

CxGGEC: Construction-Guided Grammatical Error Correction

Anonymous ACL submission

Abstract

The grammatical error correction (GEC) task aims to detect and correct grammatical errors in text to enhance its accuracy and readability. Current GEC methods primarily rely on grammatical labels for syntactic information, often overlooking the inherent usage patterns of language. In this work, we explore the potential of construction grammar (CxG) to improve GEC by leveraging constructions to capture underlying language patterns and guide corrections. We first establish a comprehensive construction inventory from corpora. Next, we introduce a construction prediction model that identifies potential constructions in ungrammatical sentences using a noise-tolerant language model. Finally, we train a CxGGEC model on construction-masked parallel data, which performs GEC by decoding construction tokens into their original forms and correcting erroneous tokens. Extensive experiments on English and Chinese GEC benchmarks demonstrate the effectiveness of our approach.

1 Introduction

Grammatical Error Correction (GEC) is the task of automatically detecting and correcting errors in text (Bryant et al., 2023), which after the advent of Transformer (Vaswani et al., 2017), has been categorized into two main types: Seq2Edit method and Seq2Seq method (Sun et al., 2021; Zhang et al., 2022b).

Seq2Edit method typically involves converting source sentences into a sequence of edit operations (Stahlberg and Kumar, 2020; Omelianchuk et al., 2020), which offers specific advantages in the GEC task due to its higher inference efficiency, while limited to manually selecting dictionaries (Awasthi et al., 2019; Malmi et al., 2019). Seq2Seq method treats GEC as a monolingual translation problem (Junczys-Dowmunt et al., 2018a; Sun et al., 2021) and demonstrates a better correction ability. Recent advances have enabled language models (LMs)

to more adequately capture syntactic phenomena (Jawahar et al., 2019; Wei et al., 2022), making them capable GEC systems when little or no data is available (Zhang et al., 2022b). However, because the use of syntactic information of prior works is limited to the application of grammatical labels, we observe that currently no method can fully leverage the syntactic information and semantic usage patterns inherent to perform the GEC task.

Construction Grammar (CxG) (Goldberg, 1995, 2003) regards constructions (i.e., form-meaning pairs) as the fundamental units of linguistic knowledge, with each construction modeled as a sequence of slot-constraints (Dunn, 2017). For example, “Subject-Verb–Object1–Object2” is a ditransitive construction (Goldberg, 1995) that represents the abstract meaning of transferring. CxG claims that our knowledge of language is captured by network of constructions (Goldberg, 2003). Grammatical errors stem from a lack of sufficient knowledge about language usage (Bryant et al., 2023), making constructions beneficial for enhancing the GEC task. Some examples are demonstrated in Table 1, which shows the improvements of the GEC task by identifying potential constructions in the sentence.

Based on the above observation, we propose the following technical approach: (1) establishing a construction inventory from corpora, (2) identifying constructions from ungrammatical sentences, and (3) training models using ungrammatical sentences augmented with constructions for the GEC task.

However, realizing the above approach presents the following three challenges:

- (Q1) What types of constructions are most effective in improving the performance of the GEC task?
- (Q2) How can constructions be identified from ungrammatical sentences?
- (Q3) How can the identified constructions be effec-

Ungrammatical Sentence	Identified Construction	Corrected Sentence
The book which I bought it yesterday is very interesting.	DET-NOUN-PRON-SUBJ-VERB	The book which I bought yesterday is very interesting.
The students in the library preparing for their exams.	DET-NOUN-ADP-DET-NOUN-AUX	The students in the library are preparing for their exams.
Some important departments need strict administration of their members.	VBP-ADJ-NOUN-ADP	Some important departments need strict administration for their members.

Table 1: Examples of three error types demonstrating the improvement of the GEC task using CxG: unnecessary, missing, replacement. (DET, NOUN, PRON, SUBJ, VERB, ADP, AUX, ADJ, and VBP denote determiner, noun, pronoun, subject, verb, preposition, auxiliary verb, adjective, and non-3rd person singular present verb, respectively.)

tively utilized to guide the GEC task?

As for (Q1), an observation is that the guiding effectiveness of constructions is maximized when they overlap with or are adjacent to grammatical errors in sentences. Current methods for construction extraction can be categorized into manual extraction and automatic extraction (Xu et al., 2023). Manual extraction is limited by scale. Two primary automatic methods exist: one calculates bidirectional association scores between adjacent words (Dunn, 2017), while the other, CxGLearner (Xu et al., 2024), leverages LM token prediction probabilities. The former produces shorter constructions with limited structural completeness due to adjacent calculation method, whereas the latter, using LMs, generates more complete constructions with well-distributed lengths because it allows extended distances when assessing slot constraints. Thus, we adopt CxGLearner for constructing the construction inventory.

Regarding (Q2), current construction generation methods are only applicable to grammatical sentences. Inspired by Jiang et al. (2021) that LMs are insensitive to subtle differences between sequences, which means LMs exhibit a certain degree of tolerance toward noise, we propose a LM-based approach to identify expected constructions from ungrammatical sentences.

To answer (Q3), we train a CxGGEC model based on a construction-augmented vocabulary. Through concatenating ungrammatical sentences with responding construction-masked sentences, CxGGEC is able to decode constructions into correct tokens by the Seq2Seq method.

Extensive experiments have been conducted to illustrate the superiority of CxGGEC on the GEC task, while multilingual experiments further indicate construction is beneficial across languages.

2 System Overview

Our CxGGEC framework can be divided into three steps: (1) construction generation, (2) construction masking, (3) CxG-guided GEC. Figure 1 displays the entire framework.

2.1 Construction Generation

Construction Inventory Establishment. Construction is represented as a sequence of slot-constraints. We annotate the part-of-speech tags in corpus from various domains, and employ CxGLearner (Xu et al., 2024) to extract constructions from annotated corpus, which assesses the association strength among slots based on LM. Therefore, we establish a well-distributed construction inventory, which will be taken as a construction vocabulary in subsequent training phase.

Identifying Construction in Ungrammatical Sentences. Because ungrammatical sentences may damage constructions, the construction inventory we obtained cannot be applied to identify expected constructions from ungrammatical sentences. Therefore, based on the tolerance of LMs for noise, we leverage the construction inventory to train a construction prediction model to identify constructions from ungrammatical sentences. The training details of the prediction model are demonstrated in Section 3.

2.2 Construction Masking

To guide the GEC task with CxG, firstly we identify expected constructions from ungrammatical sentences through the construction prediction model, and then we construct the parallel training corpus by concatenation of the ungrammatical sentences with their construction-masked counterparts and

the corresponding ground-truth sentences. Finally, we train a CxGGEC model by the Seq2Seq method.

2.3 CxG-guided GEC

For inference process, we concatenate the ungrammatical sentences with their construction-masked versions, forming a combined input just as the training phase. Specifically, construction masking serves as a context-aware signal that directs the model to locate parts requiring correction and output grammatical sentences by decoding construction tokens into original tokens and decoding error tokens into correct tokens. Through this construction-guided approach, the model aligns the grammatical error with the language usage patterns inherent in constructions, thereby improving the effects on GEC tasks.

3 Model

3.1 Construction Prediction Model

Construction Selection Strategy. Since constructions are often stored redundantly at different levels of abstractness, overlapping constructions can be captured by the grammar induction algorithm (Dunn, 2017, 2019). Xu et al. (2024) summarize the phenomenon of overlap into two scenarios: *Inclusion* and *Intersection*, which can lead to issues like redundancy and imbalanced encoding.

Based on our Seq2Seq training approach, it is essential to ensure that the constructions used to mask within the training sentences do not exhibit overlap or intersection. Drawing inspiration from RoBERTa’s (Liu, 2019) dynamic masking approach, we randomly retain the overlapping sections for each sentence, while keeping the other parts intact. This method prevents overlaps and allows the model to learn diverse combinations of constructions, helping to mitigate the risk of the construction prediction model overfitting to specific construction patterns. The algorithm is depicted in Algorithm 1. CHECKOVERLAP(\cdot) inspects whether a given construction c overlaps with any constructions in the set \mathcal{C} , returning a boolean value. We RANDOMKEEP(\cdot) resolves conflicts by stochastically retaining either c or the conflicting construction in \mathcal{C} . ADD(\cdot) appends non-overlapping constructions c to \mathcal{C} . This process is iteratively applied to all constructions in \mathcal{C} . The algorithm generates N sets of optimized constructions, $\mathcal{S} = \{\mathcal{C}_1, \mathcal{C}_2, \dots, \mathcal{C}_N\}$, by applying the dynamic masking strategy N times. Finally, \mathcal{S} cap-

Algorithm 1: Dynamic Masking for Multiple Construction Schemes

Input: The set of all constructions \mathcal{C} . Number of schemes N .

Output: A set of construction schemes $\mathcal{S} = \{\mathcal{C}_1, \mathcal{C}_2, \dots, \mathcal{C}_N\}$.

```

1  $\mathcal{S} \leftarrow \{\}$ 
2 for  $i \in \{1, 2, \dots, N\}$  do
3    $\mathcal{C}_i \leftarrow \text{INITIALIZE}()$ 
4   foreach construction  $c \in \mathcal{C}$  do
5     if CHECKOVERLAP( $c, \mathcal{C}_i$ ) then
6       RANDOMKEEP( $c, \mathcal{C}_i$ )
7     end
8     else
9       ADD( $c, \mathcal{C}_i$ )
10    end
11  end
12   $\mathcal{S} \leftarrow \mathcal{S} \cup \{\mathcal{C}_i\}$ 
13 end
14 return  $\mathcal{S}$ 

```

tures diverse valid construction schemes.

Input and Output Definition. For a given grammatical sentence S_c , constructions are extracted to produce a masked sentence S_m :

$$S_m = f_c(S_c, \mathcal{C}), \quad (1)$$

where $f_c(\cdot)$ handles dynamic construction masking.

Training. The Seq2Seq model learns the mapping:

$$\hat{S}_m = \text{Seq2Seq}(S_c), \quad (2)$$

optimizing the difference between \hat{S}_m and the target S_m .

Inference. During inference, the model inputs a sentence S , applies construction-based masking similarly, and outputs a CxG-masked sentence \hat{S}_m by aligning them with learned construction patterns:

$$\hat{S}_m = \text{Seq2Seq}(S) \quad (3)$$

3.2 CxGGEC Model

In this section, we present the training of CxGGEC models and the construction-guided GEC process. Our method includes three key steps: extending the vocabulary with constructions, preparing construction-masked parallel training data, and pre-training the model with the parallel data.

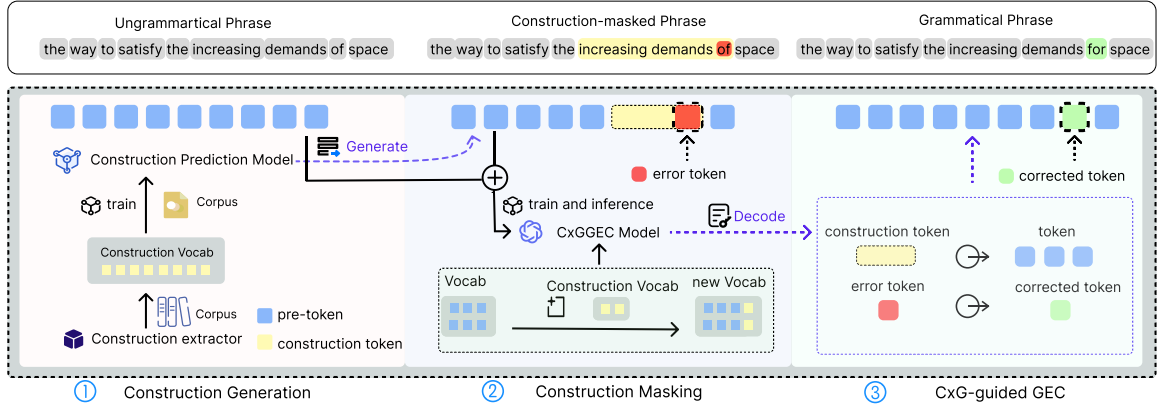


Figure 1: Overview of the proposed CxGGEC framework.

Construction Augmented Vocabulary. To integrate constructions into LMs, we explicitly extend their input vocabularies. Let \mathcal{C} denote the set of all constructions extracted during preprocessing. Each construction $c_i \in \mathcal{C}$ is treated as a new token and added to the existing vocabulary \mathcal{V} . The updated vocabulary is denoted as $\mathcal{V}' = \mathcal{V} \cup \mathcal{C}$.

For the vocabulary extension, the embedding matrix $\mathbf{E} \in \mathbb{R}^{|\mathcal{V}| \times d}$, where d is the embedding dimension, is updated to $\mathbf{E}' \in \mathbb{R}^{|\mathcal{V}'| \times d}$. All added construction embeddings are initialized randomly and fine-tuned during training. Specifically, for each construction c_i , its embedding is defined as:

$$\mathbf{e}_{c_i} = \text{Initialize}(\text{rand}(\mathbf{e}); \forall c_i \in \mathcal{C}), \quad (4)$$

where $\text{rand}(\mathbf{e})$ generates random values sampled from a uniform distribution over $[-\sqrt{d}, \sqrt{d}]$.

Construction-Augmented Input Representation.

To better leverage multiple construction predictions during training, we modify the input representation by concatenating the ungrammatical sentence \mathbf{x}_{ug} with its masking-augmented sentences generated by construction prediction model.

Let $\{\mathbf{m}_1, \mathbf{m}_2, \dots, \mathbf{m}_T\}$ denote the set of masked sentences generated by applying construction prediction model T times to \mathbf{x}_{ug} . The augmented input \mathbf{x}' is then defined as:

$$\mathbf{x}' = \mathbf{x}_{ug} \oplus \mathbf{m}_1 \oplus \mathbf{m}_2 \oplus \dots \oplus \mathbf{m}_T, \quad (5)$$

where \oplus denotes sequence concatenation.

The inclusion of multiple masked sentences allows the model to benefit from diverse masking strategies and improves generalization.

The corresponding target sentence \mathbf{y} is the standard grammatical correction for \mathbf{x}_{ug} . The parallel

training pair is defined as $\langle \mathbf{x}', \mathbf{y} \rangle$, where \mathbf{x}' is the construction-augmented input and \mathbf{y} is the grammatical ground truth. This process generates a construction-augmented parallel corpus.

Pretraining with Construction-Augmented Examples.

The pretraining phase uses the construction-augmented parallel corpus. The model's objective is to minimize the negative log-likelihood of the target sequence \mathbf{y} conditioned on the input \mathbf{x}' . Formally, the loss function is defined as:

$$\mathcal{L} = - \sum_{i=1}^N \sum_{t=1}^T \log P(y_i^t | \mathbf{x}'_i, y_i^{<t}; \Theta), \quad (6)$$

where y_i^t is the token at timestep t in the target sequence \mathbf{y}_i , T is the length of \mathbf{y}_i , and Θ are the model parameters. The probability $P(y_i^t | \cdot)$ is computed via the decoder's autoregressive output during training.

Pre-trained embeddings for vocabulary tokens remain initialized using the original model weights, while the embeddings for newly added construction tokens are learned adaptively.

4 Experiments

4.1 Experiments Setup

Datasets and Evaluation. For the English, we use the clean version of the original Lang-8 corpus (Mizumoto et al., 2011; Tajiri et al., 2012) as train sets. Specifically for the model based on Bart-Large model (Lewis et al., 2020), we use the W&I+LOCNESS train-set (Bryant et al., 2019) for model fine-tuning following Zhang et al. (2022b). Following Zhang et al. (2022b), Li et al. (2023)

and Li and Wang (2024), we use BEA-Dev (Bryant et al., 2019) as the development dataset, and use BEA-Test set and CoNLL14-Test set (Ng et al., 2014) as test datasets. For Chinese, following Li and Wang (2024), the models are fine-tuned on the Chinese Lang8 dataset (Zhao et al., 2018) and the HSK dataset (Zhang, 2009), and on the FCGEC training set (Xu et al., 2022) respectively. The models are evaluated on MuCGEC (Zhang et al., 2022a) and FCGEC test sets. For English evaluation, following Yuan et al. (2021a), we use ER-RANT and M^2 (Dahlmeier and Ng, 2012) to evaluate GEC models on BEA-Test set and CoNLL14-Test set, respectively. For Chinese experiments, following Li and Wang (2024), models are evaluated on MuCGEC and FCGEC test sets using ChERRANT (Zhang et al., 2022a; Xu et al., 2022). Precision, recall, and $F_{0.5}$ values are reported metrics for all the experiments. Dataset details are listed in Appendix A.

Implementation. We train construction prediction model based on the BART-Base model (Lewis et al., 2020). For English GEC models, we train models based on the BART-Large (Lewis et al., 2020) and T5-Large (Raffel et al., 2020) models. Specifically, for the model based on the BART-Large, we refer to the training strategy of Zhang et al. (2022b). For the T5-Large model, we adopt the training strategy of Li et al. (2023). Both take Fairseq (Ott et al., 2019) as training framework. Due to the absence of a Chinese version of the T5 model, the experiments conducted in Chinese do not incorporate the use of the T5 model. For creating Chinese construction inventory, we use Python library jieba (Feng, 2012) for sentence segmentation and part-of-speech tagging.

Baselines. (1) GECToR (Omelianchuk et al., 2020) represents the Seq2Edit models. (2) BART (Lewis et al., 2020) and T5 (Raffel et al., 2020) are backbones of Seq2Seq GEC methods. (3) SynGEC (Zhang et al., 2022b) incorporates syntactic information into the BART model. (4) Multi-Encoder (Yuan et al., 2021b) encodes error categories as auxiliary information. (5) GEC-DePend (Yakovlev et al., 2023) integrates error detection with correction by the MLM. (6) TemplateGEC (Li et al., 2023) uses the output of the GECToR model as supplementary information for Seq2Seq models. (7) DeCoGLM (Li and Wang, 2024) promotes performance of the GEC model by combining detection and collection tasks to mutually boost each other.

The performance of GECToR and BART model on the Chinese dataset is reported by Li and Wang (2024), and the results for BART on the English dataset are reported by Zhang et al. (2022b).

4.2 Main Results

The main results of our experiments are listed in Table 2. It can be observed that our CxGGEC models achieve comparable performance across various benchmarks. Our framework demonstrates improvements across all benchmarks compared to the BART and T5 backbones. We achieve better performance than existing methods on the CoNLL14-Test set and FCGEC-Test set. The results show the effectiveness of our framework. Notably, our model based on the T5 backbone outperforms BART due to the basic idea of Raffel et al. (2020) to treat every text processing problem as a “text-to-text” problem, which can easily adapt to different inputs.

CxGGEC performs well on both English and Chinese GEC tasks, showcasing its generalizability in error correction across these two major languages. Compared with SynGEC, our method achieves further improvement on English datasets with less parameters added (13M), highlighting that constructions, as sets of slots, encode more semantic and syntactic information than only grammatical labels. This enables the model to achieve a deeper understanding of language usage and further enhances its GEC performance.

4.3 Analysis Study

Analysis on construction length. To explore the impact of construction length on the performance of GEC tasks, we apply two distinct methods to establish the construction inventory to support CxGGEC. First is the method of grammarinduction algorithm (Dunn, 2017), we refer to it as GIA for simplicity. The second method is CxGLearner (Xu et al., 2024).

The construction length distribution displayed in Figure 2 originates from the construction inventory covered in the CLang8-train dataset, a widely-used dataset for GEC models to align with the distribution patterns of sentences in English. The average construction length generated by GIA is approximately 3.0, while the constructions generated by CxGLearner exhibit a higher average length of 4.1. Notably, the lengths produced by CxGLearner exhibit a more balanced distribution. As shown in Table 3, the constructions generated by CxGLearner provide more significant guidance for the LM GEC

Method	Parameters	English						Chinese					
		CoNLL-14 <i>test</i>			BEA-19 <i>test</i>			MuCGEC <i>test</i>			FCGEC <i>test</i>		
		P	R	F _{0.5}	P	R	F _{0.5}	P	R	F _{0.5}	P	R	F _{0.5}
GECToR	110M	77.5	40.1	65.3	79.2	53.9	72.4	46.72	27.14	40.83	46.11	34.35	43.16
BART	400M	73.6	48.6	66.7	74.0	64.9	72.0	41.90	29.48	38.64	38.38	37.62	38.23
T5	770M	-	-	66.1	-	-	72.1	-	-	-	-	-	-
SynGEC	110M+400M	74.7	49.0	67.6	75.1	65.5	72.9	54.69	29.10	46.51	-	-	-
Multi-Encoder	110M+107M	71.3	44.3	63.5	73.3	61.5	70.6	-	-	-	-	-	-
GEC-DePenD	253M	73.2	37.8	61.6	72.9	53.2	67.9	-	-	-	-	-	-
TemplateGEC	125M+770M	74.8	50.0	68.1	76.8	64.8	74.1	-	-	-	-	-	-
DeCoGLM	335M	75.1	49.4	68.0	77.4	64.6	74.4	45.01	31.77	41.55	55.75	37.91	50.96
CxGGEC (Bart-large)	13M+400M	73.8	50.5	67.6	74.8	65.3	72.7	47.90	29.94	42.78	59.90	35.92	52.84
CxGGEC (T5-large)	13M+770M	74.9	50.7	68.3	75.7	65.8	73.5	-	-	-	-	-	-

Table 2: Results on English and Chinese GEC benchmarks. The highest metric is indicated in bold.

Strategy	BEA-19			CoNLL-14		
	P	R	F _{0.5}	P	R	F _{0.5}
GIA	73.7	50.2	67.4	74.0	65.2	72.1
CxGLearner	74.9	50.7	68.3	75.7	65.8	73.5

Table 3: Performance of CxGGEC (T5-large) with different construction inventory establishing strategies on BEA-19 test and CoNLL-14 test benchmarks.

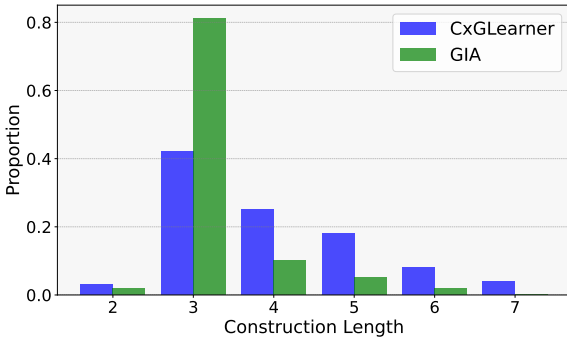


Figure 2: Length distribution of construction inventories extracted from GIA (Dunn, 2017) and CxGLearner (Xu et al., 2024).

task compared to GIA.

This observation implies that CxGLearner achieves more comprehensive coverage of constructions inherent in corpus. While both methods generate useful constructions for GEC, constructions extracted with GIA tend to be relatively short or incomplete, because GIA is prone to truncate the constructions too early. This result indicates that long and well-distributed constructions tend to perform better on GEC tasks, because they align with the usage patterns in the corpus and contain more knowledge of language usage.

Analysis on Construction Coverage. To reveal how construction coverage contributes to GEC

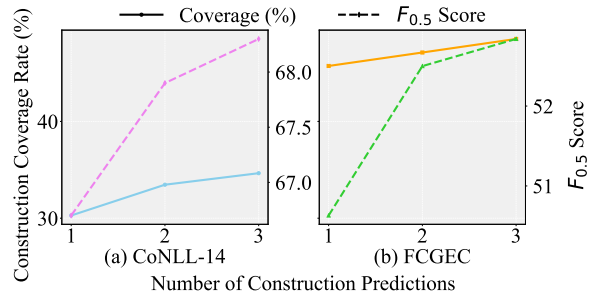


Figure 3: Construction coverage rate and $F_{0.5}$ score across prediction steps.

Strategy	BEA-19			CoNLL-14		
	P	R	F _{0.5}	P	R	F _{0.5}
CxGGEC	74.9	50.7	68.3	75.7	65.8	73.5
w/o DM	73.5	49.8	67.1	74.1	66.9	72.5

Table 4: Comparison of the performance of CxGGEC (T5-large) with and without dynamic masking.

tasks, we perform experiments on number of construction predictions in Figure 3. We observe a gradual improvement in GEC performance as the number of predictions increases. The construction coverage rate is defined as the ratio of the number of sentences of which the constructions identified cover the error positions to the total number of sentences. The result shows that increasing construction predictions enhances the model’s ability to cover sentence errors effectively, and therefore improve the overall performance of GEC tasks.

Analysis on Construction Masking Strategy. To figure out the impact of dynamic masking strategy on GEC tasks, we analyze the results of the GEC task without dynamic masking strategy compared to results of CxGGEC in Table 4. We refer to dynamic masking as DM for simplicity. The results

Type	Baseline			CxGGEC		
	P	R	F _{0.5}	P	R	F _{0.5}
M	72.4	65.8	71.0	74.2	70.0	73.3
R	72.0	60.6	69.4	73.4	63.0	71.1
U	75.0	69.5	73.8	75.2	70.6	74.2

Table 5: Results of error types in BEA-Test. Baseline is the T5-Large model. (M, R, and U stand for missing, replacement, and unnecessary errors, respectively.)

demonstrate that DM yields superior model performance compared to fixed masking. This can be attributed to the ability of DM to prevent construction prediction model from overfitting to specific masking patterns and to enhance the model’s capacity to adapt to diverse contexts.

Analysis on Error Types. To reveal what types of error can CxG guide GEC tasks better, we compare results of error types on the BEA-Test benchmark in Figure 5. The baseline is T5-large model and the CxGGEC model is based on T5-Large model. M, R, and U stand for missing, replacement, and unnecessary errors, respectively. Overall, CxGGEC demonstrates higher performance on three error types, particularly in missing and replacement errors but achieves subtle improvement in unnecessary errors. The potential reason is that constructions identified by the prediction model may fail to include unnecessary errors. This requires the model to expend effort on error detection and correction, thereby resulting in only subtle improvement.

Analysis on POS Tags. We intend to explore the impact of part-of-speech (POS) tags on the BEA-Test dataset. UPOS stands for Universal POS tags and XPOS stands for Language-Specific POS tags. We compare the results of using only UPOS, using only XPOS, and combining the two with a specified proportion during training construction prediction model to evaluate their effectiveness. As shown in Table 6, using only UPOS performs slightly worse than using only XPOS, because XPOS is better at capturing fine-grained grammatical and structural information. The combination of UPOS and XPOS yields better results because adding a certain proportion of UPOS provides high-level abstraction that aids in capturing generalized linguistic patterns. This combination enables the model to balance generalization and specificity, ultimately enhancing its overall performance.

UPOS	XPOS	BEA-19			CoNLL-14		
		P	R	F _{0.5}	P	R	F _{0.5}
✗	✗	69.2	48.4	66.5	71.1	47.7	65.1
✓	✗	71.4	63.2	69.6	73.3	47.4	66.1
✗	✓	72.8	64.4	70.9	74.5	48.7	67.4
✓	✓	75.7	65.8	73.5	74.9	50.7	68.3

Table 6: Results of POS tags.

Analysis on Visualization. To explain why CxG can effectively guide GEC tasks from the perspective of language models, we compare the attention matrices of a baseline LM (Bart-Large) and CxGGEC model based on Bart-Large model in Figure 4. Tokens identified as constructions (construction-masked segments) are highlighted in red, while the shaded area further emphasizes the attention on these tokens. The result shows that attention of CxGGEC model focuses around phrases, especially those involving constructions (highlighted parts). This reflects the ability of the CxGGEC model to incorporate constructional information from constructions, guiding the model to focus on meaningful sections of the sentence rather than isolated tokens. This allows CxGGEC to better interpret the overall context, particularly in ungrammatical sentences, where individual tokens may not provide sufficient information.

5 Related Works

GEC Methods. Two widely used approaches in GEC are Seq2Edit and Seq2Seq. In Seq2Edit methods, Seq2Edits (Stahlberg and Kumar, 2020) predicts a sequence of span-level edit operations applied to the source text, while GECToR (Omelianchuk et al., 2020) extends traditional operations with custom transformations, such as suffix changes and token merging. The advantage of the Seq2Edit approach is its faster speed compared to Seq2Seq. However, a key limitation is its reliance on manually curated editing operations, which can reduce transferability and fluency (Li et al., 2022). Seq2Seq models (Lewis et al., 2020; Raffel et al., 2020) have demonstrated high performance in GEC (Junczys-Dowmunt et al., 2018b; Choe et al., 2019; Zhao et al., 2019; Katsumata and Komachi, 2020), though their inference efficiency is lower compared to Seq2Edit. Mallinson et al. (2020) and Yakovlev et al. (2023) utilize Masked Language Models (Kenton and Toutanova, 2019) to generate corrections, aiming to benefit from self-supervised pretraining. Previous studies have also

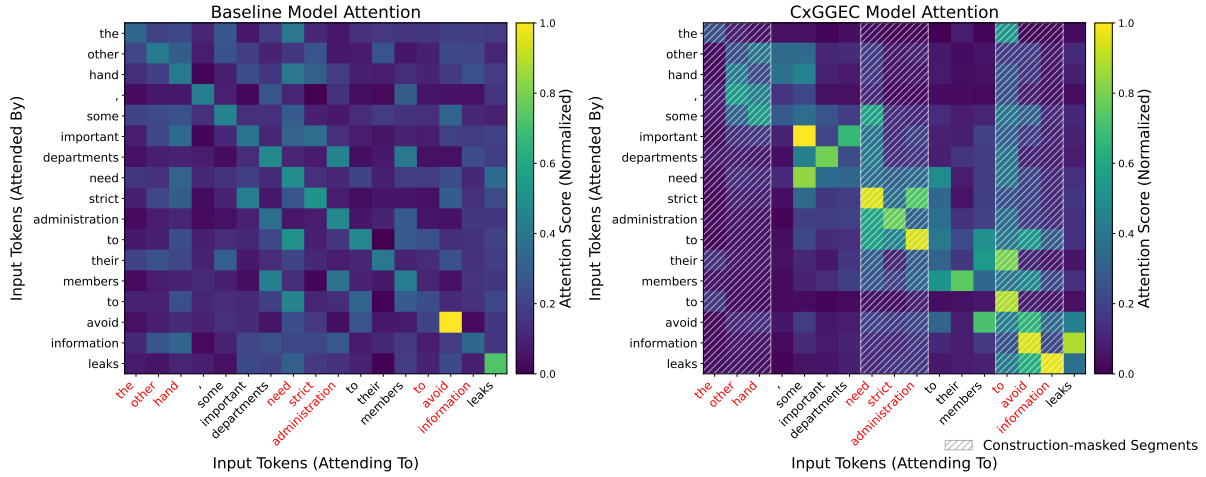


Figure 4: Comparison of attention maps based on Bart-Large and CxGGEC (Bart-Large).

505 incorporated error detection results (e.g., detection
 506 labels from a Seq2Edit model) as auxiliary informa-
 507 tion to enhance GEC performance (Kaneko et al.,
 508 2020; Yuan et al., 2021b; Li et al., 2023). State-
 509 of-the-art models further incorporate syntactic in-
 510 formation to improve performance. For example,
 511 SynGEC (Zhang et al., 2022b) integrates depen-
 512 dency syntax into GEC models, while CSynGEC
 513 (Zhang and Li, 2022) enhances GEC tasks by lever-
 514 aging constituent-based syntax. However, current
 515 methods rely on grammatical labels for syntactic
 516 information, failing to fully capture the structural
 517 and semantic usage patterns of a language. There-
 518 fore, we introduce construction grammar to address
 519 the issue.

520 **Applications of CxG in NLP.** Construction
 521 Grammar (CxG) has been explored in natural lan-
 522 guage processing tasks. Kiselev (2020) constructs
 523 a CxG-based knowledge network for a deeper un-
 524 derstanding of text. Dunn (2023) employs con-
 525 structions to model variation across and dialects.
 526 Xu et al. (2023) leverage constructional informa-
 527 tion to enrich language representation for natural
 528 language understanding tasks. Subsequently, Xu
 529 et al. (2024) encode constructions as inductive bi-
 530 ases to explicitly embed constructional semantics
 531 and guide language modeling. However, there has
 532 been no effort to ascertain whether constructions
 533 can provide benefits in guiding GEC tasks. Our
 534 work aims to bridge this gap.

535 **Construction Inventory Establishment.** An in-
 536 ventory of constructions serves as a valuable re-
 537 source for CxG-based research. Several construc-
 538 tion inventories have been created for various

539 languages (e.g., English, German) by lexicogra-
 540 phers and linguists (Lyngfelt et al., 2018), primar-
 541 ily through manual development, which is labor-
 542 intensive and depends on expert experience. Weis-
 543 sweiler et al. (2024) utilize GPT-3.5 and propose
 544 a hybrid human-LLM corpus construction method,
 545 with a focus on the caused-motion construction. To
 546 establish a comprehensive construction inventory
 547 automatically from corpora, Dunn (2017) proposes
 548 a grammar induction algorithm based on the com-
 549 putation of associations between adjacent words
 550 using a hard threshold. To generate more com-
 551 plete constructions, Xu et al. (2024) introduce a
 552 LM-based approach to assess slot constraints over
 553 longer distances. However, these methods are un-
 554 able to extract potential constructions from ungram-
 555 matical sentences. To this end, we propose a con-
 556 struction prediction model designed to identify ex-
 557 pected constructions directly from ungrammatical
 558 sentences.

6 Conclusion 559

560 In this paper, we propose a construction-guided
 561 grammatical error correction approach (CxGGEC)
 562 that leverages construction grammar (CxG) to en-
 563 hance error detection and correction. Our frame-
 564 work involves three key steps: (1) generating a
 565 comprehensive construction inventory using Cx-
 566 G Learner, (2) identifying constructions in ungram-
 567 matical sentences through a noise-tolerant language
 568 model, and (3) guiding the GEC task by integrat-
 569 ing construction-masked sentences into the training
 570 process. Extensive experiments on both English
 571 and Chinese GEC benchmarks demonstrate the ef-
 572 fectiveness of CxGGEC.

573 Limitations

574 In this study, the limitations can be summarized
575 into two major aspects:

576 (1) Increased input length and slower inference
577 speed. Incorporating constructional information
578 into the model input increases the overall input
579 length, which inevitably slows down the inference
580 speed. This trade-off between additional linguistic
581 information and computational efficiency poses a
582 challenge, especially for real-time or large-scale
583 applications.

584 (2) Randomness in construction prediction. The
585 construction-prediction model exhibits a degree
586 of randomness. Even though the use of dynamic
587 masking strategies improves the model’s ability to
588 generate diverse constructions, it cannot guarantee
589 that the generated constructions fully cover all er-
590 rors in every prediction. To address this limitation,
591 multiple rounds of inference could be applied to en-
592 hance construction coverage for uncovered errors,
593 potentially further improving GEC performance.

594 Ethics Statement

595 In this work, we use publicly available corpora and
596 benchmarks under their licenses. These publicly
597 available data are checked to ensure that they do
598 not include any offensive and illegal content.

599 References

600 Abhijeet Awasthi, Sunita Sarawagi, Rasna Goyal,
601 Sabyasachi Ghosh, and Vihari Piratla. 2019. [Par-](#)
602 [allel iterative edit models for local sequence trans-](#)
603 [duction](#). In *Proceedings of the 2019 Conference on*
604 *Empirical Methods in Natural Language Processing*
605 *and the 9th International Joint Conference on Natu-*
606 *ral Language Processing (EMNLP-IJCNLP)*, pages
607 4260–4270, Hong Kong, China. Association for Com-
608 putational Linguistics.

609 Christopher Bryant, Mariano Felice, Øistein E Ander-
610 sen, and Ted Briscoe. 2019. The bea-2019 shared
611 task on grammatical error correction. In *Proceedings*
612 *of the fourteenth workshop on innovative use of NLP*
613 *for building educational applications*, pages 52–75.

614 Christopher Bryant, Zheng Yuan, Muhammad Reza
615 Qorib, Hannan Cao, Hwee Tou Ng, and Ted Briscoe.
616 2023. [Grammatical error correction: A survey of](#)
617 [the state of the art](#). *Computational Linguistics*,
618 49(3):643–701.

619 Yo Joong Choe, Jiyeon Ham, Kyubyong Park, and Yeoil
620 Yoon. 2019. [A neural grammatical error correction](#)
621 [system built on better pre-training and sequential](#)
622 [transfer learning](#). In *Proceedings of the Fourteenth*

Workshop on Innovative Use of NLP for Building
623 *Educational Applications*, pages 213–227, Florence,
624 Italy. Association for Computational Linguistics. 625

Daniel Dahlmeier and Hwee Tou Ng. 2012. [Better](#)
626 [evaluation for grammatical error correction](#). In *Pro-*
627 *ceedings of the 2012 Conference of the North Amer-*
628 *ican Chapter of the Association for Computational*
629 *Linguistics: Human Language Technologies*, pages
630 568–572, Montréal, Canada. Association for Compu-
631 tational Linguistics. 632

Jonathan Dunn. 2017. Computational learning of
633 construction grammars. *Language and cognition*,
634 9(2):254–292. 635

Jonathan Dunn. 2019. [Frequency vs. association](#)
636 [for constraint selection in usage-based construction](#)
637 [grammar](#). In *Proceedings of the Workshop on Cogni-*
638 *tive Modeling and Computational Linguistics*, pages
639 117–128, Minneapolis, Minnesota. Association for
640 Computational Linguistics. 641

Jonathan Dunn. 2023. [Exploring the construction:](#)
642 [Linguistic analysis of a computational CxG](#). In
643 *Proceedings of the First International Workshop*
644 *on Construction Grammars and NLP (CxGs+NLP,*
645 *GURT/SyntaxFest 2023)*, pages 1–11, Washington,
646 D.C. Association for Computational Linguistics. 647

Junyi Feng. 2012. [Jieba: Chinese text segmentation](#).
648 GitHub repository. Accessed: October 2023. 649

Adele E Goldberg. 1995. Constructions: A construction
650 grammar approach to argument structure. *University*
651 *of Chicago*. 652

Adele E Goldberg. 2003. Constructions: A new theo-
653 retical approach to language. *Trends in cognitive*
654 *sciences*, 7(5):219–224. 655

Ganesh Jawahar, Benoît Sagot, and Djamé Seddah.
656 2019. What does bert learn about the structure of
657 language? In *ACL 2019-57th Annual Meeting of the*
658 *Association for Computational Linguistics*. 659

Ting Jiang, Deqing Wang, Leilei Sun, Huayi Yang,
660 Zhengyang Zhao, and Fuzhen Zhuang. 2021.
661 Lightxml: Transformer with dynamic negative sam-
662 pling for high-performance extreme multi-label text
663 classification. In *Proceedings of the AAAI Confer-*
664 *ence on Artificial Intelligence*, volume 35, pages
665 7987–7994. 666

Marcin Junczys-Dowmunt, Roman Grundkiewicz,
667 Shubha Guha, and Kenneth Heafield. 2018a. Ap-
668 proaching neural grammatical error correction as a
669 low-resource machine translation task. In *Procee-*
670 *dings of the 2018 Conference of the North American*
671 *Chapter of the Association for Computational Lin-*
672 *guistics: Human Language Technologies, Volume 1*
673 *(Long Papers)*, pages 595–606. 674

Marcin Junczys-Dowmunt, Roman Grundkiewicz,
675 Shubha Guha, and Kenneth Heafield. 2018b. [Ap-](#)
676 [proaching neural grammatical error correction as a](#)
677

793	Felix Stahlberg and Shankar Kumar. 2020. Seq2edits: Sequence transduction using span-level edit operations. In <i>Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP)</i> , pages 5147–5159.	2023. GEC-DePenD: Non-autoregressive grammatical error correction with decoupled permutation and decoding. In <i>Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)</i> , pages 1546–1558, Toronto, Canada. Association for Computational Linguistics.	849
794			850
795			851
796			852
797			853
798	Xin Sun, Tao Ge, Furu Wei, and Houfeng Wang. 2021. Instantaneous grammatical error correction with shallow aggressive decoding. In <i>Proceedings of the 59th Annual Meeting of the Association for Computational Linguistics and the 11th International Joint Conference on Natural Language Processing (Volume 1: Long Papers)</i> , pages 5937–5947, Online. Association for Computational Linguistics.	Zheng Yuan, Shiva Taslimipoor, Christopher Davis, and Christopher Bryant. 2021a. Multi-class grammatical error detection for correction: A tale of two systems. In <i>Proceedings of the 2021 Conference on Empirical Methods in Natural Language Processing</i> , pages 8722–8736.	854
799			855
800			856
801			857
802			858
803			859
804			860
805			861
806	Toshikazu Tajiri, Mamoru Komachi, and Yuji Matsumoto. 2012. Tense and aspect error correction for esl learners using global context. In <i>Proceedings of the 50th Annual Meeting of the Association for Computational Linguistics (Volume 2: Short Papers)</i> , pages 198–202.	Zheng Yuan, Shiva Taslimipoor, Christopher Davis, and Christopher Bryant. 2021b. Multi-class grammatical error detection for correction: A tale of two systems. In <i>Proceedings of the 2021 Conference on Empirical Methods in Natural Language Processing</i> , pages 8722–8736, Online and Punta Cana, Dominican Republic. Association for Computational Linguistics.	862
807			863
808			864
809			865
810			866
811			867
812	Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz Kaiser, and Illia Polosukhin. 2017. Attention is all you need. In <i>Advances in Neural Information Processing Systems</i> , volume 30. Curran Associates, Inc.	Baolin Zhang. 2009. Features and functions of the hsk dynamic composition corpus. <i>International Chinese Language Education</i> , 4:71–79.	868
813			869
814			870
815			871
816			872
817	Jason Wei, Maarten Bosma, Vincent Zhao, Kelvin Guu, Adams Wei Yu, Brian Lester, Nan Du, Andrew M. Dai, and Quoc V Le. 2022. Finetuned language models are zero-shot learners. In <i>International Conference on Learning Representations</i> .	Yue Zhang and Zhenghua Li. 2022. Csyngec: Incorporating constituent-based syntax for grammatical error correction with a tailored gec-oriented parser. <i>Preprint</i> , arXiv:2211.08158.	873
818			874
819			875
820			876
821			877
822	Leonie Weissweiler, Abdullatif Köksal, and Hinrich Schütze. 2024. Hybrid human-llm corpus construction and llm evaluation for rare linguistic phenomena. <i>Preprint</i> , arXiv:2403.06965.	Yue Zhang, Zhenghua Li, Zuyi Bao, Jiacheng Li, Bo Zhang, Chen Li, Fei Huang, and Min Zhang. 2022a. MuCGEC: a multi-reference multi-source evaluation dataset for Chinese grammatical error correction. In <i>Proceedings of the 2022 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies</i> , pages 3118–3130, Seattle, United States. Association for Computational Linguistics.	878
823			879
824			880
825			881
826	Lvxiaowei Xu, Zhilin Gong, Jianhua Dai, Tianxiang Wang, Ming Cai, and Jiawei Peng. 2024. Coelm: Construction-enhanced language modeling. In <i>Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)</i> , pages 10061–10081.	Yue Zhang, Bo Zhang, Zhenghua Li, Zuyi Bao, Chen Li, and Min Zhang. 2022b. Syngec: Syntax-enhanced grammatical error correction with a tailored gec-oriented parser. In <i>Proceedings of the 2022 Conference on Empirical Methods in Natural Language Processing</i> , pages 2518–2531.	882
827			883
828			884
829			885
830			886
831			887
832	Lvxiaowei Xu, Jianwang Wu, Jiawei Peng, Jiayu Fu, and Ming Cai. 2022. FCGEC: Fine-grained corpus for Chinese grammatical error correction. In <i>Findings of the Association for Computational Linguistics: EMNLP 2022</i> , pages 1900–1918, Abu Dhabi, United Arab Emirates. Association for Computational Linguistics.	Wei Zhao, Liang Wang, Kewei Shen, Ruoyu Jia, and Jingming Liu. 2019. Improving grammatical error correction via pre-training a copy-augmented architecture with unlabeled data. In <i>Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long and Short Papers)</i> , pages 156–165, Minneapolis, Minnesota. Association for Computational Linguistics.	888
833			889
834			890
835			891
836			892
837			893
838			894
839	Lvxiaowei Xu, Jianwang Wu, Jiawei Peng, Zhilin Gong, Ming Cai, and Tianxiang Wang. 2023. Enhancing language representation with constructional information for natural language understanding. In <i>Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)</i> , page 4685–4705. Association for Computational Linguistics.	Yuanyuan Zhao, Nan Jiang, Weiwei Sun, and Xiaojun Wan. 2018. Overview of the nlpcc 2018 shared task: Grammatical error correction. In <i>Natural Language Processing and Chinese Computing: 7th CCF International Conference, NLPCC 2018, Hohhot, China, August 26–30, 2018, Proceedings, Part II 7</i> , pages 439–445. Springer.	895
840			896
841			897
842			898
843			899
844			900
845			901
846			902
847	Konstantin Yakovlev, Alexander Podolskiy, Andrey Bout, Sergey Nikolenko, and Irina Piontkovskaya.		903
848			904
			905

Algorithm 2: Fixed Masking Using Maximum Coverage

Input: A set of construction schemes $\mathcal{S} = \{\mathcal{C}_1, \mathcal{C}_2, \dots, \mathcal{C}_N\}$. Sentence \mathcal{S}_{sent} .**Output:** The optimal set \mathcal{C}_O .

```
1  $\mathcal{C}_O \leftarrow \{\}$ 
2  $maxCoverage \leftarrow 0$ 
3 foreach  $scheme \mathcal{C}_i \in \mathcal{S}$  do
4    $coverage \leftarrow \text{CALCULATECOVERAGE}(\mathcal{C}_i,$ 
      $\mathcal{S}_{sent})$ 
5   if  $coverage > maxCoverage$  then
6      $maxCoverage \leftarrow coverage$ 
7      $\mathcal{C}_O \leftarrow \mathcal{C}_i$ 
8   end
9 end
10 return  $\mathcal{C}_O$ 
```

A Datasets Used in GEC Models

Dataset	#Sentences	%Error	Usage
CLang8	2,372,119	57.8	Pre-training (†, ‡)
W&I+LOCNESS	34,308	66.3	Fine-tuning (†)
BEA19-Dev	4,384	65.2	Validation (†, ‡)
CoNLL14-Test	1,312	72.3	Testing (†, ‡)
BEA19-Test	4,477	-	Testing (†, ‡)

Table 7: Statistics of English GEC datasets. #Sentences denotes the number of sentences. %Error refers to the proportion of erroneous sentences. †: indicates usage for model based on BART-Large model. ‡: indicates usage for model based on T5-Large model.

Dataset	#Sentences	%Error	Usage
Lang8	1,220,906	89.5	Training
HSK	15,687	60.8	Training
FCGEC-train	36,340	54.5	Training
MuCGEC-dev	1,125	95.1	Validation
MuCGEC-test	5,938	92.2	Testing
FCGEC-test	3,000	54.5	Testing

Table 8: Statistics of Chinese GEC datasets.

B Training Data Examples

We use construction-masked sentences concatenated with the original ungrammatical sentences as inputs to the GEC model and pair them with ground-truth sentences to form parallel corpora for GEC model training. Examples are shown in Table ??.

C Fixed Masking Strategy

Compared to dynamic masking to the train construction prediction model, fixed masking we use can be demonstrated in Algorithm 2. The algorithm examines a predefined set of construction schemes and selects the one that maximizes the area of constructions within the given sentence. The input to the algorithm consists of a set of construction schemes $\mathcal{S} = \{\mathcal{C}_1, \mathcal{C}_2, \dots, \mathcal{C}_N\}$ and a sentence \mathcal{S}_{sent} . The algorithm iteratively evaluates each construction scheme $\mathcal{C}_i \in \mathcal{S}$ to calculate its coverage over the input sentence, relying on the function CALCULATECOVERAGE. The goal is to identify the construction scheme \mathcal{C}_O that achieves the highest coverage with respect to the constructions inherent in the sentence. The ‘maxCoverage’ value is updated whenever a scheme \mathcal{C}_i with higher coverage is encountered, and \mathcal{C}_O is set to \mathcal{C}_i . Finally, the algorithm returns \mathcal{C}_O , which represents the optimal construction masking scheme. However, fixed masking is not conducive to improving the construction prediction model’s generalization performance. Therefore, in comparison, dynamic masking was chosen as a better alternative according to results in Table 4.

Example	Input Sentence (Original Sentence + Construction-Masked Sentence)	Ground-Truth Sentence
Example 1	About winter [SEP] <ADP><NN><NOUN>	About winter
Example 2	This is my second post . [SEP] This <VBZ><PRON><ADJ> post .	This is my second post .
Example 3	People usually get this kind of hypertesion after they become adult . [SEP] People usually get this kind of hypertesion <IN><they><VBP> adult .	People usually get this kind of hypertesion when they become adult .
Example 4	After the initial ceremony , the group photo was taken . [SEP] After <DT><JJ><NOUN> , <DET><NN><NOUN> was taken .	After the initial ceremony , the group photo was taken .
Example 5	One time , I had an Japanese examination . [SEP] One time , I had <DT><JJ><NOUN> .	One time , I had a Japanese examination .

Table 9: Examples of construction-masked sentences paired with ground-truth sentences for GEC training.