

NADIA: Diverse Paraphrase Generation with Syntax Planning and Diversity-driven Sequence Calibration

Anonymous ACL submission

Abstract

Paraphrase generation is a fundamental task in natural language processing. In this work, we study diverse paraphrase generation, and propose a novel method to increase surface-form diversity while maintaining semantic similarity for the generated paraphrase. Our method disentangles the generation into syntax structure planning and semantic realization, which first produces a syntax tree as high-level guidance and then generates surface form of paraphrase conditioned on the syntax tree. We further introduce a diversity-driven calibration loss to rank the probability of model generated sequences and enhance the output diversity. We evaluate our method on both ParaNMT dataset and a newly proposed DiverseQuora dataset, and our model outperforms strong baselines with better quality and diversity on both datasets.

1 Introduction

Paraphrase generation is an important task in natural language processing, with the goal to transform the source sentence into a different surface form while keeping the semantic meaning unchanged (Madnani and Dorr, 2010; Dou et al., 2022; Zhou and Bhat, 2021). It has various downstream applications such as question answering (Liu et al., 2020a), machine translation (Mallinson et al., 2017), and sentence simplification (Martin et al., 2022; Maddela et al., 2023).

While most studies in this domain focus on generating paraphrases with high semantic similarity, how to paraphrase with enhanced surface diversity is much less studied. Here, we define “*enhanced surface diversity*” as to generate sentences with largely different surface form compared with the original source input but still keep the semantic meaning unchanged. Surface-form diversity is an important feature for paraphrase generation because it helps to accommodate various audiences, contexts, and applications by generating multiple

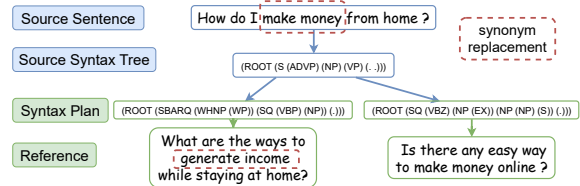


Figure 1: An example sentence and its paraphrases with different diversity. The syntax tree represents the surface-form organization of the target paraphrase.

ways to express the same idea. Diverse paraphrasing also ensures more robust and adaptable models, capable of understanding and producing a wider range of linguistic expressions. This can further benefit downstream applications by allowing more nuanced and varied outputs.

However, there remains challenges to generate diverse paraphrases with current token-level autoregressive language models. Achieving surface-form diversity while ensuring semantic fidelity to the input sentence is essential for effective paraphrase generation. Yet, the current training objective with maximum likelihood estimation (MLE) over each token in the target sequence cannot explicitly learn such disentanglement, and thus makes it hard to fulfill the aforementioned two objectives.

To overcome this, we propose a **Syntax Driven Diverse Paraphrase** framework (NADIA) to stably generate sentences with high diversity. First, syntactic structure is useful to represent the surface organization of a sentence, as shown in Figure 1. Thus, we explicitly incorporate syntactic structure as high-level guidance to control the surface-form generation and improve output diversity. Specifically, our model first produces a syntax tree as plan, and then conducts surface generation to produce paraphrases with synonym replacement conditioned on the syntax tree. By doing so, our model can effectively learn to disentangle paraphrase generation into syntax planning and semantic realization, thus generating more diverse outputs. Furthermore, to mitigate the issue of MLE training

073 that lacks sequence-level objective, we introduce
074 a diversity-driven calibration loss, which ranks
075 model outputs and aligns sequence-level likelihood
076 to both surface diversity and semantic similarity
077 in the latent space. Therefore, our model learns to
078 produce outputs with better diversity and quality.

079 To evaluate our model performance, we build
080 up a new dataset, DiverseQuora, with more di-
081 verse targets compared to the existing paraphrase
082 generation benchmarks.¹ Experiments on both
083 ParaNMT and our newly proposed DiverseQuora
084 dataset prove that NADIA with syntax planning and
085 diversity-driven sequence calibration outperforms
086 strong baselines with better quality and diversity.

087 2 Related Work

088 Paraphrase generation has received significant re-
089 search attention. Li et al., 2016 studied using mu-
090 tual information to generate more diverse responses.
091 Prakash et al., 2016 first used deep neural networks
092 to generate paraphrases. Wieting and Gimpel, 2018
093 (ParaNMT) and Kumar et al., 2020 (QQPos) built
094 up widely used datasets for paraphrase. Compared
095 to these datasets, DiverseQuora is more diverse
096 and has better quality distilled from the Large Lan-
097 guage Model. Most prior works in controllable
098 paraphrase rely on reinforcement learning(Gupta
099 et al., 2017, Li et al., 2018,Liu et al., 2020b), which
100 is difficult to train and control the diversity level.
101 Xu et al., 2018 studied using conditional embed-
102 ding to control diverse generation, and Cao and
103 Wan, 2020 studied extra loss in GAN to improve di-
104 versity. Similar to REAP(Goyal and Durrett, 2020)
105 and BRIO (Liu et al., 2022), we use ordering to
106 improve specific metrics of quality. Different from
107 their work, we incorporate ordering into planning-
108 based model to improve both surface diversity and
109 semantic fidelity. AESOP(Sun et al., 2021) and
110 GCPG(Yang et al., 2022) also use syntax informa-
111 tion to control generation, but they rely on human-
112 labeled exemplars. SGCP(Kumar et al., 2020) uses
113 a syntax tree but within a fixed human labeled set.
114 Both methods use contrastive loss to improve qual-
115 ity towards specific aspects.

116 3 Method

117 The overview of NADIA is shown in Figure 2. We
118 first describe our planning based model architec-
119 ture (§ 3.1), and then introduce the diversity-driven
120 calibration loss (§ 3.1).

¹Code and dataset will be released upon publication.

121 3.1 NADIA with Syntax-planning

122 Paraphrase generation is typically modeled as a
123 sequence-to-sequence (Seq2seq) task with the con-
124 ditional probability $P(y|x)$, where x denotes the
125 input and y is the target. In this work, we explicitly
126 incorporate target syntax feature z as high-level
127 guidance into the generation process: Instead of
128 directly generating surface target y , our model first
129 computes $p(z|x)$ to generate a syntactic plan that
130 represents the surface organization of the target,
131 and then produces the final target conditioned on
132 both input and plan with $p(y|z, x)$. In this way,
133 our planning-based modeling disentangles the syn-
134 tactic and semantic features and further improves
135 diversity with the guidance of the syntax.

136 Concretely, as shown in Figure 2, the encoder
137 takes the concatenation of the source sentence and
138 source syntactic parse as input. For the decoder,
139 instead of directly generating the target sentence, it
140 first predicts the target syntactic parse, and then pro-
141 duce the target paraphrase according to it. As the
142 generation of target depends on its syntactic plan,
143 we can manipulate the target by sampling plans
144 with desired attributes during inference, thereby
145 enabling the model to enhance output diversity.

146 3.2 Diversity-driven Calibration Loss

147 The typical training with MLE lacks sequence-level
148 objective and cannot directly optimize the model
149 towards the desired goal (Zhao et al., 2023). We
150 propose a diversity-driven calibration loss to pro-
151 vide sequence-level supervision and improve out-
152 put diversity and semantic fidelity. Following Liu
153 et al. 2022, we first train our model with the stan-
154 dard MLE objective. Then we sample multiple
155 candidates from the fine-tuned model and design
156 a multi-object calibration loss to align the model
157 towards the desired goal.

158 As the goal is to improve surface diversity while
159 maintaining semantic fidelity, we first design a
160 multi-objective based scoring function to score
161 each candidate:

$$162 \begin{aligned} S(\bar{y}) = & \lambda_{sts} \cdot STS(\bar{y}) + \lambda_b \cdot BS(\bar{y}) + \lambda_s \cdot SD(\bar{y}) \\ & - \lambda_{r1} \cdot R_1(\bar{y}) - \lambda_{r2} \cdot R_2(\bar{y}) \end{aligned} \quad (1)$$

163 where \bar{y} is the candidate, $STS(*)$ represents sen-
164 tence transformer similarity score (Reimers and
165 Gurevych, 2019) calculated based on an off-the-
166 shelf model², $BS(*)$ denotes BERT score, $SD(*)$
167

²<https://huggingface.co/sentence-transformers/all-MiniLM-L6-v2>

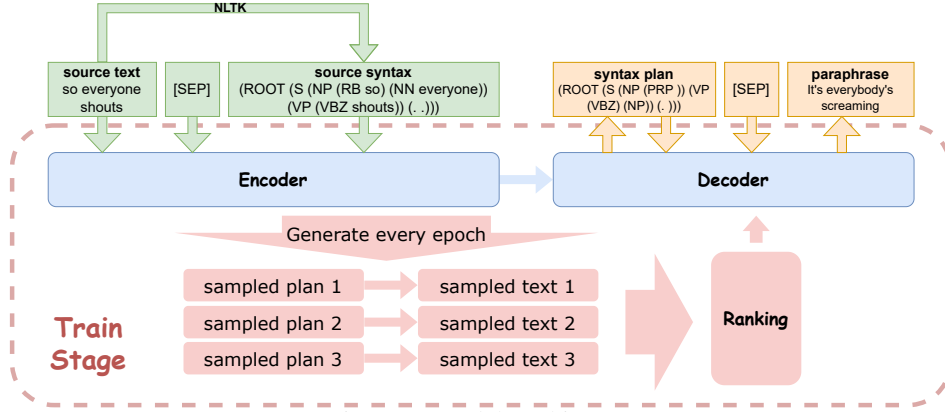


Figure 2: Model Architecture

is the syntax tree edit distance, and $R_1(*)$ and $R_2(*)$ stand for ROUGE-1/2. All scores are computed between \bar{y} and the source input x , and we omit x for simplicity. λ_* are weights of each score, tuned as hyper-parameters.

To align the model outputs with the desired object, we propose a ranking-based calibration loss to optimize the model to assign higher probability to candidates with higher scores:

$$L_{cal} = \sum \max(\log P(y_j|x) - \log P(y_i|x) + |j - i|\lambda_{cal}, 0) \quad (2)$$

where y_i and y_j are two sampled candidate paraphrases with $S(y_i) > S(y_j), \forall i, j$. λ_{cal} is chosen empirically to control the margin.

The final loss is a combination of both token-level cross-entropy (L_{ce}) and sequence-level calibration (L_{cal}): $L = L_{ce} + \alpha L_{cal}$, where α is the weight. These two losses are complementary to each other where cross entropy ensures the model not deviate significantly from the reference while the calibration loss coordinates the model for better diversity.

4 DiverseQuora Dataset

Existing work on paraphrase generation mainly adopt NLI based dataset such as Quora and ParaNMT and convert the original paraphrase to a generation task (Yang et al., 2022; Li et al., 2018). However, the target paraphrases in these datasets usually have a high surface-form similarity as the source sentences, making them less applicable in our scenario. We introduce a new dataset, DiverseQuora, for diverse paraphrase generation.

Specifically, we sample about 10K source sentences from Quora (Kumar et al., 2020) dataset, and prompt ChatGPT to produce a diverse paraphrase candidate. We then filter the low quality paraphrases by verifying their semantic similarity re-prompting ChatGPT. The detailed prompts are

| Dataset | Train | Val. | Test | Diversity |
|--------------|---------|-------|-------|-----------|
| DiverseQuora | 9,213 | 562 | 467 | 11.29 |
| Quora | 137,185 | 3,000 | 3,000 | 17.41 |
| ParaNMT | 493,081 | 500 | 800 | 18.53 |

Table 1: Statistics of DiverseQuora and existing paraphrase datasets. Diversity is measured by BLEU score between source and target, where lower score means better diversity.

in Appendix B.1. This yields 10,242 source-target pairs in total. To further validate the data quality, we randomly select 50 samples and manually check the sample quality, with the details in Appendix B.2. Finally, The data are split into train, validation and test sets, with the statistics reported in Table 1. As can be seen, the paraphrases in DiverseQuora have higher diversity compared to the existing datasets.

5 Experiments

5.1 Datasets

Following previous works (Yang et al., 2022), we include ParaNMT-small (Chen et al., 2019) which is a subset of ParaNMT-50M dataset (Wieting and Gimpel, 2018). We also evaluated our model on the DiverseQuora dataset.

5.2 Baseline

Besides Seq2Seq where we directly finetune an encoder-decoder Transformer as a baseline, we further include the following comparison methods. The implementation details of both our models and baselines are in Appendix A.1.

Control Seq2Seq. This is a diversity-controlled model. We categorize the dataset into five subsets based on the edit distance to indicate the diversity, and then prepend the diversity as control codes. During inference, we use the highest diversity as the hint to generate diverse outputs.

| ParaNMT | | | | | | |
|----------------------|--------------|---------------|--------------|--------------|--------------|--------------|
| Model | Diversity | | | | Similarity | |
| | R-2↓ | BLEU↓ | ED↑ | SD↑ | BERT↑ | STS↑ |
| Reference | 0.201 | 18.532 | 0.586 | 0.269 | 0.510 | 0.791 |
| Seq2Seq | 0.219 | 19.727 | 0.465 | 0.246 | 0.490 | 0.767 |
| Control Seq2Seq | 0.168 | 13.450 | 0.607 | 0.226 | 0.439 | 0.772 |
| Seq2Seq Post Scoring | 0.160 | 14.012 | 0.496 | 0.203 | 0.505 | 0.774 |
| NADIA | 0.102 | 8.015 | 0.558 | 0.224 | 0.436 | 0.741 |
| w/o Calibration Loss | 0.149 | 11.819 | 0.526 | 0.214 | 0.460 | 0.750 |
| w/o Plan | 0.131 | 12.025 | 0.533 | 0.218 | 0.455 | 0.750 |
| DiverseQuora | | | | | | |
| Seq2Seq | 0.396 | 23.997 | 0.498 | 0.187 | 0.643 | 0.886 |
| Control Seq2Seq | 0.327 | 18.274 | 0.605 | 0.200 | 0.598 | 0.861 |
| Seq2Seq Post Scoring | 0.307 | 16.694 | 0.563 | 0.224 | 0.619 | 0.865 |
| NADIA | 0.300 | 17.146 | 0.604 | 0.209 | 0.594 | 0.854 |

Table 2: Experimental Results. ED stands for Edit Distance, SD stands for Syntax Tree Edit Distance, R-2 stands for Rouge F-1 Score, BERT stands for BERT Score, STS stands for Semantic Textual Similarity.

| Input | Reference | NADIA |
|------------------------------|--|-------------------------------|
| It is your first own studio. | It 's the first studio you have owned. | You've got your first studio! |
| It 's a big risk for him. | The risks for him are big. | He's taking a great risk. |
| Relax. Take it easy. | Just calm down. | Calm down, buddy. |

Table 3: Three examples from NADIA output.

Seq2Seq Post Scoring. This is a post-scoring model, where we adopt Seq2Seq during training. In inference, we sample 8 outputs, and select the best one with the same scoring parameters as those used in the ordering loss, except that the syntax tree edit distance is replaced by edit distance.

5.3 Evaluation Metrics

We evaluate both diversity and semantic fidelity. For **surface diversity**, we adopt Rouge-2 and BLEU to measure the token overlap between source and generate paraphrase; we also include edit distance (ED) which calculates character level Hamming distance and syntax tree distance (SD) which computes tree edit distance (Zhang and Shasha, 1989) between the two syntax trees (keep only top 3 layers as in Figure 3) of source input and the generated paraphrase. For **semantic fidelity**, we leverage BERT score and Semantic Textual Similarity (all- MiniLM-L6-v2) to measure the similarity between source and paraphrase.

6 Results and Analysis

6.1 Automatic Results

Automatic metrics results are shown in Table 2. As can be seen, NADIA can generate outputs with both high surface-form diversity and semantic similarity, proving its effectiveness. Compared with Control Seq2Seq and Seq2Seq with Post Scoring, NADIA is able to achieve a good balance of both two objectives, demonstrating its effectiveness of generating diverse paraphrase. Notably, compared with vanilla Seq2Seq that not pursuing diversity, our

model has slightly lower BERT score and Semantic Textual Similarity. This is because changing the used words or word order inherently decreases these scores due to their order sensitivity. This is also evident in the scores between standard reference and input. Furthermore, after removing the planning or calibration loss, the results both drop, which show the effectiveness of the two components to jointly improve the model performance.

6.2 Case Study

We show sample outputs in Table 3 and Table 6. From the examples in Table 3, we can see our model output are more different to input sentence on syntax tree level (row 1 and row 2). Besides doing paraphrase, it is trying to leaking predicted information from pretrained model (row 3). This is because we use strategy to select less fine tuned checkpoint before combining with calibration loss.

7 Conclusion

We proposed a novel method to increase output diversity for the paraphrase task, which disentangles paraphrase generation into syntax planning and semantic realization. We further introduce a diversity-driven calibration loss to rank model generated outputs and enhance sequence-level diversity while maintaining semantic similarity. We propose a DiverseQuora dataset which is distilled from Large Language Model with diverse paraphrases. Experiments show that our model can generate both diverse and high-quality paraphrases compared to several strong baselines.

296 Limitations

297 Seeking diversity in paraphrase will intrinsically
298 decrease some similarity scores like BERT score
299 and Semantic Textual Similarity. Our model has
300 slightly lower similarity metric compared to base
301 seq2seq model. In the future, we will investigate
302 how to find better metrics which can evict this issue.
303 The ordering loss is hard to train on small dataset.
304 In the future, we seek to make it easier to control.

305 References

306 Yue Cao and Xiaojun Wan. 2020. [DivGAN: Towards di-](#)
307 [verse paraphrase generation via diversified generative](#)
308 [adversarial network](#). In *Findings of the Association*
309 [for Computational Linguistics: EMNLP 2020](#), pages
310 2411–2421, Online. Association for Computational
311 Linguistics.

312 Mingda Chen, Qingming Tang, Sam Wiseman, and
313 Kevin Gimpel. 2019. [Controllable paraphrase gener-](#)
314 [ation with a syntactic exemplar](#). In *Proceedings of*
315 [the 57th Annual Meeting of the Association for Com-](#)
316 [putational Linguistics](#), pages 5972–5984, Florence,
317 Italy. Association for Computational Linguistics.

318 Yao Dou, Chao Jiang, and Wei Xu. 2022. [Improv-](#)
319 [ing large-scale paraphrase acquisition and generation](#).
320 In *Proceedings of the 2022 Conference on Empiri-*
321 [cal Methods in Natural Language Processing](#), pages
322 9301–9323, Abu Dhabi, United Arab Emirates. As-
323 sociation for Computational Linguistics.

324 Tanya Goyal and Greg Durrett. 2020. [Neural syntactic](#)
325 [preordering for controlled paraphrase generation](#). In
326 *Proceedings of the 58th Annual Meeting of the Associ-*
327 [ation for Computational Linguistics](#), pages 238–252,
328 Online. Association for Computational Linguistics.

329 Ankush Gupta, Arvind Agarwal, Prawaan Singh, and
330 Piyush Rai. 2017. [A deep generative framework for](#)
331 [paraphrase generation](#).

332 Ashutosh Kumar, Kabir Ahuja, Raghuram Vadapalli,
333 and Partha Talukdar. 2020. [Syntax-guided controlled](#)
334 [generation of paraphrases](#). *Transactions of the Asso-*
335 [ciation for Computational Linguistics](#), 8:329–345.

336 Jiwei Li, Michel Galley, Chris Brockett, Jianfeng Gao,
337 and Bill Dolan. 2016. [A diversity-promoting objec-](#)
338 [tive function for neural conversation models](#).

339 Zichao Li, Xin Jiang, Lifeng Shang, and Hang Li. 2018.
340 [Paraphrase generation with deep reinforcement learn-](#)
341 [ing](#). In *Proceedings of the 2018 Conference on Em-*
342 [pirical Methods in Natural Language Processing](#),
343 pages 3865–3878, Brussels, Belgium. Association
344 for Computational Linguistics.

345 Dayiheng Liu, Jie Fu, Yidan Zhang, Chris Pal, and
346 Jiancheng Lv. 2020a. [Revision in continuous space:](#)
347 [Unsupervised text style transfer without adversarial](#)

[learning](#). *Proceedings of the AAAI Conference on*
348 [Artificial Intelligence](#), 34(05):8376–8383. 349

Mingtong Liu, Erguang Yang, Deyi Xiong, Yujie Zhang,
350 Yao Meng, Changjian Hu, Jinan Xu, and Yufeng
351 Chen. 2020b. [A learning-exploring method to gener-](#)
352 [ate diverse paraphrases with multi-objective deep](#)
353 [reinforcement learning](#). In *Proceedings of the 28th*
354 [International Conference on Computational Linguis-](#)
355 [tics](#), pages 2310–2321, Barcelona, Spain (Online).
356 International Committee on Computational Linguistics.
357 358

Yixin Liu, Pengfei Liu, Dragomir Radev, and Graham
359 Neubig. 2022. [Brio: Bringing order to abstractive](#)
360 [summarization](#). 361

Mounica Maddela, Yao Dou, David Heineman, and Wei
362 Xu. 2023. [LENS: A learnable evaluation metric for](#)
363 [text simplification](#). In *Proceedings of the 61st An-*
364 [nual Meeting of the Association for Computational](#)
365 [Linguistics \(Volume 1: Long Papers\)](#), pages 16383–
366 16408, Toronto, Canada. Association for Computa-
367 tional Linguistics. 368

Nitin Madnani and Bonnie J. Dorr. 2010. [Generat-](#)
369 [ing Phrasal and Sentential Paraphrases: A Survey](#)
370 [of Data-Driven Methods](#). *Computational Linguistics*,
371 36(3):341–387. 372

Jonathan Mallinson, Rico Sennrich, and Mirella Lapata.
373 2017. [Paraphrasing revisited with neural machine](#)
374 [translation](#). In *Proceedings of the 15th Conference of*
375 [the European Chapter of the Association for Compu-](#)
376 [tational Linguistics: Volume 1, Long Papers](#), pages
377 881–893, Valencia, Spain. Association for Computa-
378 tional Linguistics. 379

Louis Martin, Angela Fan, Éric de la Clergerie, Antoine
380 Bordes, and Benoît Sagot. 2022. [MUSS: Multilin-](#)
381 [gual unsupervised sentence simplification by mining](#)
382 [paraphrases](#). In *Proceedings of the Thirteenth Lan-*
383 [guage Resources and Evaluation Conference](#), pages
384 1651–1664, Marseille, France. European Language
385 Resources Association. 386

Aaditya Prakash, Sadid A. Hasan, Kathy Lee, Vivek
387 Datla, Ashequl Qadir, Joey Liu, and Oladimeji Farri.
388 2016. [Neural paraphrase generation with stacked](#)
389 [residual LSTM networks](#). In *Proceedings of COL-*
390 [LING 2016, the 26th International Conference on Com-](#)
391 [putational Linguistics: Technical Papers](#), pages 2923–
392 2934, Osaka, Japan. The COLING 2016 Organizing
393 Committee. 394

Nils Reimers and Iryna Gurevych. 2019. [Sentence-bert:](#)
395 [Sentence embeddings using siamese bert-networks](#). 396

Jiao Sun, Xuezhe Ma, and Nanyun Peng. 2021. [AESOP:](#)
397 [Paraphrase generation with adaptive syntactic control](#).
398 In *Proceedings of the 2021 Conference on Emperi-*
399 [cal Methods in Natural Language Processing](#), pages
400 5176–5189, Online and Punta Cana, Dominican Re-
401 public. Association for Computational Linguistics. 402

John Wieting and Kevin Gimpel. 2018. [ParaNMT-50M: Pushing the limits of paraphrastic sentence embeddings with millions of machine translations](#). In *Proceedings of the 56th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 451–462, Melbourne, Australia. Association for Computational Linguistics.

Qiongfai Xu, Juyan Zhang, Lizhen Qu, Lexing Xie, and Richard Nock. 2018. [D-page: Diverse paraphrase generation](#).

Kexin Yang, Dayiheng Liu, Wenqiang Lei, Baosong Yang, Haibo Zhang, Xue Zhao, Wenqing Yao, and Boxing Chen. 2022. [GCPG: A general framework for controllable paraphrase generation](#). In *Findings of the Association for Computational Linguistics: ACL 2022*, pages 4035–4047, Dublin, Ireland. Association for Computational Linguistics.

Kaizhong Zhang and Dennis Shasha. 1989. [Simple fast algorithms for the editing distance between trees and related problems](#). *SIAM J. Comput.*, 18:1245–1262.

Yao Zhao, Mikhail Khalman, Rishabh Joshi, Shashi Narayan, Mohammad Saleh, and Peter J Liu. 2023. [Calibrating sequence likelihood improves conditional language generation](#). In *The Eleventh International Conference on Learning Representations*.

Jianing Zhou and Suma Bhat. 2021. [Paraphrase generation: A survey of the state of the art](#). In *Proceedings of the 2021 Conference on Empirical Methods in Natural Language Processing*, pages 5075–5086, Online and Punta Cana, Dominican Republic. Association for Computational Linguistics.

A Experiment Details

A.1 Implementation Details

All models are instantiated by BART using base size. During inference, beam size is set to 5, and length penalty is set to 1.0. In training, all 8 samples are sampled with temperature=1.2, λ_{sts} , λ_b , λ_{r1} , λ_{r2} , λ_s , are set to 1.0, 0.333, 1.0, 0.0, 0.2 based on validation. α for L_{cal} is set to 1000. BERT score are calculated based on RoBERTa model, and Semantic Textual Similarity are calculated with “all-MiniLM-L6-v2”.

A.2 Explanation of Syntax Tree

We use NLTK³ to compute the syntax tree of both source and target sentences. For source syntax tree, we do not trim the tree, and concatenate the source sentence and the source syntax tree as input. For target, we trim the target syntax tree to $height = 3$ in our implementation. An concrete example is shown in Figure 3.

³<https://www.nltk.org/>

| Dataset | Avg. Quality | Std. Quality |
|--------------|--------------|--------------|
| DiverseQuora | 4.00 | 1.01 |
| Quora | 2.92 | 1.07 |

Table 4: Human evaluation of the paraphrase.

B More detailed of DiverseQuora Dataset

B.1 Prompt for DiverseQuora Construction

We leverage ChatGPT to produce diverse paraphrase given an input sentence. Concretely, we first prompt ChatGPT to produce a paraphrase candidate with the prompt:

"Given a sentence: `_input_`. Please rewrite the sentence. You need to keep the semantic meaning unchanged, while making the surface form different compared to the original sentence. You can use synonyms or/and change the sentence structure to make them different towards surface form."

To ensure the semantic similarity of the generated paraphrase, we verify the quality by prompting the ChatGPT again with the following prompt:

"sentence 1: `[_sent1_]`; sentence 2: `[_sent2_]`; Do sentence 1 and sentence 2 have the same semantic meaning? Answer "yes" or "no":"

If the candidate does not satisfy the above condition, we will repeat the process.

B.2 DiverseQuora Quality Evaluation

We further evaluated the quality of DiverseQuora and original Quora using the evaluation criteria described in Table 5. For each dataset, we randomly selected 50 examples from the training set, hide source information, merge and random shuffle them, and then evaluate them using the evaluation criteria described in Table 5. The results are shown in Table 4. Our dataset is of higher quality than the original Quora dataset (Kumar et al., 2020). We also manually selected some examples to show the difference between DiverseQuora and Quora in Table 7. Because the Quora dataset is generated by filtering negative examples from the original Quora Question Pairs dataset, some pairs are not good paraphrases but rather similar questions (such as row 9).

C More Examples Generated by NADIA

Here are some examples generated by NADIA. Through leveraging the power of BART model and Calibration Loss, we generate some examples with diversity and good quality.

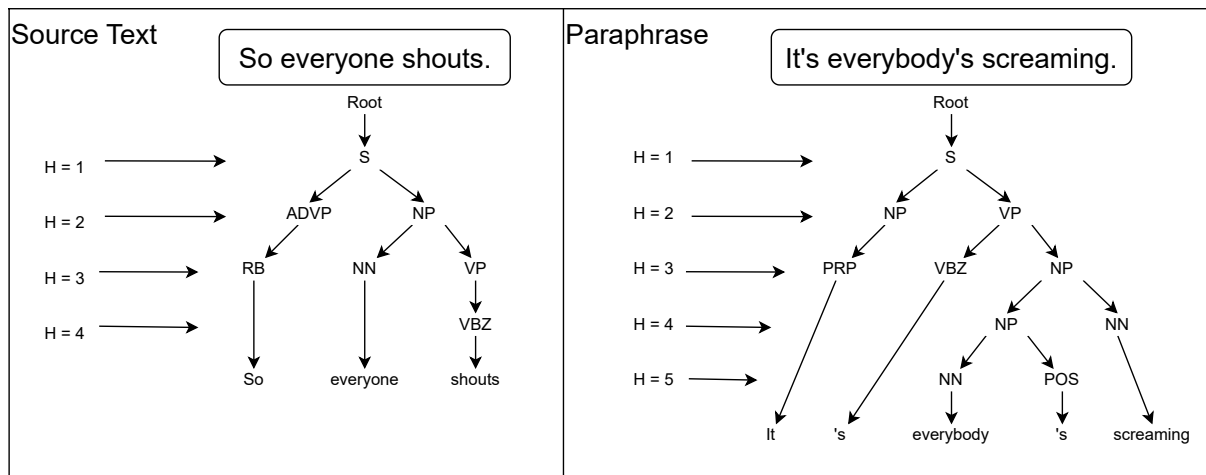


Figure 3: Showing the syntax tree examples. We select only top 3 layer in predicted syntax plan.

| Paraphrase Quality Scale | Definition |
|--------------------------|---|
| 5: Excellent | The paraphrase expresses the original meaning in a new and creative way, while still being accurate and fluent. It is clear that the paraphraser put a lot of thought into their work. |
| 4: Good Express | Expresses the original meaning accurately and fluently. It is similar to what a human would generate within a few seconds. |
| 3: Fair Generally | conveys the original meaning, but may be less fluent or original. It may contain some minor differences, such as removing unimportant information or adding well-known basic information. |
| 2: Poor | Does not accurately convey the original meaning. It may introduce new information, lack important information, use too many of the same words as the original sentence, or is poorly written. |
| 1: Very poor | Very similar to the original sentence, or expresses a very different meaning, or the paraphrase is difficult to understand. |

Table 5: The human evaluation template.

| Input | Reference | NADIA |
|---|---|---|
| It is your first own studio . | It's the first studio you have owned. | You've got your first studio! |
| The police think the bombing and today's gunfight in the apartments may be related. | The police suspect that the bombing may be tied to the apartment gunfight from earlier today. | The police think there was a connection between yesterday's explosion and today's shooting. |
| All this gold's gonna make things different. | Everything will be different with this gold. | The gold will change everything. |
| We will perform the opposite reaction to the plan he drew up . | He drew up the plan, so we 're going to perform the opposite reaction. | We're going to do a different reaction to this plan. |
| Let me show you to your seats. | I'll show you where you sit. | I'll show you the seat. |
| It's a big risk for him. | The risks for him are big. | He's taking a great risk. |
| Can't pass documents to this workspace. | Documents can not be forwarded to this workspace . | The document can not be transferred to the workspace. |
| What we want is to talk to your daddy. | We just want to talk to your dad. | We'd like to talk to your father. |
| Relax. Take it easy. | Just calm down. | Calm down, buddy. |
| Things have been getting a little weird around here. | There's some weird stuff going on. | It's getting kind of weird. |
| I have English , science , and chemistry books. | There are English books and science books and chemistry books. | I've got books in English, science and chemistry. |

Table 6: More examples generated by **NADIA**.

| Source | Quora | DiverseQuora |
|---|---|---|
| What are the benefits of a billing software? | What is the benefit of billing software? | What advantages does a billing software offer? |
| How do i get my likes and followers up on instagram? | How do i increase likes on instagram? | What strategies can I use to increase my likes and followers on Instagram? |
| How do i travel around the world without any money? | How can i travel without an id or money? | How can I journey around the globe without any funds? |
| Which laptop is best under 25000 inr? | Which is best laptop under 25000 with all features like vga and hdmi port? | What laptop is the optimal choice for under 25000 Indian rupees? |
| How do you take a screenshot on a mac laptop? | How do you take a screenshot on a mac? | What is the procedure to capture a screenshot on a mac laptop? |
| What happens if you actually drink bleach? | What happens when you swallow bleach? | What would occur if you ingested bleach? |
| How can i create a magnetic field? | How is a magnetic field created? | What steps do I need to take in order to generate a magnetic field? |
| Why was Hindi news channel NDTV India banned for one day? | What are your views on the government's decision of banning NDTV India for a day? | What was the reason for temporarily suspending the Hindi news channel NDTV India for one day? |
| Is the Aam Aadmi party's agenda economically comprehensive? | What is the agenda of Aam Aadmi party? | Does the Aam Aadmi Party's program possess an all-encompassing economic aspect? |
| How do i find ask someone to become a mentor? | How do i ask someone to be my mentor? | What would be the best way for me to request someone to be my mentor? |

Table 7: Examples from DiverseQuora.