in cash, the companies said Friday.
10	negation	Both bidders agreed to assume about \$ 90 million in debt owed on the planes	The bidders did not agree to assume about \$ 90 million in debt owed on the planes	0.0%	A positive PSA test has to be followed up with a biopsy or other procedures before cancer can be confirmed .	A positive PSA test does not have to be followed up with biopsy or other procedures before cancer can be confirmed.

Table 10: Additional one-shot paraphrase generation results. ‘Ent. score’ is the predicted demonstration reference/candidate entailment probability from the classification loss.

#	Type	Demo reference	Demo candidate	Ent. score	Reference	Generated candidate
11	summary	Both bidders agreed to assume about \$ 90 million in debt owed on the planes	They agreed to assume the debt	99.0%	Robert Stewart , a spokesman for Park Place , the parent company of Caesars Palace , said he was surprised by the court 's decision .	They said they were surprised by the decision.
12	summary	Both bidders agreed to assume about \$ 90 million in debt owed on the planes	They agreed to assume the debt	99.0%	Democrats now hope to increase the value of awards proposed by Hatch and to create a mechanism to ensure the fund remains solvent .	They hoped to increase the value of the awards.
13	summary	Both bidders agreed to assume about \$ 90 million in debt owed on the planes	They agreed to assume the debt	99.0%	Indonesia 's army has often been accused of human rights abuses during GAM 's battle for independence , charges it has generally denied while accusing the separatists of committing rights violations .	They generally accused the separatists of committing rights violations.
14	summary	Both bidders agreed to assume about \$ 90 million in debt owed on the planes	They agreed to assume the debt	99.0%	Progress Software plans to acquire privately held DataDirect Technologies for about \$ 88 million in cash , the companies said Friday .	They wanted to acquire DataDirect Technologies.
15	summary	Both bidders agreed to assume about \$ 90 million in debt owed on the planes	They agreed to assume the debt	99.0%	A positive PSA test has to be followed up with a biopsy or other procedures before cancer can be confirmed .	They had to confirm the PSA test.

Table 11: Additional one-shot paraphrase generation results. ‘Ent. score’ is the predicted demonstration reference/candidate entailment probability from the classification loss.

#	Type	Demo reference	Demo candidate	Ent. score	Reference	Generated candidate
16	passive	Both bidders agreed to assume about \$ 90 million in debt owed on the planes	\$ 90 million in debt owed on the planes was agreed to be assumed by both bidders	99.8%	Robert Stewart , a spokesman for Park Place , the parent company of Caesars Palace , said he was surprised by the court ’ s decision .	The court decision was surprised by Robert Stewart , a spokesman for Park Place , the parent company of Caesars Palace , the parent company of Caesars Palace.
17	passive	Both bidders agreed to assume about \$ 90 million in debt owed on the planes	\$ 90 million in debt owed on the planes was agreed to be assumed by both bidders	99.8%	Democrats now hope to increase the value of awards proposed by Hatch and to create a mechanism to ensure the fund remains solvent .	The value of awards proposed by Hatch is now hoped to increase by Democrats and create a mechanism to ensure the fund remains solvent .
18	passive	Both bidders agreed to assume about \$ 90 million in debt owed on the planes	\$ 90 million in debt owed on the planes was agreed to be assumed by both bidders	99.8%	Indonesia ’ s army has often been accused of human rights abuses during GAM ’ s battle for independence , charges it has generally denied while accusing the separatists of committing rights violations .	Human rights abuses in the battle for independence was often accused by Indonesia ’ s army , charges which it generally denied while accusing the separatists of committing rights violations .
19	passive	Both bidders agreed to assume about \$ 90 million in debt owed on the planes	\$ 90 million in debt owed on the planes was agreed to be assumed by both bidders	99.8%	Progress Software plans to acquire privately held DataDirect Technologies for about \$ 88 million in cash , the companies said Friday .	The privately held DataDirect Technologies plans to be acquired by Progress Software in cash for about \$ 88 million .
20	passive	Both bidders agreed to assume about \$ 90 million in debt owed on the planes	\$ 90 million in debt owed on the planes was agreed to be assumed by both bidders	99.8%	A positive PSA test has to be followed up with a biopsy or other procedures before cancer can be confirmed .	A positive PSA test has to be confirmed in biopsy or other procedures before cancer can be followed up .

Table 12: Additional one-shot paraphrase generation results. ‘Ent. score’ is the predicted demonstration reference/candidate entailment probability from the classification loss.

#	Type	Demo reference	Demo candidate	Ent. score	Reference	Generated candidate
21	conditional perfect	Both bidders agreed to assume about \$ 90 million in debt owed on the planes	Both bidders would have agreed to assume about \$ 90 million in debt owed on the planes	69.4%	Robert Stewart, a spokesman for Park Place , the parent company of Caesars Palace , said he was surprised by the court’s decision .	Robert Stewart, a spokesman for Caesars Palace, the parent company of Park Place, said he would have been surprised by the court decision.
22	conditional perfect	Both bidders agreed to assume about \$ 90 million in debt owed on the planes	Both bidders would have agreed to assume about \$ 90 million in debt owed on the planes	69.4%	Democrats now hope to increase the value of awards proposed by Hatch and to create a mechanism to ensure that the fund remains solvent .	Democrats would now hope to increase the value of awards proposed by Hatch and to create a mechanism to ensure that the fund would remain solvent.
23	conditional perfect	Both bidders agreed to assume about \$ 90 million in debt owed on the planes	Both bidders would have agreed to assume about \$ 90 million in debt owed on the planes	69.4%	Indonesia ’s army has often been accused of human rights abuses during GAM ’s battle for independence , charges it has generally denied while accusing the separatists of committing rights violations .	Indonesia’s army would have often been accused of human rights abuses during GAM’s battle for independence, and to create a mechanism to ensure that the fund would remain solvent.
24	conditional perfect	Both bidders agreed to assume about \$ 90 million in debt owed on the planes	Both bidders would have agreed to assume about \$ 90 million in debt owed on the planes	69.4%	Progress Software plans to acquire privately held DataDirect Technologies for about \$ 88 million in cash , the companies said Friday .	Progress Software plans to acquire privately held DataDirect Technologies for about \$ 88 million in cash, the companies said Friday.
25	conditional perfect	Both bidders agreed to assume about \$ 90 million in debt owed on the planes	Both bidders would have agreed to assume about \$ 90 million in debt owed on the planes	69.4%	A positive PSA test would have to be followed up with a biopsy or other procedures before cancer can be confirmed .	A positive PSA test would have to be followed up with a biopsy or other procedures before cancer can be confirmed.

Table 13: Additional one-shot paraphrase generation results. ‘Ent. score’ is the predicted demonstration reference/candidate entailment probability from the classification loss.

#	Type	Demo reference	Demo candidate	Ent. score	Reference	Generated candidate
26	unrelated	Both bidders agreed to assume about \$ 90 million in debt owed on the planes	Bush turned out a statement yesterday thanking the commission for its work , and said " Our journey into space will go on . "	0.2%	Robert Stewart , a spokesman for Park Place , the parent company of Caesars Palace , said he was surprised by the court 's decision .	A spokesman for Caesars Palace, the parent company of Park Place, said he was surprised by the court's decision.
27	unrelated	Both bidders agreed to assume about \$ 90 million in debt owed on the planes	Bush turned out a statement yesterday thanking the commission for its work , and said " Our journey into space will go on . "	0.2%	Democrats now hope to increase the value of awards proposed by Hatch and to create a mechanism to ensure the fund remains solvent .	Hatch now hopes to raise the value of awards proposed by Democrats, and create a mechanism to ensure that the fund remains solvent.
28	unrelated	Both bidders agreed to assume about \$ 90 million in debt owed on the planes	Bush turned out a statement yesterday thanking the commission for its work , and said " Our journey into space will go on . "	0.2%	Indonesia 's army has often been accused of human rights abuses during GAM 's battle for independence , charges it has generally denied while accusing the separatists of committing rights violations .	During the battle for independence, GAM's army has often been accused of human rights abuses, which it generally denied, while accusing the separatists of committing rights violations.
29	unrelated	Both bidders agreed to assume about \$ 90 million in debt owed on the planes	Bush turned out a statement yesterday thanking the commission for its work , and said " Our journey into space will go on . "	0.2%	Progress Software plans to acquire privately held DataDirect Technologies for about \$ 88 million in cash , the companies said Friday .	The company announced on Friday that Progress Software plans to acquire DataDirect Technologies, for about \$ 88 million in cash.
30	unrelated	Both bidders agreed to assume about \$ 90 million in debt owed on the planes	Bush turned out a statement yesterday thanking the commission for its work , and said " Our journey into space will go on . "	0.2%	A positive PSA test has to be followed up with a biopsy or other procedures before cancer can be confirmed .	The report states that, before cancer can be confirmed, a positive PSA test must be carried out with a biopsy or other procedures.

Table 14: Additional one-shot paraphrase generation results. ‘Ent. score’ is the predicted demonstration reference/candidate entailment probability from the classification loss.

#	Type	Demo reference	Demo candidate	Ent. score	Reference	Generated candidate
31	personal	Both bidders agreed to assume about \$ 90 million in debt owed on the planes	We agreed to assume about \$ 90 million in debt owed on the planes	95.4%	Robert Stewart , a spokesman for Park Place , the parent company of Caesars Palace , said he was surprised by the court ’s decision .	Robert Stewart said he was surprised by the court decision.
32	personal	Both bidders agreed to assume about \$ 90 million in debt owed on the planes	We agreed to assume about \$ 90 million in debt owed on the planes	95.4%	Democrats now hope to increase the value of awards proposed by Hatch and to create a mechanism to ensure the fund remains solvent .	Democrats now hope to increase the value of awards proposed by Hatch.
33	personal	Both bidders agreed to assume about \$ 90 million in debt owed on the planes	We agreed to assume about \$ 90 million in debt owed on the planes	95.4%	Indonesia ’s army has often been accused of human rights abuses during GAM ’s battle for independence , charges it has generally denied while accusing the separatists of committing rights violations .	Indonesia has often been accused of human rights abuses during the battle for independence.
34	personal	Both bidders agreed to assume about \$ 90 million in debt owed on the planes	We agreed to assume about \$ 90 million in debt owed on the planes	95.4%	Progress Software plans to acquire privately held DataDirect Technologies for about \$ 88 million in cash , the companies said Friday .	We said we planned to acquire DataDirect Technologies privately for about \$ 88 million in cash.
35	personal	Both bidders agreed to assume about \$ 90 million in debt owed on the planes	We agreed to assume about \$ 90 million in debt owed on the planes	95.4%	A positive PSA test has to be followed up with a biopsy or other procedures before cancer can be confirmed .	We have to follow a positive PSA test with biopsy or other procedures before cancer can be confirmed.

Table 15: Additional one-shot paraphrase generation results. ‘Ent. score’ is the predicted demonstration reference/candidate entailment probability from the classification loss.

#	Type	Demo reference	Demo candidate	Ent. score	Reference	Generated candidate
36	informal	Both bidders agreed to assume about \$ 90 million in debt owed on the planes	Both guys agreed to take on like \$ 90 million that they owed for some planes	98.8%	Robert Stewart , a spokesman for Park Place , the parent company of Caesars Palace , said he was surprised by the court 's decision .	Robert Stewart, a spokesman for Park Place, the parent company of Caesars Palace, said they were surprised by the court decision.
37	informal	Both bidders agreed to assume about \$ 90 million in debt owed on the planes	Both guys agreed to take on like \$ 90 million that they owed for some planes	98.8%	Democrats now hope to increase the value of awards proposed by Hatch and to create a mechanism to ensure the fund remains solvent .	Democrats now hope to get the value of awards proposed by Hatch and create a mechanism to keep the fund solvent.
38	informal	Both bidders agreed to assume about \$ 90 million in debt owed on the planes	Both guys agreed to take on like \$ 90 million that they owed for some planes	98.8%	Indonesia 's army has often been accused of human rights abuses during GAM 's battle for independence , charges it has generally denied while accusing the separatists of committing rights violations .	Indonesia's army often got accused of human rights abuses at GAM's battle for independence, which they generally denied while accusing the separatists of committing rights violations.
39	informal	Both bidders agreed to assume about \$ 90 million in debt owed on the planes	Both guys agreed to take on like \$ 90 million that they owed for some planes	98.8%	Progress Software plans to acquire privately held DataDirect Technologies for about \$ 88 million in cash , the companies said Friday .	They wanted to buy some privately held DataDirect Technologies for about \$ 88 million in cash on Friday.
40	informal	Both bidders agreed to assume about \$ 90 million in debt owed on the planes	Both guys agreed to take on like \$ 90 million that they owed for some planes	98.8%	A positive PSA test has to be followed up with a biopsy or other procedures before cancer can be confirmed .	Some PSA tests need to be followed up with a biopsy or other procedures before they get cancer confirmed.