

A Report on FCG GenChal 2022: Shared Task on Feedback Comment Generation for Language Learners

Ryo Nagata

Konan University, Japan

RIKEN, Japan

nagata-genchal@ml.hyogo-u.ac.jp.

Masato Hagiwara

Octanove Labs, USA*

masato@octanove.com

Kazuaki Hanawa

RIKEN, Japan[†]

k-hanawa@kodansha.co.jp

Masato Mita

RIKEN, Japan[‡]

mita_masato@cyberagent.co.jp

Abstract

We report on the results of the first ever shared task on *feedback comment generation for language learners* held as Generation Challenge (GenChal) in INLG 2022, which we call FCG GenChal. Feedback comment generation for language learners is a task where, given a text and a span, a system generates, for the span, an explanatory note that helps the writer (language learner) improve their writing skills. We show how well we can generate feedback comments with present techniques. We also shed light on the task properties and the difficulties in this task, with insights into the task including data development, evaluation, and comparisons of generation systems.

1 Introduction

Feedback comment generation for language learners is a task where, given a text and a span, a system generates, for the span, an explanatory note that helps the writer (language learners) improve their writing skills as exemplified in Fig. 1 (for convenience, the task will be abbreviated as *feedback comment generation*, hereafter). In this regard, feedback comment generation is related to grammatical error detection and correction. In many cases, however, it is not enough to just point out an error with its correct form in order to help language learners with writing learning. Instead, it is often essential for them to explain the underlying rules. In other words, it is essential in feedback comment generation to include more information than grammatical error detection and correction provide.

We report on the results of the first ever shared task on feedback comment generation held as Generation Challenge (GenChal) in INLG 2022, which we call FCG GenChal. One of the goals of this report is to reveal how well we can generate feedback comments with present techniques. There is a wide variety of choices for generation methods that are applicable to this task. Nevertheless, they have not yet been explored (at least, much less than in other generation tasks). Another goal is to shed a light on the task properties and the difficulties in this task. Specifically, we show, based on the results, insights into the task including data development, evaluation, and comparisons of generation systems.

2 Related Work

Generally speaking, feedback comment generation is a task of text-to-text generation. The input text, which is written by a language learner, is transformed into another text explaining the writing rules. This implies that generation methods employed in other generation tasks such as Machine Translation (MT) may be effective in the present task. For example, feedback comments often refer to words and phrases appearing in the input text, and techniques for referring to words in the source text (e.g., copy mechanisms) will likely be beneficial.

Feedback comment generation has its own unique aspects. It should be emphasized that a feedback comment is generated against a span (of the input text or sentence) whereas only a text (e.g., a sentence or utterance) is dealt with in other major text-to-text generation tasks such as MT and dialog systems. In consequence, feedback comment generation systems have to output different texts for the exact same source sentence, depending on the

*Currently also with Earth Species Project, USA

[†]Currently with KODANSHA LTD., Japan

[‡]Currently also with CyberAgent, Inc., Japan

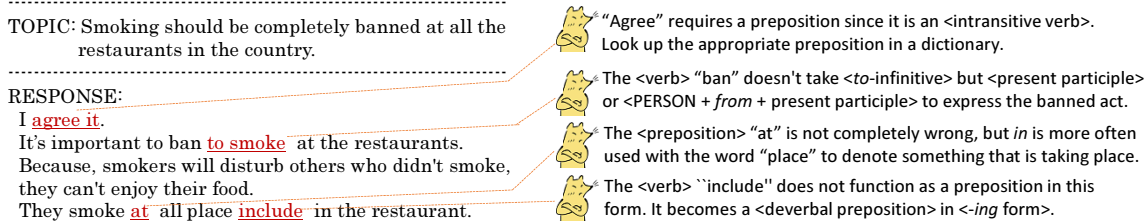


Figure 1: Example of Feedback Comments.

given spans.

The source and target languages are also unique. In this challenge, both are English, but there is room for discussion whether they fall into the same language class. The former is learner English, and inevitably it contains erroneous/unnatural words. Even within correct sentences, grammar, expressions, and style are expected to be used differently from canonical English. This brings out further research questions related to the source and target languages. For example, which is the best setting of vocabularies — only one common vocabulary for the source and target, or one for each? Does a pre-trained general (or native) language model work well to model learner English? There are a number of unaddressed research questions like these.

Feedback comment generation is also related to grammatical error detection/correction. The state-of-the-art methods typically solve the problems as sequence labeling (e.g., Kaneko et al. (2017)) or MT with DNNs (e.g., Junczys-Dowmunt et al. (2018); Napoles and Callison-Burch (2017); Rothe et al. (2021)). Recently, a DNN-based sequence labeling method is combined with symbolic transformations (Omelianchuk et al., 2020), which can be a good source of information to generate feedback comments.

Some researchers (Kakegawa et al., 2000; McCoy et al., 1996; Nagata et al., 2014) made an attempt to develop rule-based methods for diagnosing errors in line with grammatical error correction. However, this line of work suffered from the difficulty of improving coverage of errors.

More recently, researchers started to apply more modern techniques. Nagata (2019) showed that a neural-retrieval-based method was effective in preposition feedback comment generation. Lai and Chang (2019) proposed a method that used grammatical error correction and templates to generate detailed comments. Gkatzia et al. (2013) and

Gkatzia et al. (2014) proposed methods for automatically choosing feedback templates based on learning history. Hanawa et al. (2021) compared several neural-based generation methods with insights into feedback comment generation.

The availability of datasets for research in feedback comment generation has also been increasing. Nagata (2019) released a dataset consisting of feedback comments on preposition use. They marked up erroneous prepositions and annotated them with feedback comments. Nagata et al. (2020a) extended it to other grammatical errors and also other writing items such as discourse and lexical choice. Pilan et al. (2020) released a unique dataset where feedback comments on linking words were annotated.

3 Task Definition

3.1 General Definition

This subsection describes the general task definition of feedback comment generation, which is somewhat different from the one used in FCG Genchal. The task definition that was actually used is described in Subsec. 3.2, which is a reduced version of the general definition.

In the general task definition, a unit of the input in feedback comment generation consists of a text and spans of the text. Spans, which are counted by 1-based index based on characters, correspond to where to comment. An example input text would be:

(1) *I agree it.*

as shown on the left-hand side of Figure 1. A span would be 3 to 10, which will be abbreviated as 3:10, hereafter.

The output for a span is a string that explains why the span is not good, together with the underlying rule. To make the task different from grammatical error detection/correction, the output string has to

contain more information than what grammatical error detection/correction provide. In other words, just indicating the error position, the erroneous word(s), and/or the correct form are not enough as a valid feedback comment, details of which are discussed in Subsection 3.2.

3.2 Task Definition Used in FCG GenChal

The above task definition is too general and abstract to be a practical one. For this reason, we put some constraints on it.

First, the target language(s) can be any language, but we limit ourselves to English input texts and English feedback comments in this challenge. As shown in Figure 1, a feedback comment is typically made about erroneous, unnatural, or problematic words in a given text so that the writer can understand why the present form is not good together with the underlying rule.

Second, we limit the target only to errors related to preposition usages, as in the examples in Figure 1. It should be emphasized that the target preposition errors involve a much wider range of errors than in the conventional definition of preposition errors (such as the one provided by ERRANT (Bryant et al., 2017)). Examples include verb phrases used as a subject (e.g., **Lean English is difficult.*) and comparison between a phrase and a clause (e.g., **because an error → because of an error*); see the work (Nagata et al., 2020b) for the details.

Third, we also limit the input to a narrower unit. Specifically, the input text always consists of only one sentence with one span. Also, they are pre-tokenized where tokens are separated by whitespace. For example, the first sentence in Figure 1 would give an input:

(2) *I agree it . \t 3:10*

where \t stands for the tab character. If a sentence contains more than one preposition error, it appears two or more times with different spans.

Under these settings, participants develop a system that automatically generates an appropriate feedback comment in English for an input sentence and a span. The length of a generated feedback comment should be less than 100 tokens. If a system cannot generate an appropriate feedback comment for a given span, it may generate the special token <NO_COMMENT>, which is not counted as a system output. This allows us to calculate recall, precision, and F_1 , as explained below. An example output would be:

(3) *I agree it . \t 3:10* \t **“agree” is an intransitive verb and thus it requires a preposition before its object.**

Also note that the input sentence and its span are included in the system output for evaluation convenience.

Evaluation is probably the hardest challenge in this task. We adopt automated and manual evaluation methods. In the former, we simply take BLEU between a system output and its corresponding reference (manually created feedback comment)¹. In the latter, human evaluators examine whether a system output and its corresponding reference are equivalent in meaning. To be precise, a system output is regarded as appropriate if (1) it contains information similar to the reference and (2) it does not contain information that is irrelevant to the span; it may contain information that the reference does not contain as long as it is relevant to the span. This way of manual evaluation inevitably involves human subjectivity to some extent. In practice, however, the results of a pilot study show that inter-evaluator agreement is high.

The final manual evaluation measures are recall, precision, and F_1 . Recall is defined as the number of appropriate system outputs divided by the number of target spans. Similarly, precision is defined as the number of appropriate system outputs divided by the number of system outputs where the special output <NO_COMMENT> is excluded. F_1 is the harmonic mean of recall and precision.

We can do the same for BLEU. Simply, we replace the binary human judgment with the normalized, continuous BLEU value.

4 Data

Based on the work (Nagata, 2019; Nagata et al., 2020a), we created two versions of new datasets for this generation challenge: feedback comments written in the same language as the target (input) text (i.e., English) and in a different language (specifically, Japanese). The input texts (written by learners) are excerpts from the essays in ICNALE (Ishikawa, 2011). We had experts, who had experience in English teaching, manually annotate all preposition errors in the input texts with feedback comments in English and Japanese.

¹An official score is available at the FCG GenChal Official webpage: https://nagata-github.github.io/fcg_genchal/

Split	No. of feedback comments
Training	4,868
Development	170
Test	215

Table 1: Statistics on Datasets.

After having finished all annotations, we looked into the results. It turned out that the overall quality of the obtained data was much higher in the Japanese version than in the English. For this reason, we decided to use the Japanese version in this FCG GenChal; we translated the Japanese Feedback comments into English. Overall, it took us approximately three years to create the final datasets.

The results were split into training, development, and test sets. If a sentence contains more than one preposition error, it appears two or more times with different spans (in different lines). The split sets were provided for the participants, which are also available on the official FCG GenChal web site. Table 1 shows their statistics.

5 Participants and Results

5.1 Timeline and Summary of FCG GenChal

As shown in Figure 2, we initially had 12 registrations from seven countries. After registration, we released the training and development sets on 28 January, 2022. We let the participants have approximately four months to prepare their system.

After four months, we released the test set on 2 May, 2022. The participants had one week to prepare their generation results for final submission. In the end, seven teams submitted their results². Four out of the seven systems are available on the

²Probably, feedback comment generation is a relatively new task and we guess four months were not enough for some teams to develop their systems.

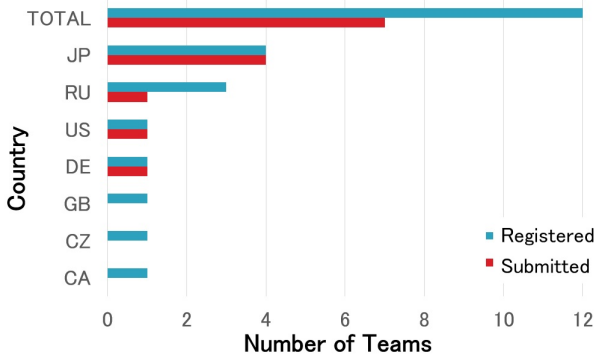


Figure 2: Statistics on Registration and Submission.

official FCG GenChal website³. Table 2 summarizes the seven systems. Also, a short description of each system is as follows:

ihmana: This system consists of three modules: retrieval, masking, and generation modules. The first module retrieves the instance most similar to the input learner sentence from the training data. Then, the second module masks tokens in the retrieved feedback comment that do not fit the input sentence well. Finally, the generation module generates a feedback comment given the input sentence and the retrieved, masked feedback comment. The retrieval and masking modules are based on BERT while the generation module uses a pre-trained T5 (Raffel et al., 2020). This system is capable of generating <NO_COMMENT>.

nigula: This system is based on a pre-trained T5. The generator is trained on the official training set and also on an extra set obtained by data augmentation. Data augmentation is done by completing clipped input learner sentences using a language model. This system is capable of generating <NO_COMMENT>.

TMUUED: This system is also based on a pre-trained T5. It takes the Part-Of-Speech labels for the input learner sentence as an extra source of information. It also uses a synonym dictionary to determine if the generation result is appropriate or not. This system is capable of generating <NO_COMMENT>.

kjimichi: This system also uses a pre-trained T5 as a generator. It also uses RoBERTa (Liu et al., 2020) as a classifier to obtain grammatical term labels such as noun and preposition. The predicted grammatical term labels are used as an additional source of information in the T5 generator. This system is not capable of generating <NO_COMMENT>.

shotakoyama: This system is based on GPT-2. Its approach is unique compared to the other systems in that it focuses on cleaning training data rather than improving the generation module itself. Specifically, it automatically corrects errors in feedback comment annotation such as incorrect spans. Also, it uses error type tags obtained via GECToR (Omelianchuk et al., 2020) as an extra source of information. This system is capable of generating <NO_COMMENT>.

stahl: This system uses BART (Lewis et al., 2020) as a generator. It is also unique in that only

³https://nagata-github.github.io/fcg_genchal/

Participant ID	Generator	Other Information
ihmana	T5 (t5-base)	Retrieve and masking modules: BERT (bert-base-cased)
nigula	T5 (t5-large)	Data augmentation: T5 (t5-large)
TMUUED	T5 (t5-base)	NLTK to obtain POS tags.
GU	T5 (t5-large)	Data augmentation: EleutherAI/gpt-neo-1.3B
kjimichi	T5 (t5-base)	Grammar term prediction: RoBERTa (roberta-large)
shotakoyama	GPT-2 (gpt2-large)	Data cleaning, error correction operation tags (GECToR)
stahl	BART	Clustering of training instances (k -means clustering)
Baseline	BiLSTM	—

Table 2: Summary of Participating Systems.

it exploits clustering. Specifically, before training, feedback comments in the training set are automatically grouped by clustering where TF-IDF vectors are used. This system is not capable of generating `<NO_COMMENT>`.

We ourselves implemented a baseline system for comparison. It was a text generation system based on a BiLSTM-based encoder-decoder with a copy mechanism (Hanawa et al., 2021). It is also available on the official website.

We initially had two months for manual evaluation. It actually took us approximately one month to evaluate the results of all systems including our baseline system. It took some more time to double check the evaluation results and to perform related tasks such as summarizing the results. We released the results on 25 June, 2022 as planned.

5.2 Results

Table 3 and Table 4 show the manual and automatic evaluation results, respectively. Both tables show a similar overall tendency. However, if we look at the details, we can see differences between them.

In BLEU-based evaluation, the system rankings are reversed compared to manual evaluation in some cases. This means that we cannot use BLEU to obtain strict system rankings as in shared tasks. We will get back to this point in Sect. 6.

In manual evaluation, the performance values tend to be larger than the corresponding automatic evaluation values. This suggests that even if n -gram overlap rate is not so high between a generated feedback and its reference, it can be judged to be appropriate by human evaluators. In other words, a feedback comment can be described by different words and phrases as expected.

Participant ID	Precision	Recall	$F_{1.0}$
ihmana	0.6244	0.6186	0.6215
nigula	0.6093	0.6093	0.6093
TMUUED	0.6132	0.6047	0.6089
GU	0.5860	0.5860	0.5860
kjimichi	0.5628	0.5628	0.5628
shotakoyama	0.5756	0.5488	0.5619
stahl	0.3581	0.3581	0.3581
Baseline	0.3116	0.3116	0.3116

Table 3: Results of Manual Evaluation.

Participant ID	Precision	Recall	$F_{1.0}$
ihmana	0.486	0.482	0.484
TMUUED	0.477	0.471	0.474
GU	0.471	0.471	0.471
nigula	0.463	0.463	0.463
kjimichi	0.460	0.460	0.460
stahl	0.437	0.437	0.437
shotakoyama	0.444	0.424	0.434
Baseline	0.334	0.334	0.334

Table 4: Results of Automatic Evaluation (BLEU).

6 Discussion

As shown in Sect. 5, all participating systems are based on a pre-trained, transformer-based generator while the baseline system uses a non-pre-trained BiLSTM. This partially answers one of the research questions raised in Sect. 2 (i.e., Does a pre-trained general (native) language model work on learner writings?). The results show that pre-training and/or the architectures (likely both) contribute to performance improvement, although we need more investigation to confirm this argument.

The top five systems use T5 as a generator while the rest use either GPT-2 or BART. The results prefer T5 as a generator for feedback comment

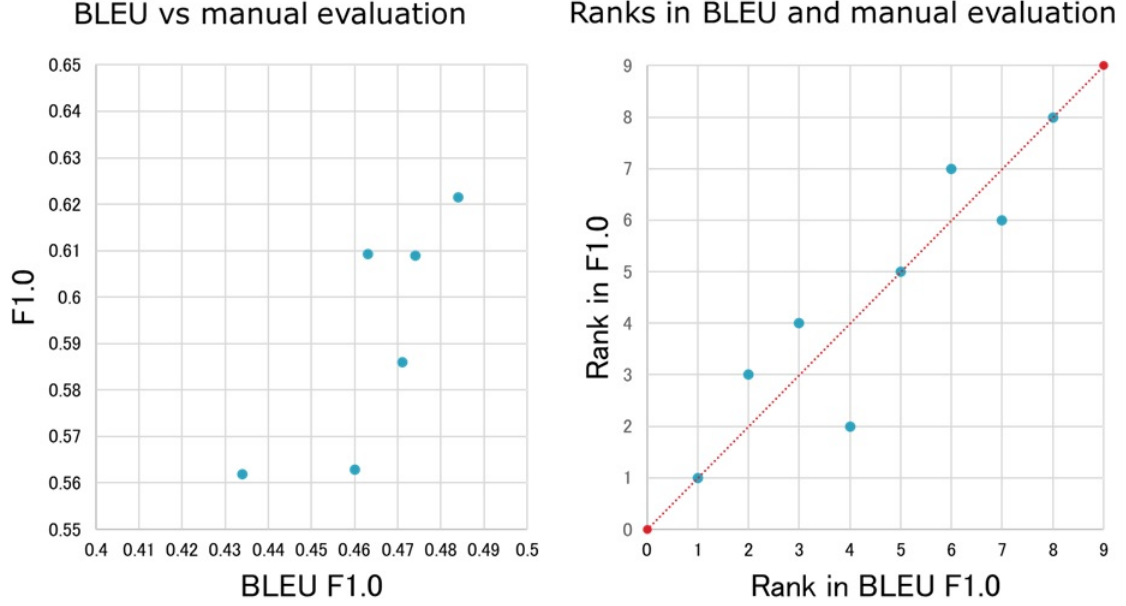


Figure 3: Comparison between Automatic (BLEU-based $F_{1.0}$) and Manual ($F_{1.0}$) Evaluation Results.

generation. Having said that, we need more investigations to confirm this argument, considering the amount of training, development, and test data.

Interestingly, some teams report that smaller models perform equal to, or even better than, the corresponding larger model (see their reports for the details). For example, they report that under the same condition, ‘t5-base’ achieves a better BLEU-based $F_{1.0}$ than ‘t5-large’ while ‘bart-base’ and ‘bart-large’ achieve a similar BLEU-based $F_{1.0}$. A possible reason for this is that the training set is not so large and that the amount is not enough to fine-tune a large model properly. Here, one thing we should note is that these comparisons are based on BLEU because manual evaluation was applied only to the final submission results (thus, one has to use automatic evaluation (e.g., BLEU) to compare their system variants). Manual evaluation may lead to a different conclusion.

Comparison between BLEU and manual evaluation results provide an interesting insight into this task, which is summarized in Fig. 3. BLEU and manual evaluation results correlate well (correlation coefficient: 0.85). However, the system rankings differ from those by manual evaluation when the difference in BLEU is small. Specifically, even if the difference is more than 0.01 (i.e., TMUUED: 0.474 vs. nigula: 0.463), a reversal of a system ranking occurs. According to the obtained results, when the difference is larger than a certain

value (e.g., 0.02 in this case), BLEU might be a reliable measure to choose a better system (or a better method, or a better hyper-parameter setting). We need more investigations to confirm that this argument is correct. For the time-being, we do not have enough data to do so and we need manual evaluation to obtain reliable system rankings. At the same time, manual evaluation is costly and time-consuming. One of the necessary research directions is to explore more efficient ways of evaluation.

Another challenging direction is to pursue methods for generating `<NO_COMMENT>` (i.e., *not possible to generate a reliable feedback comment*). Considering practical use, it is important to decide not to generate when the system is not confident enough. In FCG GenChal, four out of the seven systems are capable of generating `<NO_COMMENT>`. Their implementations are rather simple (e.g., simple rule-based) and their effects are limited; the difference between precision and recall is rather small as shown in Table 3.

7 Conclusions

In this paper, we have reported on the results of a new generation challenge called *feedback comment generation for language learners*. The best-performing system achieves an $F_{1.0}$ of 0.62 in manual evaluation. The results suggest that pre-

training and/or transformer-based methods are effective. They also suggest that smaller models within transformer-based methods perform better with the training data available. We have also reported insights into automatic and manual evaluation in feedback comment generation.

Acknowledgements

We thank all participants for their efforts. Thanks to their work, we now have a unique testbed including datasets, system outputs, and system source codes. Without their contributions, we could not have held this generation challenge. We also thank the GenChal and INLG organizers for their support. This work was supported by JSPS KAKENHI Grant-in-Aid for Scientific Research(C) Grant Number JP22K12326 and Japan Science and Technology Agency (JST), PRESTO Grant Number JPMJPR1758, Japan. This work was partly conducted by using computational resource of AI Bridging Cloud Infrastructure (ABCI) provided by National Institute of Advanced Industrial Science and Technology (AIST).

References

- Christopher Bryant, Mariano Felice, and Ted Briscoe. 2017. [Automatic annotation and evaluation of error types for grammatical error correction](#). In *Proceedings of the 55th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 793–805, Vancouver, Canada. Association for Computational Linguistics.
- Dimitra Gkatzia, Helen Hastie, Srinivasan Janarthnam, and Oliver Lemon. 2013. Generating student feedback from time-series data using reinforcement learning. In *Proc. of 14th European Workshop on Natural Language Generation*, pages 115–124.
- Dimitra Gkatzia, Helen Hastie, and Oliver Lemon. 2014. Comparing multi-label classification with reinforcement learning for summarisation of time-series data. In *Proc. of 52nd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 1231–1240.
- Kazuaki Hanawa, Ryo Nagata, and Kentaro Inui. 2021. [Exploring methods for generating feedback comments for writing learning](#). In *Proceedings of the 2021 Conference on Empirical Methods in Natural Language Processing*, pages 9719–9730, Online and Punta Cana, Dominican Republic. Association for Computational Linguistics.
- Shinichiro Ishikawa. 2011. *A new horizon in learner corpus studies: The aim of the ICNALE project*, pages 3–11. University of Strathclyde Publishing, Glasgow.
- Marcin Junczys-Dowmunt, Roman Grundkiewicz, Shubha Guha, and Kenneth Heafield. 2018. Approaching neural grammatical error correction as a low-resource machine translation task. In *Proc. of 2018 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long Papers)*, pages 595–606.
- Jun’ichi Kakegawa, Hisayuki Kanda, Eitaro Fujioka, Makoto Itami, and Kohji Itoh. 2000. Diagnostic processing of Japanese for computer-assisted second language learning. In *Proc. of 38th Annual Meeting of the Association for Computational Linguistics*, pages 537–546.
- Masahiro Kaneko, Yuya Sakaizawa, and Mamoru Komachi. 2017. Grammatical error detection using error- and grammaticality-specific word embeddings. In *Proc. of 8th International Joint Conference on Natural Language Processing (Volume 1: Long Papers)*, pages 40–48.
- Yi-Huei Lai and Jason Chang. 2019. TellMeWhy: Learning to explain corrective feedback for second language learners. In *Proc. of 2019 Conference on Empirical Methods in Natural Language Processing and 9th International Joint Conference on Natural Language Processing (EMNLP-IJCNLP): System Demonstrations*, pages 235–240.
- Mike Lewis, Yinhan Liu, Naman Goyal, Marjan Ghazvininejad, Abdelrahman Mohamed, Omer Levy, Veselin Stoyanov, and Luke Zettlemoyer. 2020. [BART: Denoising sequence-to-sequence pre-training for natural language generation, translation, and comprehension](#). In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, pages 7871–7880, Online. Association for Computational Linguistics.
- Yinhan Liu, Myle Ott, Naman Goyal, Jingfei Du, Mandar Joshi, Danqi Chen, Omer Levy, Mike Lewis, Luke Zettlemoyer, and Veselin Stoyanov. 2020. [Ro{bert}a: A robustly optimized {bert} pretraining approach](#).
- Kathleen F. McCoy, Christopher A. Pennington, and Linda Z. Suri. 1996. English error correction: A syntactic user model based on principled “mal-rule” scoring. In *Proc. of 5th International Conference on User Modeling*, pages 69–66.
- Ryo Nagata. 2019. Toward a task of feedback comment generation for writing learning. In *Proc. of 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing*, pages 3197–3206.
- Ryo Nagata, Kentaro Inui, and Shin’ichiro Ishikawa. 2020a. Creating Corpora for Research in Feedback Comment Generation. In *Proc. of the 12th Language Resources and Evaluation Conference*, pages 340–345.

- Ryo Nagata, Kentaro Inui, and Shin'ichiro Ishikawa. 2020b. [Creating corpora for research in feedback comment generation](#). In *Proceedings of the Twelfth Language Resources and Evaluation Conference*, pages 340–345, Marseille, France. European Language Resources Association.
- Ryo Nagata, Mikko Vilenius, and Edward Whittaker. 2014. Correcting preposition errors in learner English using error case frames and feedback messages. In *Proc. of 52nd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 754–764.
- Courtney Napoles and Chris Callison-Burch. 2017. Systematically adapting machine translation for grammatical error correction. In *Proc. of 12th Workshop on Innovative Use of NLP for Building Educational Applications*, pages 345–356.
- Kostiantyn Omelianchuk, Vitaliy Atrasevych, Artem Chernodub, and Oleksandr Skurzhashnyi. 2020. [GECToR – grammatical error correction: Tag, not rewrite](#). In *Proceedings of the Fifteenth Workshop on Innovative Use of NLP for Building Educational Applications*, pages 163–170, Seattle, WA, USA → Online. Association for Computational Linguistics.
- Ildiko Pilan, John Lee, Chak Yan Yeung, and Jonathan Webster. 2020. A Dataset for Investigating the Impact of Feedback on Student Revision Outcome. In *Proc. of 12th Language Resources and Evaluation Conference*, pages 332–339.
- Colin Raffel, Noam Shazeer, Adam Roberts, Katherine Lee, Sharan Narang, Michael Matena, Yanqi Zhou, Wei Li, and Peter J. Liu. 2020. Exploring the limits of transfer learning with a unified text-to-text transformer. *Journal of Machine Learning Research*, 21(1).
- Sascha Rothe, Jonathan Mallinson, Eric Malmi, Sebastian Krause, and Aliaksei Severyn. 2021. [A Simple Recipe for Multilingual Grammatical Error Correction](#). In *Proc. of 59th Annual Meeting of the Association for Computational Linguistics and 11th International Joint Conference on Natural Language Processing*, pages 702–707.