

# Building Sequence-to-Sequence Document Revision Models from Matched and Multiple Partially-Matched Datasets

Anonymous ACL submission

## Abstract

This paper defines the document revision task and proposes a novel modeling method that can utilize not only a matched dataset but also multiple partially-matched datasets. In the document revision task, we aim to simultaneously consider multiple perspectives for writing supports. To this end, it is important not only to correct grammatical errors but also to improve readability and perspicuity, through means such as conjunction insertion and sentence reordering. However, it is difficult to prepare enough the matched dataset for the document revision task since this task has to consider multiple perspectives simultaneously. To mitigate this problem, our idea is to utilize not only a limited matched dataset but also various partially-matched datasets that handles individual perspectives, e.g., correcting grammatical errors or inserting conjunctions. Since suitable partially-matched datasets have either been published or can easily be made, we expect to prepare a large amount of these partially-matched datasets. To effectively utilize these multiple datasets, our proposed modeling method incorporates “on-off” switches into sequence-to-sequence modeling to distinguish the matched datasets and individual partially-matched datasets. Experiments using our created document revision datasets demonstrate the effectiveness of the proposed method.

## 1 Introduction

With the advance of natural language processing technology using deep learning, applications for writing support systems have been developed (Tsai et al., 2020; Ito et al., 2020). Such writing support systems often implement a grammatical error correction task that correct errors such as typos and mistakes in inflected verbs forms (Rothe et al., 2021). To advance writing support, it is important not only to correct grammatical errors but also to improve readability and perspicuity. For example, when we manually perform document revision,

we attempt not only to correct grammatical errors but also to split a long sentence into sentences to improve the readability and the perspicuity. In addition, we also consider the relationships between sentences, such as reordering to obtain a consistent order and conjunction insertion. Accordingly, this paper defines a document revision task that simultaneously considers these multiple perspectives for writing support.

In natural language processing area, the document revision task has been studied by breaking it down into partial tasks. The most common partial task is grammatical error correction, and various methods have been proposed to model this task (Sawai et al., 2013; Mizumoto and Matsumoto, 2016; Junczys-Dowmunt and Grundkiewicz, 2016). In recent studies, the sequence-to-sequence (seq2seq) modeling methods has achieved high performance with the advance of deep learning (Yuan and Briscoe, 2016; Junczys-Dowmunt et al., 2018; Rothe et al., 2021). In addition, other famous partial tasks are the sentence ordering (Yin et al., 2019) or discourse relation classification (Liu et al., 2016; Dai and Huang, 2018). Most of these tasks have also been studied with the seq2seq modeling (Wang and Wan, 2019). On the other hand, there are few studies that address multiple perspectives in the document revision task. Lin et al. (2021) addressed the sentence ordering and sentence paraphrasing tasks, and Ihori et al. (2020) addressed multiple perspectives for spoken-to-written style conversion such as style unification, disfluency deletion, punctuation restoration at the same time. However, to the best our knowledge, the document revision task that comprehensively handle multiple perspectives has not well examined. Therefore, we aim to model such document revision task using the promising seq2seq modeling.

There are two difficulties in building the document revision seq2seq models.

- The first difficulty is that the document revision

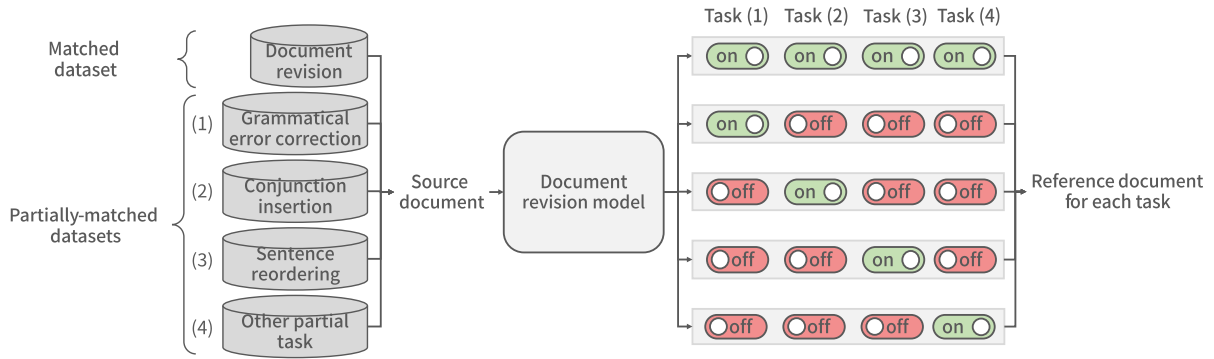


Figure 1: Document revision model using both a dataset for main task and datasets for partial tasks.

084 sion model has to handle multiple perspectives  
 085 simultaneously. Although seq2seq models can  
 086 address any problems that convert a source se-  
 087 quence into a target sequence, handling multi-  
 088 ple perspectives is considered as a difficult  
 089 task.

- 090 • The second difficulty is that improving the  
 091 readability and perspicuity requires precisely  
 092 handling long-range contexts of multiple sen-  
 093 tences. While the conventional grammatical  
 094 error correction tasks take contexts within a  
 095 sentence into consideration, our document re-  
 096 vision task must handle a set of sentences, i.e.,  
 097 document-level information.

098 These two difficulties induce us to prepare a lot of  
 099 training datasets so as to robustly model the docu-  
 100 ment revision task; however, it is difficult to prepare  
 101 enough matched training data because these two  
 102 difficulties also affect the data creation cost.

103 Our key idea to mitigate this problem is to utilize  
 104 not only a limited matched dataset but also various  
 105 partially-matched datasets that handle individual  
 106 perspectives for building the document revision  
 107 models. The partially-matched datasets can be re-  
 108 garded as datasets for the partial tasks. There are  
 109 several existing datasets for grammatical error cor-  
 110 rection (Dahlmeier et al., 2013; Tajiri et al., 2012)  
 111 and sentence ordering (Chen et al., 2016; Huang  
 112 et al., 2016). In addition, datasets can be gener-  
 113 ated heuristically for noisy sentence deletion and  
 114 conjunction insertion tasks. For example, for the  
 115 conjunction insertion task, we can construct paired  
 116 data by deleting and restoring conjunctions from  
 117 existing documents. We expect that these partially-  
 118 matched datasets will be effective for improving  
 119 our document revision task. The important issue  
 120 is how we exactly utilize both a limited matched

121 dataset and various partially-matched datasets for  
 122 building a document revision model.

123 In this paper, we propose a novel modeling  
 124 method that simultaneously utilize both a matched  
 125 dataset and multiple partially-matched datasets. In  
 126 the proposed method, we incorporate multiple “on-  
 127 off” switches into seq2seq modeling so as to distin-  
 128 guish the matched datasets and individual partial-  
 129 matched datasets. Figure 1 shows an example of  
 130 how the proposed method uses multiple switches.  
 131 It is implemented by using switching tokens, which  
 132 were previously proposed by (Ihori et al., 2021b).  
 133 The switching tokens have the role of switching  
 134 the “on” or “off” state for each task. By introduc-  
 135 ing the switching tokens into the seq2seq model-  
 136 ing, the main document revision task and each par-  
 137 tial task can be explicitly distinguished within one  
 138 modeling. We expect that our proposed modeling  
 139 method effectively improves the main document  
 140 revision task by appropriately leveraging knowl-  
 141 edge from partially-matched datasets. Furthermore,  
 142 our proposed method can be combined with self-  
 143 supervised pre-training, which is the most success-  
 144 ful approach in recent modeling methods (Kenton  
 145 and Toutanova, 2019). In this approach, unpaired  
 146 text datasets are used for building a base model in  
 147 a pre-training phase and the model is fine-tuned  
 148 by paired datasets. In natural language generation  
 149 tasks using seq2seq models, several successful self-  
 150 supervised pre-training methods had been proposed  
 151 (Song et al., 2019; Ihori et al., 2021a). We expect  
 152 that our proposed method can be effectively applied  
 153 after performing the self-supervised pre-training.

154 For evaluation, we newly construct a Japanese  
 155 document revision dataset (see Sec. 3). In our ex-  
 156 periments, we used the new dataset as the matched  
 157 dataset, and grammatical error correction and con-  
 158 junction insertion datasets as the partially-matched

159 datasets. Our experimental results demonstrate  
160 that our proposed modeling method effectively im-  
161 proves the document revision performance by using  
162 not only the matched dataset but also the partially-  
163 matched datasets.

164 Our main contributions are as follows:

- 165 • We define a document revision task that si-  
166 multaneously considers multiple perspectives  
167 for writing support, and specify the relation-  
168 ship between our document revision task and  
169 conventional related tasks.
- 170 • We create a novel dataset for a Japanese doc-  
171 ument revision task and detail how we create  
172 it.
- 173 • We present a novel modeling method that can  
174 utilize not only a matched dataset but also  
175 multiple partially matched datasets, and show  
176 the effectiveness of the proposed method in  
177 our experiments.

## 178 2 Related Work

179 The partial tasks that compose a document revision  
180 task have been studied as individual tasks. The  
181 most typical task is the grammatical error correc-  
182 tion task, which corrects the errors in an input text  
183 by deleting, inserting, and replacing words. Many  
184 studies on this task focused on sentence-level er-  
185 rors, and they performed error correction by using a  
186 seq2seq model to achieve high performance (Yuan  
187 and Briscoe, 2016; Junczys-Dowmunt et al., 2018;  
188 Rothe et al., 2021). In addition, recent studies  
189 have introduced the seq2seq pre-training to utilize  
190 a large amount of unpaired data to improve the  
191 performance with a limited amount of paired data  
192 (Lewis et al., 2020; Song et al., 2019; Ihuri et al.,  
193 2021a). Thus, in this work, we investigated the  
194 combination of such pre-training methods and our  
195 proposal. For the grammatical error correction task,  
196 synthetic training data generation is also introduced  
197 as another way to deal with paired-data scarcity  
198 (Grundkiewicz et al., 2019; Kiyono et al., 2020;  
199 Rothe et al., 2021). For the document revision task,  
200 however, it is difficult to generate synthetic data be-  
201 cause the task involves multiple partial tasks such  
202 as grammatical error correction, sentence reorder-  
203 ing, and conjunction insertion.

204 In addition, certain tasks handle multiple sen-  
205 tences, such as a discourse relation classification  
206 task (Liu et al., 2016; Dai and Huang, 2018) and a

207 sentence reordering task (Wang and Wan, 2019). In  
208 the discourse relation classification task, the model  
209 predicts the relation class (e.g., contrast and causal-  
210 ity) of two arguments. In this work, we adopted a  
211 conjunction insertion task that is similar to the dis-  
212 course relation classification task but directly com-  
213 pletes conjunctions according to the relationship  
214 between sentences. Sentence ordering is another  
215 task that considers the document-level coherence  
216 where an input set of sentences are re-arranged into  
217 a logically consistent order. For this task, seq2seq  
218 models like pointer-network were mainly used (Cui  
219 et al., 2018). In a recent work, graph network is  
220 also introduced and achieved high performance  
221 (Yin et al., 2019). Although these studies improved  
222 readability in terms of sentence order, they did not  
223 cover other aspects of the document revision.

224 There are few studies to cover multiple aspects  
225 of document revision at the same time. Lin et al.  
226 (2021) proposed document-level paraphrase genera-  
227 tion task that simultaneously performs the sentence  
228 reordering and sentence rewriting tasks. In this  
229 study, a pseudo dataset for document-level para-  
230 phrase generation task was created and the task was  
231 performed with a specific model architecture. To  
232 perform multiple tasks, the task-specific model ar-  
233 chitecture and matched dataset were needed. Thus,  
234 it is difficult to add a new task for document-level  
235 paraphrase generation task.

## 236 3 Dataset for Document Revision Task

### 237 3.1 Dataset construction

238 In this paper, we present a new dataset for Japanese  
239 document revision task. The dataset contains  
240 paired data consisting of source and reference doc-  
241 uments in Japanese. The source documents were  
242 written by Japanese crowd workers. Also, the  
243 reference documents were revised by Japanese  
244 two labelers. Each document contained multiple  
245 Japanese sentences to enable the consideration of  
246 contextual information. Below, we explain the de-  
247 tails of creating source and reference documents.

248 **Source documents:** To make the source docu-  
249 ments, we employed crowd workers and they wrote  
250 essays consisted of a single paragraph document  
251 in Japanese. The documents have an essay-style  
252 structure, because Japanese schools teach how to  
253 write essays; thus, we expected that many of the  
254 workers could write the essays at the same level.  
255 Specifically, we employed 161 workers whose na-

(1)	Correct the following mistakes. typos, punctuation, kanji, syntax and grammatical errors, spoken-style text, and redundant expressions
(2)	Split long sentences containing more than 60 characters.
(3)	Unify words with different expressions that have the same meaning.
(4)	If there is no subject, restore the subject by using words that have already been mentioned.
(5)	Change the sentence order if it is not appropriate.
(6)	Delete sentences that describe unrelated topics.
(7)	Insert correct conjunctions for the relationships between sentences.

Table 1: Guidelines for document revision.

256 tive language was Japanese. First, we showed the  
 257 workers 48 possible themes, and they individually  
 258 selected 1-15 themes. The 48 themes were chosen  
 259 by the crowdsourcing company from actual themes  
 260 that were used for exam essays in Japan. Next, the  
 261 workers wrote single paragraph documents, each  
 262 of which contained 200-300 characters and four  
 263 or more sentences. These multiple sentences are  
 264 needed to conduct the revision by considering the  
 265 relationship between sentences. Each worker wrote  
 266 1-15 documents per person, and took up to 15 min-  
 267 utes to write each document. Although the workers  
 268 were asked to be careful about typos, they were not  
 269 asked to compose the essay perfectly.

270 **Reference documents:** To revise the source doc-  
 271 uments, we employed two labelers whose native  
 272 language was Japanese. One labeler was licensed  
 273 as a Japanese language teacher, while the other la-  
 274 beler received guidance of revision for a document.  
 275 In the document revision task, we should handle  
 276 multiple perspectives to improve the readability of  
 277 a document. Thus, we asked them to follow the  
 278 revision guidelines listed in Table 1, to ensure that  
 279 the labelers can consider revising from the multi-  
 280 ple perspectives. Table 1 shows the guidelines for  
 281 document revision. In the table, (1) shows the er-  
 282 ror correction task and (2-7) shows the other tasks  
 283 for improving the readability and the perspicuity.  
 284 Since it is difficult to clearly define the readabil-  
 285 ity and the perspicuity, we told labelers specific  
 286 examples of each task. For example, for (2), it is  
 287 possible to divide the sentences according to the  
 288 number of characters, and for (7), we represented  
 289 the list of conjunctions that shows their kinds and  
 290 roles, and asked them to select from this list. We  
 291 expected that the labelers would be able to revise  
 292 documents with equivalent quality by following the

		# of documents	# of sentences
Training	Input	5,000	26,477
	Output	5,000	28,158
Validation	Input	554	2,922
	Output	554	3,128
Test	Input	1,121	6,054
	Output	2,242	12, 831

Table 2: Details of the dataset for document revision.

guidelines. Note that they do not necessarily have  
 293 to consider all the perspectives simultaneously, but  
 294 only made these revisions if there were any mis-  
 295 takes or unnatural points.  
 296

### 3.2 Details 297

Table 2 lists that the details of the resulting dataset  
 298 for document revision task. The dataset is divided  
 299 into a training set, validation set and test set. The  
 300 training and validation sets have one reference doc-  
 301 ument, while the test set has two reference docu-  
 302 ments for each source document. Figure 2 shows an  
 303 example from the dataset. As this example demon-  
 304 strates, the dataset was created while considering  
 305 multiple perspectives simultaneously. For example,  
 306 typo correction, too-long sentence splitting, and  
 307 conjunctions insertion tasks are performed at the  
 308 same as shown in Table 1. To the best of our knowl-  
 309 edge, this is the first dataset to address such multi-  
 310 ple perspectives of the document revision task.  
 311

## 4 Document Revision Models 312

### 4.1 Strategy 313

To build document revision model, we utilize a  
 314 matched dataset for document revision task (cre-  
 315 ated in chapter 3) and multiple partially-matched  
 316 datasets. In this paper, the document revision task  
 317 is referred to as the main task and tasks that handle  
 318 each perspective in the main task are referred to  
 319 as the partial tasks. The partially-matched datasets  
 320 can be regarded as datasets for the partial tasks.  
 321

Our strategy is to incorporate multiple “on-off”  
 322 switches into seq2seq modeling to distinguish the  
 323 matched datasets and individual partially-matched  
 324 datasets. It is implemented by using switching  
 325 tokens (Ihori et al., 2021b). A switching token  
 326 represents the “on” state (the target task) or “off”  
 327 state (not the target task) for each perspective. By  
 328 introducing the switching tokens into the seq2seq  
 329 modeling, the main document revision task and  
 330 each partial task can be explicitly distinguished  
 331 within one modeling.  
 332



Source	ソーシャルメディアの発達により、私たちは好きな情報を簡単に得られるようになった。その一方で、日々膨大な情報にさらされることの弊害も出てきている。ソーシャルメディアでは、自身の趣味嗜好にあ沿った情報を得ようとすることが多い。これまでは新聞やテレビが情報源だったので全員が同じ情報に触れて対等に対話していたが、現在は無意識に自分に近い意見を全体の情報として捉えてしまう人もおり、情報が偏る。ソーシャルメディアは一見情報の宝庫に見えるが、視点を変えれば物事を都合よくとらえる為のツールなのかもしれない。
Reference	ソーシャルメディアの発達で、私たちは好きな情報を簡単に得られるようになった。その一方で、日々膨大な情報にさらされる弊害も出てきている。また、ソーシャルメディアでは、自身の趣味嗜好に沿った情報を得ようとする人が多い。これまでは新聞やテレビが情報源だったので、全員が同じ情報に触れて対等に対話していた。しかし、現在は、無意識に自分に近い意見を全体の情報として捉える人もおり、情報が偏る。そのため、ソーシャルメディアは一見情報の宝庫に見えるが、視点を変えれば、物事を都合よく捉える為のツールなのかもしれない。
Translation	The development of social media has made it easier to get information. On the other hand, there can be difficulties in handling vast amounts of information. Also, in most cases, we only use social media to access our favorite types and sources of information. Previously, many people got the same information from newspapers and television, and thus, they could talk on an equal footing. Now, however, some people unknowingly treat their closely held opinions as complete information, so their information is biased. Therefore, social media seems to be a treasure trove of information, but it may also be a tool for maintaining biased information.

Figure 2: Example from the document revision task dataset.

Figure 3 shows an example of our strategy using switching tokens. In this example, we use grammatical error correction (GEC) dataset, conjunction insertion (CI) dataset, and the main task dataset to build the document revision model. In this case, we use six switching tokens [gec\_on], [ci\_on], [other\_on], [gec\_off], [ci\_off], and [other\_off]. Here, we specify the “other” token because the main task handle other perspectives that are not considered in the grammatical error correction and conjunction insertion tasks, as listed in Table 1. The seq2seq models are split into an encoder network and a decoder network. These switching tokens are utilized for inputs of the decoder network as given contexts. In a training phase, we use all datasets for building a seq2seq model while distinguishing each task using above switching tokens. In an inference phase, we expect to perform the main task by feeding [gec\_on], [ci\_on], and [other\_on]. Note that we can also perform the grammatical error correction or conjunction insertion by feeding appropriate switching tokens.

## 4.2 Proposed modeling method

In this paper, we propose a novel modeling method that simultaneously utilize both a matched dataset and multiple partially-matched datasets. In the proposed method, we incorporate multiple “on-off” switches into seq2seq modeling so as to distinguish the matched datasets and individual partially-matched datasets.

**Modeling:** We define the source document as  $\mathbf{X} = \{x_1, \dots, x_m, \dots, x_M\}$  and the reference document as  $\mathbf{Y} = \{y_1, \dots, y_n, \dots, y_N\}$ , where  $M$  and  $N$  are the numbers of tokens in source and reference documents, respectively.  $x_m$  and  $y_n$  are

tokens which include not only characters or words but also punctuation marks. Note that  $X$  and  $Y$  involves multiples sentences.

Our proposed document revision model predicts the generation probabilities of a reference document  $\mathbf{Y}$  given a source document  $\mathbf{X}$  and switching tokens  $s_{1:T} = \{s_1, \dots, s_t, \dots, s_T\}$ , where  $T$  is the number of “on-off” switches. The generation probability of  $\mathbf{Y}$  is defined as

$$P(\mathbf{Y}|\mathbf{X}, s_{1:T}; \Theta) = \prod_{n=1}^N P(y_n|y_{1:n-1}, \mathbf{X}, s_{1:T}; \Theta), \quad (1)$$

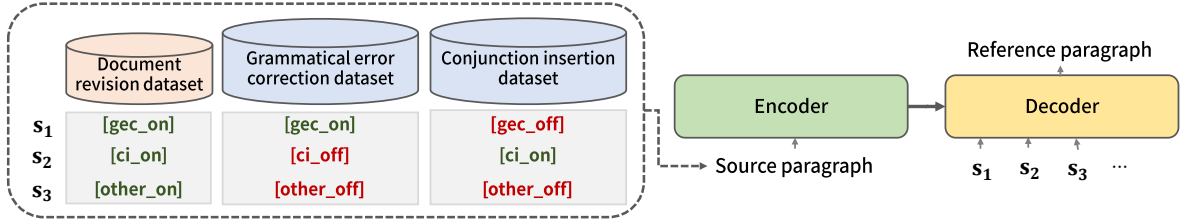
where  $\Theta$  represents the trainable parameters.  $s_t$  is the  $t$ -th switching token represented as

$$s_t \in \{[t\text{-th task\_on}], [t\text{-th task\_off}]\}. \quad (2)$$

In this paper, we use Transformer pointer-generator networks (Deaton, 2019) for this modeling. Transformer pointer-generator networks are effective for monolingual translation tasks because they contain a copy mechanism that copies tokens from a source text to help generate infrequent tokens. Note that our method does not change the architecture of a transformer pointer-generator network, but merely adds switching tokens to the model input.

**Pre-training:** In this paper, we use a MAsked Pointer-Generator Network (MAPGN) (Ihori et al., 2021a) because it is a suitable pre-training method for pointer-generator networks. In MAPGN, the pointer-generator network is pre-trained by predicting a sentence fragment  $y_{a:b}$  giving a masked sequence  $\mathbf{Y}_{/a:b}$  and. Here,  $\mathbf{Y}_{/a:b}$  denotes a fragment in which positions a to b are masked, and

### Joint modeling of a matched dataset and partially-matched datasets:



### Decoding for document revision task:

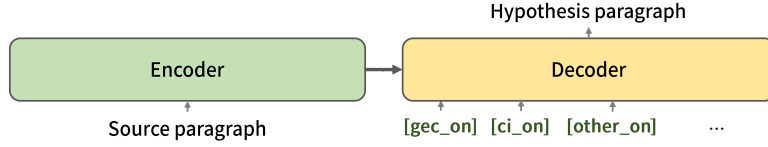


Figure 3: Example of joint modeling based on switching-token.

$y_{a:b}$  denotes a sentence fragment of  $\mathbf{Y}$  from  $a$  to  $b$ . The model parameter set can be optimized from unpaired dataset  $\mathcal{D}^u$ . The training loss function  $\mathcal{L}$  is defined as

$$\begin{aligned} \mathcal{L} &= - \sum_{(\mathbf{Y}) \in \mathcal{D}^u} \log P(y_{a:b} | y_{a-1}, \mathbf{Y}_{/a:b}; \Theta), \quad (3) \\ &= - \sum_{(\mathbf{Y}) \in \mathcal{D}^u} \sum_{t=a}^b \log P(y_t | y_{a-1:t-1}, \mathbf{Y}_{/a:b}; \Theta). \end{aligned}$$

Note that all switching tokens have to be included in the vocabulary in the pre-training.

**Fine-tuning:** In our proposed method, the matched dataset  $\mathcal{D}^m$ , and multiple partially-matched datasets  $\{\mathcal{D}_1^{\text{pm}}, \dots, \mathcal{D}_t^{\text{pm}}, \dots, \mathcal{D}_T^{\text{pm}}\}$  are trained jointly in a single model. The training loss function  $\mathcal{L}$  is defined as

$$\mathcal{L} = \mathcal{L}^m + \sum_{t=1}^T \mathcal{L}_t^{\text{pm}}, \quad (4)$$

where  $\mathcal{L}^m$  is the loss function against the main task and it is computed from

$$\mathcal{L}^m = - \sum_{(\mathbf{X}, \mathbf{Y}) \in \mathcal{D}^m} \log P(\mathbf{Y} | \mathbf{X}, \hat{s}_{1:T}; \Theta), \quad (5)$$

where  $\hat{s}_{1:T} = \{\hat{s}_1, \dots, \hat{s}_T\}$  are switching tokens and  $\hat{s}_t$  is represented as

$$s_t = [t\text{-th task\_on}]. \quad (6)$$

$\mathcal{L}_t^{\text{pm}}$  is the loss function against the  $t$ -th partial task and it is computed from

$$\mathcal{L}_t^{\text{pm}} = - \sum_{(\mathbf{X}, \mathbf{Y}) \in \mathcal{D}_t^{\text{pm}}} \log P(\mathbf{Y} | \mathbf{X}, \bar{s}_{1:T}; \Theta), \quad (7)$$

where  $\bar{s}_{1:T} = \{\bar{s}_1, \dots, \bar{s}_T\}$  are switching tokens and  $\bar{s}_{t'}$  is represented as

$$\bar{s}_{t'} = \begin{cases} [t'\text{-th task\_on}] & \text{if } t' = t, \\ [t'\text{-th task\_off}] & \text{otherwise.} \end{cases} \quad (8)$$

**Decoding:** The decoding problem using switching tokens is defined as

$$\hat{\mathbf{Y}} = \arg \max_{\mathbf{Y}} P(\mathbf{Y} | \mathbf{X}, s_{1:T}; \Theta). \quad (9)$$

The model can perform the document revision task or each partial task according to the given switching tokens.

## 5 Experiments

We experimentally evaluated the effectiveness of the proposed modeling method that can utilize both matched and multiple partially-matched datasets.

### 5.1 Dataset

For preparing the partially-matched datasets, we adopted the grammatical error correction task (gec) and the conjunction insertion task (ci) as partial tasks. Accordingly, we used three datasets: document revision dataset described in section 3, a Japanese grammatical error correction dataset (Tanaka et al., 2020), and a conjunction insertion dataset. The Japanese grammatical error correction dataset was obtained from revision history on Wikipedia. It contained four categories of Japanese typos: erroneous substitution, deletion, insertion, and kanji-conversion. The conjunction insertion dataset was constructed based on Japanese Wiki-40B dataset (Guo et al., 2020), which is a high

		# of documents	# of sentences
Training	a).	5,000	26,477
	b).	-	506,786
	c).	90,000	533,422
Validation	a).	554	2,922
	b).	-	8,542
	c).	10,000	59,396
Test	a).	1,121	6,054
	b).	-	8,542
	c).	1,000	6,026
Switchs	a).	[gec_on][ci_on][other_on]	
	b).	[gec_on][ci_off][other_off]	
	c).	[gec_off][ci_on][other_off]	

a. Document revision dataset  
b. Japanese grammatical error correction dataset  
c. Conjunction insertion dataset

Table 3: Details of document revision task datasets

quality processed Wikipedia dataset. To construct this dataset, first, we divided the Wiki-40B dataset into single paragraph documents and selected the documents that contained conjunctions. Next, we deleted the conjunctions from each document, and we used the resulting and original documents as paired data.

For unpaired data which is used for self-supervised pre-training, we prepared 880k single paragraph documents from Wiki-40B dataset that were not used in the conjunction insertion dataset. The details of these datasets are listed in Table 3, where “Switch” refers to switching tokens. We use six switching tokens [gec\_on], [ci\_on], [other\_on], [gec\_off], [ci\_off], and [other\_off] for training and decoding. In decoding, we can also perform the grammatical error correction or conjunction insertion tasks by feeding appropriate switching tokens. Thus, we use test set for each partial task to evaluate each partial task performance. For example, when the model performs the grammatical error correction task in the decoding, the switching tokens [gec\_on], [ci\_off], and [other\_off] are given in the decoder. Moreover, we compare each partial task performance using the joint modeling with a individual model performance. Note that the the number of documents corresponds to the number of sentences in the Japanese grammatical error correction dataset because the dataset is consisted not of documents but of single sentence.

## 5.2 Setup

For evaluation purposes, we constructed 11 Transformer-based pointer-generator networks. (1) a document revise model, (2) (1) with pre-training, and (3) a grammatical error correction model, (4) a

conjunction insertion model, (5) a modeling of the document revision and grammatical error correction datasets, (6) (5) with switching tokens, (7) a modeling of the document revision and conjunction insertion datasets, (8) (7) with switching tokens, (9) a modeling of all three datasets, (10) (9) with switching tokens, (11) (10) with pre-training. (1), (3), and (4) are trained using only each task dataset. We use the unpaired data for pre-training in these models. Note that the Transformer-based pointer-generator network architecture is the same in all of these models.

As for the model details, we used the following configurations. The encoder had a 4-layer transformer encoder block with 512 units, while the decoder had a 2-layer transformer decoder block with 512 units. The output unit size (corresponding to the number of tokens in the pre-training data) was set to 12,773. To train the Transformer pointer-generator networks, we used the RAdam optimizer (Liu et al., 2019) and label smoothing (Lukasik et al., 2020) with a smoothing parameter of 0.1. We set the mini-batch size to 32 documents and the dropout rate in each Transformer block to 0.1. All trainable parameters were initialized randomly, and we used characters as tokens. The pre-training and fine-tuning were the same setups. For decoding, we used the beam search algorithm with a beam size of 4.

For evaluation, we calculated automatic evaluation scores in terms of two metrics: GLEU (Napoles et al., 2015), and  $F_{0.5}$ . Specifically, we calculated these metrics for characters and used 4-grams for GLEU.  $F_{0.5}$  score is calculated using the characters in the generated documents. In addition, we also calculated the F1 score for conjunction insertion, denoted as C-F1, to evaluate the performance of conjunction insertion task. Note that multiple conjunctions can have the same meaning (e.g., “but”, and “however”). We thus evaluated whether the system could insert conjunctions with the correct meaning.

## 5.3 Results

Table 4 shows the results of 11 Transformer pointer-generator networks. In the table, the models of (1)-(11) are described in section 5.2, and (6), (8), (10), and (11) are our proposals. The columns “Switch” and “Pre-train” indicate whether the proposed switching tokens and the pre-training are introduced or not, respectively. The row “Source”

				Document revision			GEC		CI		
	Dataset	Switch	Pre-train	GLEU	F <sub>0.5</sub>	C-F1	GLEU	F <sub>0.5</sub>	GLEU	F <sub>0.5</sub>	C-F1
Source	-	-	-	0.886	0	0	-	-	-	-	-
(1)	a	w/o	w/o	0.857	0.198	0.193	-	-	-	-	-
(2)		w/o	w/	0.884	0.321	0.211	-	-	-	-	-
(3)	b	w/o	w/o	-	-	-	<b>0.943</b>	<b>0.635</b>	-	-	-
(4)	c	w/o	w/o	-	-	-	-	-	0.964	0.198	0.230
(5)	a + b	w/o	w/o	0.863	0.189	0.164	-	-	-	-	-
(6)		w/	w/o	0.887	0.278	0.163	-	-	-	-	-
(7)	a + c	w/o	w/o	0.881	0.155	0.101	-	-	-	-	-
(8)		w/	w/o	0.888	0.234	0.214	-	-	-	-	-
(9)	a + b + c	w/o	w/o	0.883	0.236	0.205	0.932	0.613	0.966	0.207	0.222
(10)		w/	w/o	0.889	0.282	0.270	<b>0.943</b>	0.630	<b>0.967</b>	<b>0.239</b>	<b>0.263</b>
(11)		w/	w/	<b>0.892</b>	<b>0.333</b>	<b>0.274</b>	-	-	-	-	-

a. Document revision dataset    b. Japanese grammatical error correction dataset    c. Conjunction insertion dataset

Table 4: Results of document revision, grammatical error correction (GEC), and conjunction insertion (CI) tasks.

indicate the results for source documents in the document revision task dataset.

First, we describe the results of the document revision task. The scores of the task with the switching tokens were higher than those without the switching tokens as shown in lines (5) v.s. (6), (7) v.s. (8), and (9) v.s. (10) in the table. In addition, the scores with switching tokens of lines (6), (8), and (10) were higher than the score of the system trained only with the main task data (1). Among the system (6), (8) and (10), the system trained with all three datasets (10) performed the best. These results indicate that the switching tokens are effective for the joint modeling, and the more partial tasks we use, the better the performance of the main task is. In addition, when we compare the results of lines (10) with (11), the results with pre-training outperformed those without pre-training. This indicate that our proposed method can be effectively applied after performing the self-supervised pre-training.

Next, we focus on the results of the performance of each partial task. The switching-token-based joint modeling can perform each partial task by feeding appropriate switching tokens. Thus, we compare the results of a model using each task dataset individually with using a matched dataset and each task dataset simultaneously. In grammatical error correction, the performance of individual modeling and switching-token-based joint modeling were not significantly different. On the other hand, the performance of joint modeling without switching tokens under-performed that of the individual modeling. For conjunction insertion task, the results of joint modeling outperformed that of individual modeling. Also, the results of joint modeling with switching tokens outperformed those

without switching tokens. Therefore, these results indicated that switching-token-based joint modeling can improve the performance of the main task without impairing the performance of each task.

## 6 Conclusion

In this paper, we examined the document revision task with a novel modeling method that can that can utilize not only a matched dataset but also multiple partially-matched datasets. In our document revision task, we revise document descriptions by considering not only to correct grammatical errors but also to improve readability and perspicuity. In our proposed modeling method, we incorporate multiple “on-off” switches into seq2seq modeling so as to distinguish the matched datasets and individual partially matched datasets. The key strength is to effectively improve main document revision task by appropriately leveraging knowledge from partially-matched dataset. The experimental results using our created Japanese document revision dataset demonstrated that our proposed modeling method can improves the document revision performance by utilizing datasets for the grammatical error correction task and the conjunction insertion task. In addition, our proposed method can be effectively applied after performing the self-supervised pre-training. In our future work, we will develop a model architecture that is suitable for handling much longer documents.

## References

- Xinchi Chen, Xipeng Qiu, and Xuanjing Huang. 2016. Neural sentence ordering. *arXiv preprint arXiv:1607.06952*.
- Baiyun Cui, Yingming Li, Ming Chen, and Zhongfei



607	Zhang. 2018. Deep attentive sentence ordering network. In <i>Proceedings of the 2018 Conference on Empirical Methods in Natural Language Processing</i> , pages 4340–4349.	
608		
609		
610		
611	Daniel Dahlmeier, Hwee Tou Ng, and Siew Mei Wu. 2013. Building a large annotated corpus of learner english: The nus corpus of learner english. In <i>Proceedings of the eighth workshop on innovative use of NLP for building educational applications</i> , pages 22–31.	
612		
613		
614		
615		
616		
617	Zeyu Dai and Ruihong Huang. 2018. Improving implicit discourse relation classification by modeling inter-dependencies of discourse units in a paragraph. In <i>Proc. Conference of the North American Chapter of the Association for Computational Linguistics (NAACL)</i> , pages 141–151.	
618		
619		
620		
621		
622		
623	John Deaton. 2019. Transformers and pointer-generator networks for abstractive summarization.	
624		
625	Roman Grundkiewicz, Marcin Junczys-Dowmuntz, and Kenneth Heafield. 2019. Neural grammatical error correction systems with unsupervised pre-training on synthetic data. In <i>Proc. Workshop on Innovative Use of NLP for Building Educational Applications (BEA Workshop)</i> , pages 252–263.	
626		
627		
628		
629		
630		
631	Mandy Guo, Zihang Dai, Denny Vrandečić, and Rami Al-Rfou. 2020. Wiki-40b: Multilingual language model dataset. In <i>Proc. Language Resources and Evaluation Conference (LREC)</i> , pages 2440–2452.	
632		
633		
634		
635	Ting-Hao Huang, Francis Ferraro, Nasrin Mostafazadeh, Ishan Misra, Aishwarya Agrawal, Jacob Devlin, Ross Girshick, Xiaodong He, Pushmeet Kohli, Dhruv Batra, et al. 2016. Visual storytelling. In <i>Proc. Conference of the North American Chapter of the Association for Computational Linguistics (ACL)</i> , pages 1233–1239.	
636		
637		
638		
639		
640		
641		
642	Mana Ihuri, Naoki Makishima, Tomohiro Tanaka, Akihiko Takashima, Shota Orihashi, and Ryo Masumura. 2021a. Mapgn: Masked pointer-generator network for sequence-to-sequence pre-training. In <i>Proc. International Conference on Acoustics, Speech and Signal Processing (ICASSP)</i> , pages 7563–7567.	
643		
644		
645		
646		
647		
648	Mana Ihuri, Naoki Makishima, Tomohiro Tanaka, Akihiko Takashima, Shota Orihashi, and Ryo Masumura. 2021b. Zero-Shot Joint Modeling of Multiple Spoken-Text-Style Conversion Tasks Using Switching Tokens. In <i>Proc. International Speech Communication Association (Interspeech)</i> , pages 776–780.	
649		
650		
651		
652		
653		
654	Mana Ihuri, Akihiko Takashima, and Ryo Masumura. 2020. Parallel corpus for Japanese spoken-to-written style conversion. In <i>Proc. Language Resources and Evaluation Conference (LREC)</i> , pages 6346–6353.	
655		
656		
657		
658	Takumi Ito, Tatsuki Kuribayashi, Masatoshi Hidaka, Jun Suzuki, and Kentaro Inui. 2020. Langsmith: An interactive academic text revision system. In <i>Proc. Conference on Empirical Methods in Natural Language Processing (EMNLP)</i> , pages 216–226.	
659		
660		
661		
662		
	Marcin Junczys-Dowmunt and Roman Grundkiewicz. 2016. Phrase-based machine translation is state-of-the-art for automatic grammatical error correction. In <i>Proc. Conference on Empirical Methods in Natural Language Processing (EMNLP)</i> , pages 1546–1556.	663
		664
		665
		666
		667
	Marcin Junczys-Dowmunt, Roman Grundkiewicz, Shubha Guha, and Kenneth Heafield. 2018. Approaching neural grammatical error correction as a low-resource machine translation task. In <i>Proc. Conference of the North American Chapter of the Association for Computational Linguistics (NAACL)</i> , pages 595–606.	668
		669
		670
		671
		672
		673
		674
	Jacob Devlin Ming-Wei Chang Kenton and Lee Kristina Toutanova. 2019. Bert: Pre-training of deep bidirectional transformers for language understanding. In <i>Proc. Annual Conference of the North American Chapter of the Association for Computational Linguistics (NAACL-HLT)</i> , pages 4171–4186.	675
		676
		677
		678
		679
		680
	Shun Kiyono, Jun Suzuki, Masato Mita, Tomoya Mizumoto, and Kentaro Inui. 2020. An empirical study of incorporating pseudo data into grammatical error correction. In <i>Proc. Conference on Empirical Methods in Natural Language Processing and International Joint Conference on Natural Language Processing (EMNLP-IJCNLP)</i> , pages 1236–1242.	681
		682
		683
		684
		685
		686
		687
	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 <i>Proc. Annual Meeting of the Association for Computational Linguistics (ACL)</i> , pages 7871–7880.	688
		689
		690
		691
		692
		693
		694
	Zhe Lin, Yitao Cai, and Xiaojun Wan. 2021. Towards document-level paraphrase generation with sentence rewriting and reordering. In <i>Findings of the Association for Computational Linguistics: EMNLP 2021</i> , pages 1033–1044.	695
		696
		697
		698
		699
	Liyuan Liu, Haoming Jiang, Pengcheng He, Weizhu Chen, Xiaodong Liu, Jianfeng Gao, and Jiawei Han. 2019. On the variance of the adaptive learning rate and beyond. In <i>Proc. International Conference on Learning Representations (ICLR)</i> .	700
		701
		702
		703
		704
	Yang Liu, Sujian Li, Xiaodong Zhang, and Zhifang Sui. 2016. Implicit discourse relation classification via multi-task neural networks. In <i>Proceedings of the AAAI Conference on Artificial Intelligence</i> .	705
		706
		707
		708
	Michal Lukasik, Srinadh Bhojanapalli, Aditya Menon, and Sanjiv Kumar. 2020. Does label smoothing mitigate label noise? In <i>Proc. the International Conference on Machine Learning (ICML)</i> , pages 6448–6458.	709
		710
		711
		712
		713
	Tomoya Mizumoto and Yuji Matsumoto. 2016. Discriminative reranking for grammatical error correction with statistical machine translation. In <i>Proc. Conference of the North American Chapter of the Association for Computational Linguistics (NAACL)</i> , pages 1133–1138.	714
		715
		716
		717
		718
		719

720 Courtney Napoles, Keisuke Sakaguchi, Matt Post, and  
721 Joel Tetreault. 2015. Ground truth for grammatical  
722 error correction metrics. In *Proc. Annual Meeting*  
723 *on Association for Computational Linguistics (ACL)*,  
724 pages 588–593.

725 Sascha Rothe, Sebastian Krause, Jonathan Mallinson,  
726 Eric Malmi, and Aliaksei Severyn. 2021. A simple  
727 recipe for multilingual grammatical error correction.  
728 In *Proc. Annual Meeting of the Association for Com-*  
729 *putational Linguistics and the International Joint*  
730 *Conference on Natural Language Processing (ACL-*  
731 *IJCNLP)*, pages 702–707.

732 Yu Sawai, Mamoru Komachi, and Yuji Matsumoto.  
733 2013. A learner corpus-based approach to verb sug-  
734 gestion for ESL. In *Proc. Annual Meeting of the As-*  
735 *sociation for Computational Linguistics (ACL)*, pages  
736 708–713.

737 Kaitao Song, Xu Tan, Tao Qin, Jianfeng Lu, and Tie-  
738 Yan Liu. 2019. Mass: Masked sequence to sequence  
739 pre-training for language generation. In *Proc. Inter-*  
740 *national Conference on Machine Learning (ICML)*,  
741 pages 5926–5936.

742 Toshikazu Tajiri, Mamoru Komachi, and Yuji Mat-  
743 sumoto. 2012. Tense and aspect error correction  
744 for ESL learners using global context. In *Proc. An-*  
745 *ual Meeting of the Association for Computational*  
746 *Linguistics (ACL)*, pages 198–202.

747 Yu Tanaka, Yugo Murawaki, Daisuke Kawahara, and  
748 Sadao Kurohashi. 2020. Building a Japanese typo  
749 dataset from Wikipedia’s revision history. In *Proc.*  
750 *Annual Meeting of the Association for Computational*  
751 *Linguistics (ACL)*, pages 230–236.

752 Chung-Ting Tsai, Jhih-Jie Chen, Ching-Yu Yang, and  
753 Jason S. Chang. 2020. LingleWrite: a coaching  
754 system for essay writing. In *Proc. Annual Meeting of*  
755 *the Association for Computational Linguistics (ACL)*,  
756 pages 127–133.

757 Tianming Wang and Xiaojun Wan. 2019. Hierarchical  
758 attention networks for sentence ordering. In *Proc.*  
759 *AAAI Conference on Artificial Intelligence (AAAI)*,  
760 pages 7184–7191.

761 Yongjing Yin, Linfeng Song, Jinsong Su, Jiali Zeng,  
762 Chulun Zhou, and Jiebo Luo. 2019. Graph-  
763 based neural sentence ordering. *arXiv preprint*  
764 *arXiv:1912.07225*.

765 Zheng Yuan and Ted Briscoe. 2016. Grammatical error  
766 correction using neural machine translation. In *Proc.*  
767 *Conference of the North American Chapter of the*  
768 *Association for Computational Linguistics (NAACL)*,  
769 pages 380–386.