

# Instruct-SCTG: Guiding Sequential Controlled Text Generation through Instructions

Anonymous ACL submission

## Abstract

001 Instruction-tuned large language models have  
002 shown remarkable performance in aligning gener-  
003 ated text with user intentions across various  
004 tasks. However, maintaining human-like dis-  
005 course structure in the generated text remains  
006 a challenging research question. In this pa-  
007 per, we propose Instruct-SCTG, a flexible and  
008 effective sequential framework that harnesses  
009 instruction-tuned language models to generate  
010 structurally coherent text in both fine-tuned and  
011 zero-shot setups. Our framework generates arti-  
012 cles in a section-by-section manner, aligned  
013 with the desired human structure using natural  
014 language instructions. Furthermore, we intro-  
015 duce a new automatic metric that measures dis-  
016 course divergence in a fuzzy manner. Extensive  
017 experiments on three datasets from representa-  
018 tive domains of news and recipes demonstrate  
019 the state-of-the-art performance of our frame-  
020 work in imposing discourse structure during  
021 text generation, as verified by both automatic  
022 and human evaluation. Our code will be avail-  
023 able on Github.

## 024 1 Introduction

025 The recent progress in Language Models (LMs)  
026 have attracted widespread attention from both  
027 academia and industry. These models, pow-  
028 ered by massive corpora and advanced hardware,  
029 have demonstrated improving performance across  
030 various NLP benchmarks, ranging from genera-  
031 tive tasks, such as Machine Translation or Data-  
032 to-Text generation, to understanding tasks, e.g.  
033 GLUE (Wang et al., 2018). In particular, Large  
034 Language Models (LLMs) designed for instruction-  
035 following, such as ChatGPT<sup>1</sup> and Flan-T5 (Chung  
036 et al., 2022), exhibit impressive capabilities in com-  
037 prehending instructions expressed in natural lan-  
038 guage and precisely aligning the model outputs  
039 with human intentions.

<sup>1</sup><https://openai.com/blog/chatgpt>

040 Generating high-quality text is essential for var-  
041 ious Natural Language Generation (NLG) tasks.  
042 However, certain tasks, such as news report genera-  
043 tion, require more than just textual fluency. Effec-  
044 tively organizing the underlying discourse structure  
045 of an article can help readers quickly grasp key in-  
046 formation, enhancing engagement and readability.  
047 For example, an experienced journalist can coher-  
048 ently structure the core event, background, conse-  
049 quence, critics’ evaluations and other elements of a  
050 news report. As shown in Fig. 1, a well-structured  
051 report can efficiently deliver event information, cap-  
052 ture readers’ attention and even convey opinions.  
053 The task of text generation with specific discourse  
054 structure constraints has long been a research focus  
055 in the field covering various domains, including  
056 stories, news, recipes and question answering. We  
057 address this challenge as the task of Sequential  
058 Controlled Text Generation (SCTG), previously  
059 formulated by Spangher et al. (2022). In SCTG,  
060 the goal is to generate coherent text following an  
061 input prompt and a sequence of control code.

062 In this paper, we propose Instruct-SCTG, a  
063 simple yet effective framework that harnesses  
064 instruction-following LMs to generate structurally  
065 coherent text. Specifically, our framework breaks  
066 down the generation task into a sequence of sub-  
067 tasks and guides the Supervised Fine-tuned (SFT)  
068 LMs sequentially to produce content section by  
069 section through natural language instructions. This  
070 approach effectively aligns the resulting articles  
071 with the given discourse structures, enhancing the  
072 overall coherence and readability of the generated  
073 text. We also investigate crucial factors to con-  
074 sider during the SFT stage, such as different levels  
075 of discourse information exposure. Furthermore,  
076 to evaluate the adherence of generated articles to  
077 the input control codes, we introduce a novel auto-  
078 matic metric that measures discourse divergence in  
079 a fuzzy positional manner.

080 We conducted extensive experiments using three

## "Britain's Vision for 2100: Spaceports and Sky Farms Propel the Nation's Innovation"

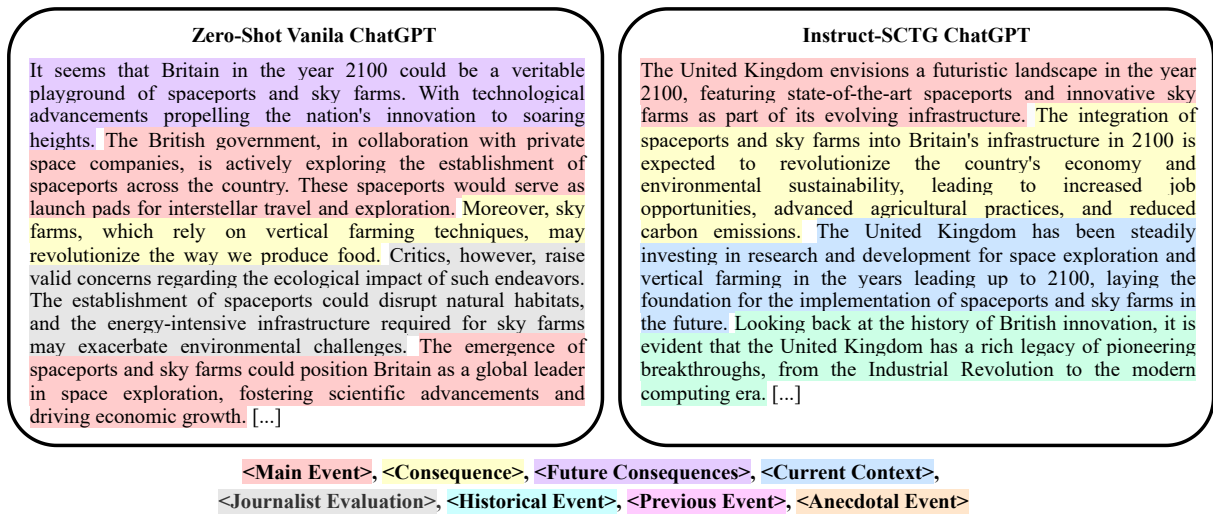


Figure 1: Comparing examples with discourse role labels: left (zero-shot ChatGPT) vs. right (Instruct-SCTG framework utilizing zero-shot ChatGPT as backbone generator). The right article exhibits improved content flow and enhanced discourse structure.

081 datasets from two representative domains, i.e. news  
 082 and recipes. For news articles, we utilized the All-  
 083 The-News dataset<sup>2</sup> from Kaggle and the News Dis-  
 084 course dataset (Choubey et al., 2020). For recipe  
 085 generation, the experiments were performed on the  
 086 Recipe1M+ dataset (Marin et al., 2019). We assess  
 087 the textual fluency and structural coherence of the  
 088 generated text with both automatic and human eval-  
 089 uations. The results demonstrate the effectiveness  
 090 of our framework in controlling LMs to generate  
 091 text adhering to the given discourse structures.

092 In summary, our contributions are three-folds:  
 093 Firstly, we introduce a straightforward yet effec-  
 094 tive framework that leverages instruction-following  
 095 LMs to generate structurally coherent texts in the  
 096 task of SCTG, achieving state-of-the-art (SOTA)  
 097 performance on three datasets from two representa-  
 098 tive domains. Secondly, we introduce a novel auto-  
 099 matic metric that can effectively measure the fuzzy  
 100 adherence of discourse structure. Lastly, our work  
 101 is the first one that explore the design of instruc-  
 102 tions to exert control over the underlying discourse  
 103 structure during text generation.

## 104 2 Background and Related Works

### 105 2.1 Instruction Fine-tuned Language Model

106 Instruction-following LMs are language models  
 107 specially optimized to comprehend and execute nat-  
 108 ural language instructions. These models leverage

<sup>2</sup>kaggle.com/snapcrack/all-the-news.

109 large-scale Pre-trained Language Models (PLM)  
 110 like GPT-3 and incorporate an additional super-  
 111 vised aligning fine-tuning process. Their recent  
 112 emergence has significantly advanced the under-  
 113 standing of human intentions and the generation  
 114 process conditioning on those intentions.

115 For instance, InstructGPT (Ouyang et al., 2022)  
 116 fine-tunes GPT-3 (Brown et al., 2020) to achieve hu-  
 117 man desired model behavior through reinforcement  
 118 learning from human feedback (RLHF, Christiano  
 119 et al. (2017); Stiennon et al. (2020)). Similarly,  
 120 Flan-T5 (Chung et al., 2022) fine-tunes the T5 lan-  
 121 guage model (Raffel et al., 2020) using a diverse  
 122 range of instruction templates from a collection of  
 123 data sources. Another example is Alpaca, proposed  
 124 by Taori et al. (2023), which is an instruction fine-  
 125 tuned Language model based on LLaMA (Touvron  
 126 et al., 2023), using an instruction dataset gener-  
 127 ated in the style of self-instruct (Wang et al., 2022).  
 128 These instruction-following LLMs showcase the  
 129 progress in leveraging instructions to guide lan-  
 130 guage generation, facilitating a more interactive  
 131 and controllable generation process.

### 132 2.2 Discourse Structure

133 Discourse structure investigates the organization  
 134 of language into larger units like paragraphs, sec-  
 135 tions, and complete articles. In this work, we fo-  
 136 cus on the communicative functions within entire  
 137 articles served by those linguistic units. There-  
 138 fore, texts from different domains are characterized

by different discourse schemas, as their linguistic units also play different functional roles. The discourse roles of scientific papers or experimental abstracts (Liddy, 1991; Mizuta et al., 2006) include background, methodology, experiments and findings. In the domain of long-form question answering Xu et al. (2022), the discourse function of each sentence can be answer, summary, example and so on. Liu et al. (2022) developed a discourse schema for recipes based on actions and controlled the generation process according to the predicted discourse sequences. The explicit functional discourse structure of news reports was addressed (Van Dijk, 2013; Choubey et al., 2020) by defining roles based on their relations with the main event, such as consequence and journalist evaluation.

Multiple established frameworks also proposed different definition of discourse structure, which focus on how each linguistic unit relates to each other through discourse connectives, such as causal, temporal, etc. For instance, Rhetorical Structure Theory, RST (Mann and Thompson, 1988), seeks to identify rhetorical relations between text segments and form a hierarchical organization of discourse. The Penn Discourse Treebank, PDTB (Prasad et al., 2008), defines its schema based on low-level discourse connectives presented in the text.

### 2.3 Sequential Controlled Text Generation

Extensive research has been conducted on Controlled Text Generation (CTG) to enable the control of attributes such as lexical constraints, style and length in the output of PLM. One notable example is prefix-tuning, introduced by Li and Liang (2021), which only optimizes a short task-specific vector (prefix) while keeping the rest of the PLM frozen, thereby controlling the domain of generation. Another representative work is PPLM by Dathathri et al., which uses gradients from an attribute discriminant model to steer the text generation.

In this work, we focus specifically on the task of Sequential Controlled Text Generation (SCTG), recently formalized by Spangher et al. (2023). In SCTG, a model is provided with an input prompt and a sequence of control codes, and the output is a text sequence comprising multiple sentences. Each control code specifies the desired content or style of the corresponding output sentence, enabling control over the inter-sentence structure of the generated text. The task of SCTG is different from the conventional CTG tasks, which focuses on controlling isolated local attributes at a time. However,

SCTG tackles a more intricate challenge. The generation conditions not only on the discourse of the current sentence or paragraph but also on previous text and contextual discourse structure to maintain contextual coherence throughout the articles.

Previous works relevant to this task include Liu et al. (2022), who proposed a plug-and-play guided decoding method that predicts content plans to control the generation process accordingly. For coherent text generation that considers discourse, Bosselut et al. (2018) modeled discourse structure as cross-sentence ordering. Furthermore, Spangher et al. (2023) introduced a pipeline method that improves discourse through guided generation and an overall editing process.

## 3 Methodology

### 3.1 Overview

We propose a novel framework called Instruct-SCTG (Instruction Sequential Control Text Generation) to incorporate discourse structure into generated articles, by decomposing the generation process into a series of sub-tasks. Each sub-task is designed to generate a single specific text section, such as a main event section or journalist evaluation section, based on the given discourse sequence. In this section, we explain our framework in details and how we design the SFT instruction for our generator LM. Additionally, we introduce an automatic metric that measures the adherence of the discourse structure.

### 3.2 Instruct-SCTG Framework

**Task Formulation.** The goal of SCTG is to generate a coherent article represented by a sequence of linguistic units, e.g. sentences,  $\mathbf{x} = \{x_1, \dots, x_{|\mathbf{x}|}\}$ . Each unit  $x_i$  is denoted as  $x_i = \{x_{i,1}, \dots, x_{i,|x_i|}\}$ , where  $x_{i,j}$  is the  $j$ -th token of  $x_i$ . In the formulation of SCTG, we assume that the input information  $\mathcal{I}$ , such as news headlines or recipe title and ingredients, and discourse structure are provided. The discourse structure is represented as a control code sequence  $\mathbf{c} = \{c_1, \dots, c_s\}$ , where each code denotes the expected discourse role for its corresponding unit  $x_s$ . Hence, the objective of generation is to model the conditional distribution of the document  $\mathbf{x}$ , expressed by Equation 1.

$$P(\mathbf{x}|\mathbf{c}, \mathcal{I}) = \prod_{i=1}^s p(x_i | \mathbf{x}_{<i}, \mathbf{c}, \mathcal{I}) \quad (1)$$

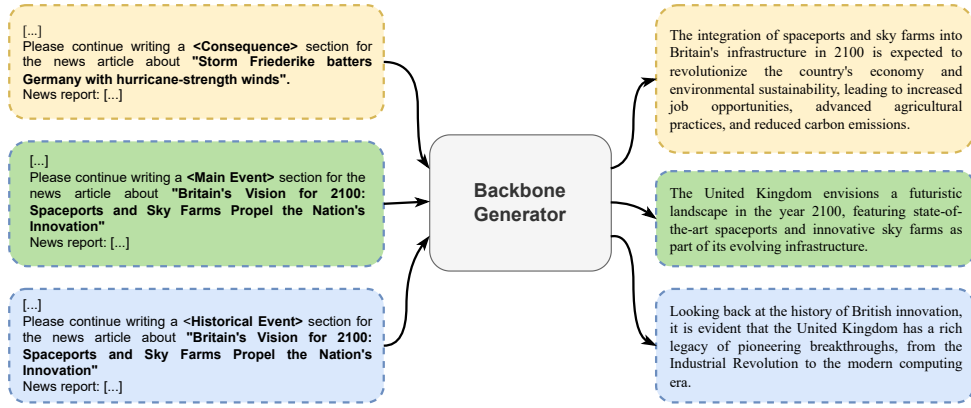


Figure 2: Overview of the instruction tuning of the backbone generator for the Instruct-SCTG.

**Our Framework.** We decompose the document-level conditional distribution into a series of unit-level sub-tasks. During each iteration, we instruct the backbone generator to continue writing for the current linguistic unit according to the specified control code. We can use either a fine-tuned LM with task-specific instructions or a zero-shot large LM as the backbone generator. The control code  $c$ , or discourse roles, are predefined categories based on a specific discourse schema designed for different domains and tasks. In Section 3.3, we explain the design of our SFT instructions.

### 3.3 Task-specific Instruction Tuning

To prepare the backbone generator for our sequential framework, we design task-specific instructions for fine-tuning LMs. As shown in Figure 2, our approach segments articles into sentences or paragraphs. We then create instruction–paragraph pairs as the Supervised Fine-tuning data.

In this section, we also explore the impact of different instruction designs on the resulting fine-tuned generator. The instructions, as shown in the example in Table 1, consists of three main components: (i) discourse context, (ii) input information and (iii) textual context.

In the discourse context, we specifically explore the influence of various facets of contextual discourse information on the generator’s control performance. While exposure to extensive discourse context offers more information, it can potentially introduce additional noise to the current generation process. Previous research (Spangher et al., 2022) employed three levels of discourse dependency assumptions (local, past-aware and Full-sequence) when setting up the discriminator in their post-processing controlling algorithm. In contrast, in this work, we include diverse levels of discourse

context in our instructions. This variation enables us to simulate the those dependency approximations, such that we can directly condition the text generation process on them.

**Local discourse.** If we assume the generation of the current linguistic unit only depends on its corresponding discourse role, but not the contextual discourse structure, the conditional distribution Equation 1 can be simplified as below.

$$P(\mathbf{x}|\mathbf{c}, \mathcal{I}) \approx \prod_{i=1}^{|\mathbf{x}|} \prod_{j=1}^{|\mathbf{x}_i|} p(x_{i,j}|x_{i,<j}, \mathbf{x}_{<i}, c_i, \mathcal{I})$$

**Past-aware.** If we relax the complete independence assumption and allow the previous discourse structure to influence the generation of the current sentence, the Equation 1 will be simplified as below. The discourse context in the instruction template includes only previous discourse sequence but not the future.

$$P(\mathbf{x}|\mathbf{c}, \mathcal{I}) \approx \prod_{i=1}^{|\mathbf{x}|} \prod_{j=1}^{|\mathbf{x}_i|} p(x_{i,j}|x_{i,<j}, \mathbf{x}_{<i}, c_{\leq i}, \mathcal{I})$$

**Full-structure.** If we make no assumption and provide the full discourse structure, the Equation 1 should be expressed as below. Articles generated under different discourse information exposure are compared to determine the optimal instruction template. Experimental results are presented in Section 5.3.

$$P(\mathbf{x}|\mathbf{c}, \mathcal{I}) = \prod_{i=1}^{|\mathbf{x}|} \prod_{j=1}^{|\mathbf{x}_i|} p(x_{i,j}|x_{i,<j}, \mathbf{x}_{<i}, \mathbf{c}, \mathcal{I})$$

In the input information section, we specify the input prompt of the overall generation task and the

---

**Instruction template**

---

The previous discourse structure is :

<Main Event> <Main Event>

The future discourse structure is:

<Journalist Evaluation> <Anecdotal Event> [...]

Please continue writing a <Consequence> section for the news article about "Storm Friederike batters Germany with hurricane-strength winds"

News report:

The United Kingdom envisions a futuristic landscape in the year 2100, featuring state-of-the-art spaceports and innovative sky farms as part of its evolving infrastructure. [...]

---

Table 1: Instruction example. The **discourse context**, **current discourse role** and **headline** are dynamically adjusted based on the context and position of the current sentence. The **textual context** is all previous text before the current target sentence.

discourse role of the current generation unit. For example, in Figure 2, the instruction shown is for generating news report, where the input prompt is the news headline. In the case of the recipe domain, dish title and ingredient list serve as the input, populating the corresponding template.

In the final component, we incorporate textual context to guide the generator in continuing writing the current text segment. During SFT, preceding segments of the article up to the current target one are aggregated to form the previous text, while for inference, all previously generated texts are used.

### 3.4 Zero-shooting LLMs

While we fine-tuned LMs with above-mentioned instructions as the backbone generators of our sequential framework. However, we also explored the option of using zero-shot prompting LLMs with minor modifications to the instruction template. Specifically, for the SFT paradigm, we fine-tuned Flan-T5-base (Chung et al., 2022) and GPT-2 base (Radford et al.). In the case of the zero-shot setup, we opted GPT-3.5-turbo and Flan-T5-xxl. These models have exhibited strong performance in general tasks but are either expensive or not readily available for further training.

To enhance the LLMs’ comprehension of the discourse schema, we introduced a natural language definition of the target discourse role at the beginning of the instruction template. Further details on the discourse definition are listed in Section A.4. In Section 5.3, we present the results achieved using backbone generators under both fine-tuned and

zero-shot paradigms. The results demonstrate the effectiveness and applicability of our framework across different settings.

### 3.5 Measuring the discourse structure

Intuitively, for texts of a certain genre, they tend to follow similar discourse sequences while allowing for some degree of local flexibility. In other words, the distributions of discourse roles in similar areas of the articles are expected to be roughly similar. For instance, in news reports, it is common to have a sentence introducing the main event or consequence at the beginning to quickly capture readers’ attention, but the exact position may vary. In Figure 3, we present the disparity between the discourse distributions of the articles generated by the zero-shot LLM and the reference texts written by humans is evident.

Therefore, to measure the positional difference between the discourse distributions in a fuzzy manner, we introduce the Positional Divergence  $D_{pos}$  as an automatic metric. Equation 2 demonstrates the calculation of the Positional Divergence.

$$D_{pos} = \frac{1}{N} \sum_{n=1}^N D_{KL}(p^n(r) || q^n(r)) \quad (2)$$

Here,  $p^n(r)$  represents the distribution of discourse role  $r$  for the reference data in the  $n$ -th position bin and  $q^n(r)$  represents the distribution for the generated articles. To compute this metric, we firstly segment the reference and generated articles from the evaluation set into  $N$  bins based on their relative positions in the articles. Then, for each bin  $n$ , we calculate the KL divergence  $D_{KL}(p^n(r) || q^n(r))$  between the discourse distributions with add-one smoothing to avoid zero probabilities.

Because the divergence is calculated based on their relative positions in the articles, it mitigates the impact of variations in segmentation styles or the total number of sentences, which cannot be solved by simply calculating the exact match rate. We further elaborate the difference and show that our positional divergence has high correlation with human evaluations in Section A.5. We note that, for this metric, a discourse role classifier is required to label the generated articles.

## 4 Dataset and Schema

In this work, we demonstrate the application of our framework in two representative domains: News

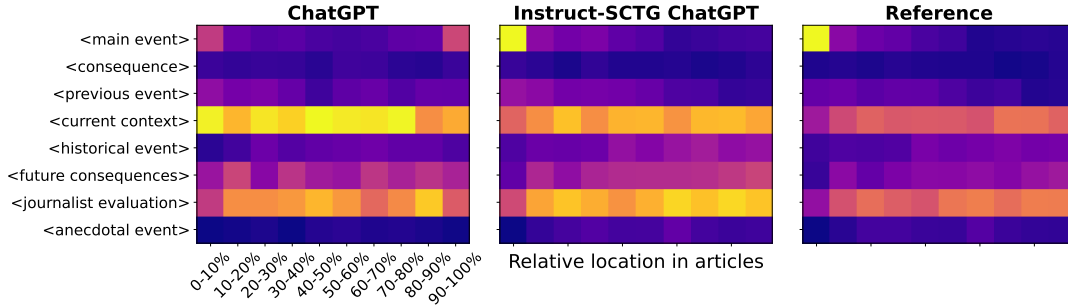


Figure 3: Comparison of discourse distributions at each relative position within news articles. The x-axis represents the relative position from the beginning to the end (0-9), while the y-axis represents different discourse roles based on the news schema of Choubey et al. (2020). Our framework (mid) demonstrates closer discourse distributions to the human-written articles (right), compared with the vanilla baseline (left).

and Recipe. News generation is considered an open-ended task, where there is no fixed predefined answer, allowing more room for creative variations. Whereas Recipe generation is regarded as a closed-ended task, where there exists a correct reference recipe for a given input title. To create the training data, we segment articles into sentences and label them with the assistance of discourse role classifiers.

For news domain, we adopt an existing theory of functional discourse schema proposed by Van Dijk (1988, 2013), which defines a discourse schema based on eight types of relations between each sentence and the main event. A recent News Discourse dataset (Choubey et al., 2020) is manually annotated following the functional discourse schema, which contains 802 documents spanning over four domains and three media sources. We utilize the training set of this dataset to train our discourse role classifier and the test set for evaluating the performance of our framework. In addition, we label the Kaggle All-The-News dataset using our trained discourse role classifier, creating silver-labelled data. Our backbone generators are fine-tuned on the All-The-News training set and evaluated on the News Discourse test set and All-The-News validation set.

For the domain of Recipe, we adopt the discourse schema proposed by Liu et al. (2022) which includes seven discourse roles based on cooking actions specifically designed for recipes. We reimplement their discourse role classifier trained on a subset of the Recipe1M+ validation set (Marin et al., 2019), where the discourse annotations are generated using a rule-based system. We apply this classifier to the remaining Recipe1M+ dataset to generate the silver discourse labels. The fine-tuning of backbone generators for the Recipe do-

main is performed on the Recipe1M+ training set, and the evaluation is conducted on the Recipe1M+ test set. Before using these datasets, we apply pre-processing and filtering based on specified conditions, as elaborated in Appendix A.6. For evaluation, we randomly sample 200 examples from each evaluation set to assess the performance of our framework and the baseline models, and the results are reported in Table 2 and 4.

## 5 Experiments

### 5.1 Implementation Details

In the news domain, the Flan-T5-base backbone generator is trained on the Kaggle All-The-News pre-processed training set for 200k steps, using a batch size of 4. For recipe domain, training is conducted on the processed Recipe1M+ training set for 100k steps with a batch size of 8. Both generators are optimized using the Adam optimizer (Kingma and Ba, 2014) with a learning rate of  $3e - 5$  and an L2 decay rate of 0.05. As for the zero-shot backbone generators, we employ the GPT-3.5-turbo (ChatGPT) and Flan-T5-xxl. During inference, a temperature value of 0.7 is set for news generation and 0.2 for recipes. For news generation, we utilize the top-p sampling method with a value of  $p = 0.8$ , while for recipe generation, we employ beam search decoding with a beam size of 5.

Regarding the Flan-T5, it has limits on the maximum sequence length for both input and output. Therefore, we truncate the input textual context from the beginning to ensure the instruction prompt does not exceed 1024 tokens. The maximum output length is set at 256 tokens.

For the discourse role classifiers, we fine-tune a DistilBERT model (Sanh et al., 2019) using the News Discourse training set for the news domain

Model		Kaggle All-the-news						News Discourse					
		Fluency			Structure			Fluency			Structure		
		PPL. $\downarrow$	R-L $\uparrow$	C-F $\uparrow$	Acc. $\uparrow$	Pos. $\downarrow$	C-S $\uparrow$	PPL. $\downarrow$	R-L $\uparrow$	C-F $\uparrow$	Acc. $\uparrow$	Pos. $\downarrow$	C-S $\uparrow$
F-T	GPT2 <sub>Base</sub>	87.8	21.1	2.7	25.3	0.31	3.0	91.2	19.7	3.1	20.2	0.36	2.2
	FT5 <sub>Base</sub>	100.4	21.4	2.4	24.9	0.37	2.7	108.1	20.4	2.9	21.3	0.32	2.4
CTG	GPT2 <sub>Base</sub>	80.3	22.8	3.0	46.8	0.16	3.2	86.3	20.4	3.2	44.9	0.18	2.9
Z-S	GPT3.5	5.9	22.7	4.7	35.2	0.19	4.0	<b>4.9</b>	21.6	4.4	32.5	0.21	4.4
	FT5 <sub>XXL</sub>	10.0	22.4	4.2	30.6	0.20	3.9	11.8	21.0	4.3	36.1	0.22	4.2
<b>I-SCTG</b>													
F-T	FT5 <sub>Base-L</sub>	78.6	22.2	3.1	60.0	0.10	3.6	74.5	21.8	3.3	63.2	0.13	3.6
	FT5 <sub>Base-P</sub>	63.5	<b>23.5</b>	3.2	<b>63.5</b>	<b>0.08</b>	3.7	65.1	21.1	3.5	<b>67.6</b>	0.10	3.2
	FT5 <sub>Base-F</sub>	65.1	22.4	3.5	61.4	0.11	3.5	67.9	20.9	3.3	65.7	<b>0.09</b>	3.5
Z-S	GPT3.5-P	<b>5.7</b>	22.4	<b>4.8</b>	48.5	0.16	<b>4.8</b>	6.7	22.5	<b>4.5</b>	40.0	0.17	<b>4.7</b>
	FT5 <sub>XXL-P</sub>	9.1	22.6	4.4	42.1	0.18	4.3	10.8	<b>23.1</b>	4.2	42.5	0.19	4.4

Table 2: Results of automatic evaluations conducted on the News domain. The top half shows the outcomes for three types of baseline methods, while the bottom half for various model settings within our Instruct-SCTG (I-SCTG) framework. In the table, "L" denotes the setting for local-discourse, "P" for past-aware, and "F" for full-structure. Our framework shows better ability in controlling the discourse structure of the generated text. For fine-tuned backbone generators, our framework also achieves better surface fluency.

and the recipe1M+ training set for recipes. Both classifiers are trained for 10k steps with a batch sizes of 32. The remaining hyper-parameters are the same with the settings of the backbone generators. The DistilBERT model, being relatively lightweight, demonstrates promising performance as discussed in Appendix A.1.

## 5.2 Experimental setup

**Metrics** We assess our framework from two main perspectives: Surface fluency and adherence to the discourse structure. To measure the surface fluency, we utilize established metrics such as BLEU (B) (Papineni et al., 2002), ROUGE-L R-L (Lin, 2004) and perplexity (PPL.) by another language model OPT-2.7B (Zhang et al., 2022). As for discourse structure control, we measure the exact match accuracy (Acc.), which is the average percentage of matched discourse sequences between the generated text and the reference. Additionally, we use the previously described positional discourse divergence (Pos.) with the number of bins  $N = 10$ .

Traditional automatic metrics often struggle to capture inter-sentence coherence, especially in open-ended generation tasks. Following a recent work (Kocmi and Federmann, 2023), we employ ChatGPT to perform evaluation on both textual fluency (C-F) and structural coherence (C-S) with a scale from 1 to 5. Furthermore, we also perform human evaluations on these aspects by hiring three

native English speakers. They evaluate a randomly selected subset of 100 examples for each evaluation dataset, producing ratings on a scale of 1 to 5. Detailed information about the evaluation prompts for ChatGPT and the setup for human evaluation can be found in Appendix A.2 and A.3.

**Baselines** We evaluate our framework against three types of baselines: 1) Vanilla Fine-tuned LMs (F-T): We fine-tune a GPT-2-base and a Flan-T5-base using only input headlines and reference text pairs from the All-The-News and Recipe1M+ training sets, without incorporating discourse information. We employ the top-k sampling decoding method with a value of  $k = 5$ . 2) Controlled Text Generation methods (CTG): We compare against the approach proposed by Liu et al. (2022), which utilizes a discourse classifier to guide the decoding of a fine-tuned GPT-2-base backbone decoder. 3) Zero-shot large language models (Z-S): We experiment with the GPT-3.5-turbo and Flan-T5-xxl models, prompting them only with the input but no discourse information. The remaining hyper-parameters remain consistent with our framework.

## 5.3 Results

### 5.3.1 News articles

Experimental results were obtained using a randomly selected subset of 200 samples for each dataset. Table 2 displays the averaged experimental results over 5 runs with different random seeds for news generation using various methods.

Model	Fluency	Coherence
FT5 <sub>Base</sub> -FT	2.4	2.5
GPT3.5-ZS	<b>4.5</b>	4.0
FT5 <sub>Base</sub> -P	2.5	3.5
GPT3.5-P	4.3	<b>4.2</b>

Table 3: Results of human evaluations on the News Discourse test set comparing baselines with our Instruct-SCTG framework. Our framework demonstrates improved structural coherence while maintaining a comparable level of surface fluency.

The results demonstrate that our framework outperforms all baseline models on surface fluency and structural coherence metrics when using fine-tuned backbone generators. Among the different contextual discourse information settings, past-aware exhibits better performance. This could be attributed to the fact that subsequent discourse structures might not provide informative enough guidance and could distract the attentions from the more important current discourse roles. When employing zero-shot generators, our framework only utilizes past-aware discourse structure setup to minimize the computational cost. Although vanilla zero-shot LLMs achieve satisfactory surface fluency, our framework can still further enhance the structural coherence of the generated text.

In terms of human evaluations, our framework is compared to two representative baseline models, and the results are presented in Table 3. The human evaluations align with the findings from automatic metrics, confirming that our Instruct-SCTG framework can effectively control the generation process to adhere to the provided discourse structure, resulting in improved structural coherence, while maintaining comparable surface fluency.

### 5.3.2 Recipes

Having the same experiment setup as the news domain, we present the results in Table 4. We observe similar trend with the results on news datasets: Our framework improves the structure coherence for both types of generators, while only fine-tuned generators exhibit better surface fluency. This can be attributed to fact that the recipes generated by latest large-scale LLMs already achieve satisfactory fluency, leaving limited room for further improvement. By applying our framework, the order of actions can be adjusted to better align with the input discourse sequence, while the fluency level remains comparable due to the strong generation capabilities of LLMs. On the other hand, for the fine-tuned

		Recipe1M+				
Model		Fluency			Structure	
		B $\uparrow$	PPL $\downarrow$	R-L $\uparrow$	Acc. $\uparrow$	Pos. $\downarrow$
F-T	GPT2 <sub>Base</sub>	13.1	28.1	38.0	29.3	0.36
	FT5 <sub>Base</sub>	12.7	27.4	37.3	27.9	0.41
CTG	GPT2 <sub>Base</sub>	15.8	26.8	39.1	50.8	0.14
Z-S	GPT3.5	<b>19.2</b>	7.7	<b>44.5</b>	35.2	0.25
	FT5 <sub>XXL</sub>	17.7	9.5	43.2	32.3	0.27
<b>I-SCTG</b>						
F-T	FT5 <sub>Base</sub> -L	16.5	19.8	40.3	66.0	0.10
	FT5 <sub>Base</sub> -P	16.3	24.0	40.5	<b>68.3</b>	<b>0.08</b>
	FT5 <sub>Base</sub> -F	15.8	23.2	39.8	67.5	<b>0.08</b>
Z-S	GPT3.5-P	19.0	<b>6.1</b>	44.2	47.6	0.15
	FT5 <sub>XXL</sub> -P	18.1	8.2	43.5	49.2	0.15

Table 4: Automatic evaluation results for the Recipe domain. Our framework exhibits excellent performance in controlling discourse structure. Improvements in textual fluency are observed when applied our framework to the fine-tuned generators.

generators, incorporating more natural discourse structures can effectively enhance fluency.

## 6 Conclusion

In this work, we address the task of controlling the discourse structure during the generation process. We propose a sequential framework, the Instruct-SCTG, which decomposes article generation into sentence-level tasks. Our framework effectively leverages supervised fine-tuned LMs or zero-shot LLMs as backbone generators to produce structurally more coherent text. We also propose the automatic metric, positional discourse divergence, measuring the discrepancy in discourse distributions across relative positions within the articles. Extensive evaluations demonstrate that our framework can effectively leverage instruction-following LMs to align the discourse structures and achieve SOTA performance on SCTG tasks in both News and Recipe domains.

## Limitations

**Hallucination of news content** In our experimental setup, our primary focus is on controlling the discourse structure of the generated text, rather than the content itself. Consequently, there is a potential for hallucination or the generation of inaccurate information. We acknowledge that in the domain of news reports, the presence of unfactual content can pose problems for readers, as it may compromise the credibility and reliability of the



582	generated articles.		
583	<b>Length limitations</b> News articles are typically		
584	lengthy, but current LLMs often have constraints		
585	on maximum input or output token length. We ac-		
586	knowledge that the truncation method employed		
587	in our study may not be optimal, and alternative		
588	approaches for encoding/decoding extra-long arti-		
589	cles could be explored to capture more contextual		
590	information.		
591	<b>Granularity of discourse annotations</b> When ap-		
592	plying our framework on the zero-shot backbone		
593	generators, we observe instances of local repeti-		
594	tion where consecutive sentences conveyed similar		
595	meanings. This may be attributed to the LLMs’		
596	differing understanding of the granularity of dis-		
597	course structure compared to the reference annota-		
598	tions. LLM-generated articles tend to have fewer		
599	sentences, resulting in shorter discourse sequences.		
600	We recognize that this issue could potentially be		
601	improved by employing more suitable granularity		
602	when annotating the discourse labels.		
603	<b>Data leakage in LLMs</b> Modern LLMs use enor-		
604	mous corpora during pre-training stage, some of		
605	which may not be publicly disclosed. News data,		
606	in particular, has a tendency to be easily acces-		
607	sible, because for an event there might be multi-		
608	ple source of reporting, which makes them easily		
609	scraped for the pre-training. As a result, experi-		
610	ments conducted on news datasets may not be as		
611	indicative as before due to potential data leakage		
612	concerns.		
613	<b>References</b>		
614	Antoine Bosselut, Asli Celikyilmaz, Xiaodong He, Jian-		
615	feng Gao, Po-Sen Huang, and Yejin Choi. 2018.		
616	Discourse-aware neural rewards for coherent text gen-		
617	eration. In <i>Proceedings of the 2018 Conference of</i>		
618	<i>the North American Chapter of the Association for</i>		
619	<i>Computational Linguistics: Human Language Tech-</i>		
620	<i>nologies, Volume 1 (Long Papers)</i> , pages 173–184.		
621	Tom Brown, Benjamin Mann, Nick Ryder, Melanie		
622	Subbiah, Jared D Kaplan, Prafulla Dhariwal, Arvind		
623	Neelakantan, Pranav Shyam, Girish Sastry, Amanda		
624	Askeel, et al. 2020. Language models are few-shot		
625	learners. <i>Advances in neural information processing</i>		
626	<i>systems</i> , 33:1877–1901.		
627	Prafulla Kumar Choubey, Aaron Lee, Ruihong Huang,		
628	and Lu Wang. 2020. <a href="#">Discourse as a function of event:</a>		
629	<a href="#">Profiling discourse structure in news articles around</a>		
630	<a href="#">the main event.</a> In <i>Proceedings of the 58th Annual</i>		
	<i>Meeting of the Association for Computational Lin-</i>	631	
	<i>guistics</i> , pages 5374–5386, Online. Association for	632	
	Computational Linguistics.	633	
	Paul F Christiano, Jan Leike, Tom Brown, Miljan Mar-	634	
	tic, Shane Legg, and Dario Amodei. 2017. Deep	635	
	reinforcement learning from human preferences. <i>Ad-</i>	636	
	<i>vances in neural information processing systems</i> , 30.	637	
	Hyung Won Chung, Le Hou, Shayne Longpre, Barret	638	
	Zoph, Yi Tay, William Fedus, Yunxuan Li, Xuezhi	639	
	Wang, Mostafa Dehghani, Siddhartha Brahma, Al-	640	
	bert Webson, Shixiang Shane Gu, Zhuyun Dai,	641	
	Mirac Suzgun, Xinyun Chen, Aakanksha Chowdh-	642	
	ery, Alex Castro-Ros, Marie Pellat, Kevin Robinson,	643	
	Dasha Valter, Sharan Narang, Gaurav Mishra, Adams	644	
	Yu, Vincent Zhao, Yanping Huang, Andrew Dai,	645	
	Hongkun Yu, Slav Petrov, Ed H. Chi, Jeff Dean, Ja-	646	
	cob Devlin, Adam Roberts, Denny Zhou, Quoc V. Le,	647	
	and Jason Wei. 2022. <a href="#">Scaling instruction-finetuned</a>	648	
	<a href="#">language models.</a>	649	
	Sumanth Dathathri, Andrea Madotto, Janice Lan, Jane	650	
	Hung, Eric Frank, Piero Molino, Jason Yosinski, and	651	
	Rosanne Liu. Plug and play language models: A	652	
	simple approach to controlled text generation. In <i>In-</i>	653	
	<i>ternational Conference on Learning Representations.</i>	654	
	Diederik P Kingma and Jimmy Ba. 2014. Adam: A	655	
	method for stochastic optimization. <i>arXiv preprint</i>	656	
	<i>arXiv:1412.6980.</i>	657	
	Tom Kocmi and Christian Federmann. 2023. Large	658	
	language models are state-of-the-art evaluators of	659	
	translation quality. <i>arXiv preprint arXiv:2302.14520.</i>	660	
	Xiang Lisa Li and Percy Liang. 2021. Prefix-tuning:	661	
	Optimizing continuous prompts for generation. In	662	
	<i>Proceedings of the 59th Annual Meeting of the Asso-</i>	663	
	<i>ciation for Computational Linguistics and the 11th</i>	664	
	<i>International Joint Conference on Natural Language</i>	665	
	<i>Processing (Volume 1: Long Papers)</i> , pages 4582–	666	
	4597.	667	
	Elizabeth DuRoss Liddy. 1991. The discourse-level	668	
	structure of empirical abstracts: An exploratory study.	669	
	<i>Information Processing &amp; Management</i> , 27(1):55–	670	
	81.	671	
	Chin-Yew Lin. 2004. <a href="#">ROUGE: A package for auto-</a>	672	
	<a href="#">matic evaluation of summaries.</a> In <i>Text Summariza-</i>	673	
	<i>tion Branches Out</i> , pages 74–81, Barcelona, Spain.	674	
	Association for Computational Linguistics.	675	
	Yinhong Liu, Yixuan Su, Ehsan Shareghi, and Nigel	676	
	Collier. 2022. <a href="#">Plug-and-play recipe generation with</a>	677	
	<a href="#">content planning.</a> In <i>Proceedings of the 2nd Work-</i>	678	
	<i>shop on Natural Language Generation, Evaluation,</i>	679	
	<i>and Metrics (GEM)</i> , pages 223–234, Abu Dhabi,	680	
	United Arab Emirates (Hybrid). Association for Com-	681	
	putational Linguistics.	682	
	William C Mann and Sandra A Thompson. 1988.	683	
	Rhetorical structure theory: Toward a functional the-	684	
	ory of text organization. <i>Text-interdisciplinary Jour-</i>	685	
	<i>nal for the Study of Discourse</i> , 8(3):243–281.	686	

687	Javier Marin, Aritro Biswas, Ferda Ofli, Nicholas Hynes, Amaia Salvador, Yusuf Aytar, Ingmar Weber, and Antonio Torralba. 2019. Recipe1m+: A dataset for learning cross-modal embeddings for cooking recipes and food images. <i>IEEE Trans. Pattern Anal. Mach. Intell.</i>	Rohan Taori, Ishaan Gulrajani, Tianyi Zhang, Yann Dubois, Xuechen Li, Carlos Guestrin, Percy Liang, and Tatsunori B. Hashimoto. 2023. Stanford alpaca: An instruction-following llama model. <a href="https://github.com/tatsu-lab/stanford_alpaca">https://github.com/tatsu-lab/stanford_alpaca</a> .	742
688			743
689			744
690			745
691			746
692			
693	Yoko Mizuta, Anna Korhonen, Tony Mullen, and Nigel Collier. 2006. Zone analysis in biology articles as a basis for information extraction. <i>International journal of medical informatics</i> , 75(6):468–487.	Hugo Touvron, Thibaut Lavril, Gautier Izacard, Xavier Martinet, Marie-Anne Lachaux, Timothée Lacroix, Baptiste Rozière, Naman Goyal, Eric Hambro, Faisal Azhar, Aurelien Rodriguez, Armand Joulin, Edouard Grave, and Guillaume Lample. 2023. <i>Llama: Open and efficient foundation language models</i> .	747
694			748
695			749
696			750
697	Long Ouyang, Jeffrey Wu, Xu Jiang, Diogo Almeida, Carroll Wainwright, Pamela Mishkin, Chong Zhang, Sandhini Agarwal, Katarina Slama, Alex Ray, et al. 2022. Training language models to follow instructions with human feedback. <i>Advances in Neural Information Processing Systems</i> , 35:27730–27744.	Teun A Van Dijk. 1988. News analysis. <i>Case Studies of International and National News in the Press</i> . New Jersey: Lawrence.	751
698			752
699			753
700			754
701			755
702		Teun A Van Dijk. 2013. <i>News as discourse</i> . Routledge.	756
703	Kishore Papineni, Salim Roukos, Todd Ward, and Wei-Jing Zhu. 2002. <a href="#">Bleu: a method for automatic evaluation of machine translation</a> . In <i>Proceedings of the 40th Annual Meeting of the Association for Computational Linguistics</i> , pages 311–318, Philadelphia, Pennsylvania, USA. Association for Computational Linguistics.	Alex Wang, Amanpreet Singh, Julian Michael, Felix Hill, Omer Levy, and Samuel Bowman. 2018. <a href="#">GLUE: A multi-task benchmark and analysis platform for natural language understanding</a> . In <i>Proceedings of the 2018 EMNLP Workshop BlackboxNLP: Analyzing and Interpreting Neural Networks for NLP</i> , pages 353–355, Brussels, Belgium. Association for Computational Linguistics.	757
704			758
705			759
706			760
707			761
708			762
709			763
710	Rashmi Prasad, Nikhil Dinesh, Alan Lee, Eleni Milt-sakaki, Livio Robaldo, Aravind K Joshi, and Bonnie L Webber. 2008. The penn discourse treebank 2.0. In <i>LREC</i> .	Yizhong Wang, Yeganeh Kordi, Swaroop Mishra, Alisa Liu, Noah A. Smith, Daniel Khashabi, and Hannaneh Hajishirzi. 2022. Self-instruct: Aligning language model with self generated instructions.	764
711			765
712			766
713			767
714	Alec Radford, Jeffrey Wu, Rewon Child, David Luan, Dario Amodei, Ilya Sutskever, et al. Language models are unsupervised multitask learners.	Fangyuan Xu, Junyi Jessie Li, and Eunsol Choi. 2022. How do we answer complex questions: Discourse structure of long-form answers. <i>arXiv preprint arXiv:2203.11048</i> .	768
715			769
716			770
717	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. <i>The Journal of Machine Learning Research</i> , 21(1):5485–5551.	Susan Zhang, Stephen Roller, Naman Goyal, Mikel Artetxe, Moya Chen, Shuohui Chen, Christopher Dewan, Mona Diab, Xian Li, Xi Victoria Lin, et al. 2022. Opt: Open pre-trained transformer language models. <i>arXiv preprint arXiv:2205.01068</i> .	771
718			772
719			773
720			774
721			775
722			776
723	Victor Sanh, Lysandre Debut, Julien Chaumond, and Thomas Wolf. 2019. Distilbert, a distilled version of bert: smaller, faster, cheaper and lighter. <i>arXiv preprint arXiv:1910.01108</i> .		777
724			
725			
726			
727	Alexander Spangher, Xinyu Hua, Yao Ming, and Nanyun Peng. 2023. <a href="#">Sequentially controlled text generation</a> . <i>arXiv preprint arXiv:2301.02299</i> .	<b>A Appendix</b>	778
728			
729			
730	Alexander Spangher, Yao Ming, Xinyu Hua, and Nanyun Peng. 2022. <a href="#">Sequentially controlled text generation</a> . In <i>Findings of the Association for Computational Linguistics: EMNLP 2022</i> , pages 6848–6866, Abu Dhabi, United Arab Emirates. Association for Computational Linguistics.	<b>A.1 Discourse Classifier Results</b>	779
731			
732			
733			
734			
735			
736	Nisan Stiennon, Long Ouyang, Jeffrey Wu, Daniel Ziegler, Ryan Lowe, Chelsea Voss, Alec Radford, Dario Amodei, and Paul F Christiano. 2020. Learning to summarize with human feedback. <i>Advances in Neural Information Processing Systems</i> , 33:3008–3021.	For the News domain, the discourse role classifier is trained on the News Discourse training set and evaluated on the validation set using human-annotated gold labels. The classifier achieves an accuracy of 67%.	780
737			781
738			782
739			783
740			784
741			785
		In the Recipe domain, the discourse role classifier is trained on the Recipe1M+ training set and evaluated on the validation set using silver annotations generated by the rule-based system proposed by Liu et al. (2022). The classifier achieves an accuracy of 92%.	786
			787
			788
			789
			790

791  
792  
793  
794  
795  
796  
797  
798  
799  
800  
801  
  
802  
803  
804  
805  
806  
807  
808  
809  
810  
811  
812  
813  
814  
815  
816  
817  
818  
819  
820  
821  
822  
823  
824  
825  
  
826  
827  
828  
829  
830  
831  
832  
833  
834  
835  
836  
837  
838  
839

## A.2 ChatGPT Evaluation Templates

We use the following instruction to prompt the ChatGPT to rate the textual fluency **C-F** and structural coherence **C-S** of the generated texts.

“You are a helpful virtual journalist. Please rate the textual fluency of the following news report with a score from 1 to 5. Only return the value:”

“You are a helpful virtual journalist. Please rate the structural coherence and the discourse structure quality of the following new report with a score from 1 to 5. Only return the value:”

## A.3 Human Evaluation Guidance Questions

Please rate the following article from two aspects: 1) Textual fluency and 2) structural coherence with score 1 to 5. When evaluating the article, please consider the following guidance.

- **Introduction and lead:** Does the article have a clear and engaging introduction that effectively presents the main topic and captures the reader’s attention?
- **Structure organizatio:** Do the sections and paragraphs follow a clear structure that contributes to the overall understanding of the topic? Are the paragraphs well-structured, with clear topic sentences and appropriate supporting details? Do the paragraphs transition smoothly, maintaining a consistent flow of ideas?
- **Clarity and precision:** Is the language clear, concise, and precise? Are the ideas expressed in a way that is easy to understand for the target audience?
- **Use of evidence and sources:** Are relevant sources and evidence used to support the article’s claims and arguments?

## A.4 Discourse Schema

The definition of the discourse schema we used for news articles:

- **Main Event:** The major subject of the news article.
- **Consequence:** An event or phenomenon that is caused by the main event.
- **Previous Event:** A specific event that occurred shortly before the main event.
- **Current Context:** The general context or world state immediately preceding the main event.
- **Historical Event:** An event occurring much earlier than the main event.

- **Future Consequences:** An analytical insight into future consequences or projections.
- **Journalist Evaluation:** A summary, opinion or comment made by the journalist.
- **Anecdotal Event:** An event that is uncertain and cannot be verified. The primary purpose is to provide more emotional resonance to the main event.

The definition of the discourse schema we used for recipes:

- **Pre-processing** means the preparations of ingredients or cooker.
- **Mixing** includes actions of combining one or more ingredients together.
- **Transferring** is for the actions of moving or transferring food or intermediate food to a specific place.
- **Cooking** represents the actual cooking actions, which could vary drastically across different recipes.
- **Post-processing** usually refers to the following up actions after the ‘cooking’ stage, such as ‘cooling down’, ‘garnish’.
- **Final** refers to the last few actions before serving the food or the serving action itself.
- **General** includes the rest of actions which cannot be classified into the above categories.

## A.5 Further details on Positional Divergence

**Metric Necessity.** We clarify two main practical benefits of our proposed metric:

- For open-ended generation tasks, **it is common for the generated text to have different length (different total number of sentences) or different paragraph layout (different number of sentences for each paragraph) as compared to reference text.** However, these variations do not necessarily mean a substantial deviation in discourse structure. To address this, our proposed positional divergence focuses only on comparing discourse role distributions based on the corresponding relative positions. The continual labels merging strategy couldn’t provide a correct paragraph segmentation due to the aforementioned discontinuity.
- The discourse labels for **the existing dataset don’t usually have multi-sentence continuity, either because the labels are noisy or the flexible nature of the text from open-ended domains.** For instance, below we show the

	$\rho(\text{Acc.}, \text{H.C.})$	$\rho(\text{Pos.}, \text{H.C.})$
FT5 <sub>base</sub> -FT	0.19	0.32
FT5 <sub>base</sub> -P	0.28	0.36
GPT3.5-ZS	0.26	0.33
GPT3.5-P	0.24	0.36

Table 5: The correlations between Human Coherence (H.C) and Exact Match (Acc.) and between H.C. and Positional Divergence. Our proposed metric has shown better correlation with human evaluation.

- Having total number of words over 300 or below 50.
- Duplicate recipes.

929  
930  
931  
932

discourse labels for the sentences in the first paragraph of the Number 18 datapoint of the News Discourse dataset test set: ['main', 'previous\_event', 'main', 'journalist\_evaluation', 'main', 'main', 'main', 'main', 'main', 'consequence']. While the main role of the paragraph is to describe the <main event>, sentences within it might be assigned different role labels such as <evaluation> or <consequence>. In such cases, a simplistic strategy like merging continual sentences cannot effectively handle the evaluation unless guided by a sophisticated merging policy.

**Metric Effectiveness.** We conducted supplementary evaluations to further justify the effectiveness of our metric. We compare the correlation of our metric and the exact match rate to the human evaluation results. In Table 5, we show correlations on the same 100 examples from the News Discourse dataset as shown in Table 3. The results show that **our positional divergence has generally higher correlations than the exact match.**

## A.6 Data Preprocessing

For Kaggle All-The-News, we filtered the dataset based on the following conditions:

- Containing special characters: @, [, +.
- Having total number of words over 800 or below 100.
- Containing random comments.
- Containing more than two reports.

Then we pre-process the data by

- Removing extra space.
- Removing reporting source.
- Removing journalist names.
- Removing emoji.

For Recipe1M+, we filter it based on the following conditions:

- Containing irrelevant information, such as ad-