# Topic Shift Detection in Chinese Dialogues: Corpus and Benchmark

Jiangyi Lin[1], Yaxin Fan[1], Feng Jiang[2], Xiaomin Chu[1], and Peifeng Li[⋆1]

[1] School of Computer Science and Technology, Soochow University, Suzhou, China
{jylin, yxfansuda}@stu.suda.edu.cn, {xmchu, pfli}@suda.edu.cn
[2] School of Data Science, The Chinese University of Hong Kong, Shenzhen, China
jeffreyjiang@cuhk.edu.cn

**Abstract.** Dialogue topic shift detection is to detect whether an ongoing topic has shifted or should shift in a dialogue, which can be divided into two categories, i.e., response-known task and response-unknown task. Currently, only a few investigated the latter, because it is still a challenge to predict the topic shift without the response information. In this paper, we first annotate a Chinese Natural Topic Dialogue (CNTD) corpus consisting of 1308 dialogues to fill the gap in the Chinese natural conversation topic corpus. And then we focus on the response-unknown task and propose a teacher-student framework based on hierarchical contrastive learning to predict the topic shift without the response. Specifically, the response at high-level teacher-student is introduced to build the contrastive learning between the response and the context, while the label contrastive learning is constructed at low-level student. The experimental results on our Chinese CNTD and English TIAGE show the effectiveness of our proposed model.

**Keywords:** Dialogue topic shift detection · Hierarchical contrastive learning · Chinese natural topic dialogues corpus.

## 1 Introduction

Dialogue topic shift detection is to detect whether a dialogue's utterance has shifted in the topic, which can help the dialog system to change the topic and guide the dialogue actively. Although dialog topic shift detection is a new task, it has become a hotspot due to its remarkable benefit to many downstream tasks, such as response generation [1] and reading comprehension [2,3], and can help those real-time applications produce on-topic or topic-shift responses which perform well in dialogue scenarios [4,5,6].

The task of dialogue topic shift detection can be divided into two lines, i.e., response-known task and response-unknown task, as shown in Fig. 1. The former can gain the response information and obtain a better result, while the latter is the opposite. Moreover, both of them are not accessible to future information. This is the biggest difference from the task of text topic segmentation, in which

---

[⋆] Corresponding author.

all the basic utterances are visible to each other. That is, those existing topic segmentation models cannot be applied to dialogue topic shift detection since it depends on the response and its subsequent utterances heavily. Therefore, it is more difficult to discern differences between utterances in the task of dialogue topic shift detection. Due to the absence of future utterances, dialogue topic shift detection is still a challenging task.

In this paper, we focus on the response-unknown task of topic shift detection in Chinese dialogues. There are two issues in the response-unknown task of topic shift detection in Chinese dialogues, i.e., lack of annotated corpus in Chinese and how to predict the response.
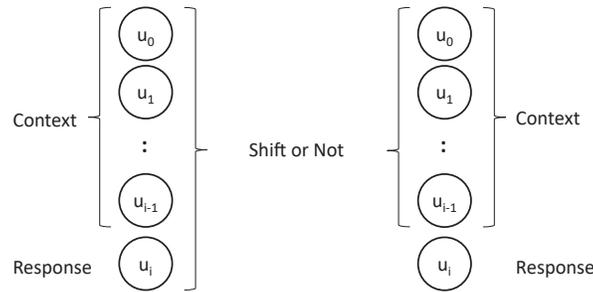


**Fig. 1.** Two lines of dialogue topic shift detection tasks to detect whether it exists topic shift between the utterances $u_{i-1}$ and $u_i$, where the response-known task (left) can use the response $u_i$, while the response-unknown task (right) can be regarded as topic shift prediction without the response $u_i$.

There are only a few publicly dialogue topic shift corpus available and most of them are provided for the segmentation task, which does not satisfy natural conversation. Xie et al. [7] provided a detailed definition of the dialogue topic shift detection task, and annotated an English dialogue topics corpus TIAGE. Although it can fill the gap in the corpus of English conversation topics, its scale is still too small. In Chinese, Xu et al. [8] annotated a Chinese dialogue topic corpus. However, due to its small size and poor quality, this is detrimental to the further research and development of Chinese dialogue topic shift tasks. To fill the gap in the Chinese natural dialogue topic corpus, we first annotated a Chinese Natural Topic Dialogue (CNTD) corpus which consists of 1308 dialogues with high quality.

Xie et al. [7] also established a benchmark for this response-unknown task based on the T5 model [9] and this benchmark only used the context to predict topic shift and performed poorly due to the lack of the response information. Thus, it is more challenging to predict the topic shift in natural dialogue without useful response information.

The teacher-student framework has been used widely to obtain information that is not available to the model [1]. To solve the issue of the lack of response information, we propose a teacher-student framework to introduce the response

information. The teacher can obtain the response information, and the student can learn the response information from the teacher through knowledge distillation. To facilitate knowledge transfer, the student mimics the teacher on every layer instead of just the top layer, which alleviates the delayed supervised signal problem using hierarchical semantic information in the teacher [10].

Besides, we construct hierarchical contrastive learning in which we consider the teacher-student as high-level and the student as low-level. At high-level, we build an information simulation loss between the context and the response to improve the semantic information of the student model with more reliable predictive information. At low-level, we design a semantic coherence-aware loss to better distinguish the different shift cases and produce more reliable prediction results.

Finally, the experimental results on our Chinese CNTD and the English TIAGE show that our proposed model outperforms the baselines. The contributions of this paper are as follows.

- We manually annotate a corpus with 1308 dialogues based on NaturalConv to fill the gaps in the Chinese natural dialogues topic corpus.
- We propose a teacher-student framework to learn the response information for the topic shift detection task.
- We introduce hierarchical contrastive learning to further improve performance.
- The experimental results both on the CNTD and TIAGE datasets show that our model outperforms the baselines.

## 2   Related Work

### 2.1   Corpus

Previous studies explored the dialogue topic tasks and published the annotated topic dialogue corpus. For English, Xie et al. [7] annotated the TIAGE consisting of 500 dialogues with 7861 turns based on PersonaChat [11]. Xu et al. [8] built a dataset including 711 dialogues by joining dialogues from existing multi-turn dialogue datasets: MultiWOZ Corpus [12], and Stanford Dialog Dataset [13]. Both corpora are either small or limited to a particular domain, and neither applies to the study of the natural dialogue domain.

For Chinese, Xu et al. [8] annotated a dataset including 505 phone records of customer service on banking consultation. However, this corpus is likewise restricted to a few specialized domains while natural dialogues are more complicated. Natural dialogues have a range of topic shift scenarios, unrestricted topics, and more free colloquialisms in the utterances. The above corpus is insufficient to fill the gap in the Chinese natural dialogue topic corpus.

### 2.2   Topic Shift Detection in Dialogues

The task of dialogue topic shift detection is also in its initial stage and only a few studies focused on this task. As we mentioned in Introduction, topic segmentation is a similar task. Hence, we first introduce the related work of topic

segmentation in dialogues. Due to the lack of training data on dialogue, early approaches of dialogue topic segmentation usually adopted an unsupervised approach using word co-occurrence statistics [14] or sentence topic distributions [15] to measure sentence similarity between turns to achieve detection of thematic or semantic changes. Recently, with the availability of large-scale corpora sampled from Wikipedia, supervised methods for monologic topic segmentation have grown rapidly by using partial tokens as ground-truth segmentation boundaries, especially neural-based methods [16,17,18]. These supervised solutions are favored by researchers due to their more robust performance and efficiency.

Dialogue topic shift detection is strongly different from dialogue topic segmentation. For the dialogue topic detection task, Xie et al. [7] proposed a detailed definition with two lines: the response-known task considering both the context and the response, and the response-unknown task considering the context only. However, methods based solely on the context are still scarce. Only Xie et al. [7] predict the topic shift or not based on the T5 model. Sun et al. [19] introduce structural and semantic information to help the model detect topic shifts in online discussions, which is similar to response-known task. It is imperative to address the dialogue topic shift detection. In general, the dialogue topic shift detection task is still a challenge, as it can only rely on the context information of the dialogue. In this paper, we solve the lack of response information by utilizing knowledge distillation and hierarchical contrastive learning.

## 3   Corpus

The existing corpus of Chinese dialogue topic detection [8] is small and does not satisfy natural conversation. Although the English dialogue topic corpora can be converted into Chinese by machine translation, they lack natural conversation colloquiality and are small in size. Therefore, we annotate a Chinese dialogue topic detection corpus CNTD based on NaturalConv dataset [20].

In this section, we show our annotation guidelines and outline the reasons for our selection of corpus sources, as well as the manual annotation procedure and data statistics. We also analyze the topic shift distribution in CNTD.

### 3.1   Strengths

Each dialogue in our corpus has a piece of news as a base document, which is not available in other corpus and can be used as additional information for further research and expansion. The news is from six domains, which brings our conversations closer to natural dialogue. Besides, the speakers in our corpus are not restricted in any way, which also makes it closer to natural dialogues. In addition, we annotated the fine-grained dialogues topics, refer to Section 3.2. Fine-grained labels are beneficial to promote further research on dialogue topics.

Compared with the existing Chinese topic corpus annotated by Xu et al. [8], the dialogues in our corpus do not have meaningless and repetitive turns. Also, the corpus is more than twice the size of the other corpus. In addition, the news in the corpus can be studied as additional information for the dialogues.

### 3.2 Annotation Guidelines

Following the annotation guidelines in TIAGE[7], we distinguish each dialogue turns whether changed the topic compared with the context. The response of a speaker to the dialogue context usually falls into one of the following cases in dialogues where the examples can be found in Table 1.

**Table 1.** Different scenarios of response in dialogues.

| | |
|---|---|
| B | "Recently playing handheld games, "dunking master" to look back on the period of watching anime." |
| A | "I see this game so many platforms are pushing ah, but I have not played." |
| B | "Yeah, it's pretty fun, you can go play it when the time comes. There are anime episodes."→ not a topic shift |

(a) Commenting on the previous context.

| | |
|---|---|
| A | "What grade is your child in?" |
| B | "He's a freshman." → not a topic shift |

(b) Question Answering.

| | |
|---|---|
| A | "The Laval Cup is about to start, and the European team has two kings, Federer and Nadal." |
| B | "I know Federer, he is one of the best in the tennis world."→ not a topic shift |

(c) Developing The Conversation to Sub-topics

| | |
|---|---|
| B | "Haha, so what do you usually like to do sports ah? Do you usually go out for a run?" |
| A | "Rarely, usually lying at home watching TV, running and so on is to see their fat can not pretend to look." |
| B | "I also, then what TV are you watching lately? Have you been watching "Elite Lawyers"?"→ topic shift |

(d) Introducing A Relevant But Different Topic.

| | |
|---|---|
| B | "This movie I saw crying, Iron Man died, really moved." |
| A | "Yes, the special effects of this movie are very good." |
| B | "Who is your favorite actor?"→topic shift |

(e) Completely Changing The Topic.

- Commenting on the previous context: The response is a comment on what is said by the speaker previously;
- Question answering: The response is an answer to the question that comes from the speaker previously;
- Developing the dialogue to sub-topics: The response develops to a sub-topic compared to the context;
- Introducing a relevant but different topic: The response introduces a relevant but different topic compared to the context;

– Completely changing the topic: The response completely changes the topic compared to the context.

Among them, we uniformly identify the two cases of greeting and farewell specific to CNTD as the topic shift.

### 3.3 Data Source

We chose the NaturalConv dataset [20] as the source corpus, which contains about 400K utterances and 19.9K dialogues in multiple domains. It is designed to collect a multi-turn document grounded dialogue dataset with scenario and naturalness properties of dialogue.

We consider NaturalConv as a promising dataset for dialogue topic detection for the following reasons: 1) NaturalConv is much closer to human-like dialogue with the natural property, including a full and natural setting such as scenario assumption, free topic extension, greetings, etc.; 2) NaturalConv contains about 400K utterances and 19.9K dialogues in multiple domains; 3) The average turn number of this corpus is 20, and longer dialogue contexts tend to exhibit a flow with more topics; 4) The corpus has almost no restrictions or assumptions about the speakers, e.g., no explicit goal is proposed [21].

### 3.4 Annotation Process

We have three annotators for coarse-grained annotations and two for fine-grained annotations. Both annotations are divided into three stages as follows.

**Co-annotation Stage** First, for coarse-grained annotations, we draw a total of 100 dialogues from each domain of the NaturalConv dataset proportionally for a total of 2014 dialogue turns. In this stage, three annotators are asked to discuss every 20 dialogues they annotated, and each annotator is asked to give a reason for the annotation during the discussion. Finally, the Kappa value of all annotators for coarse-grained annotations at this stage is 0.7426. In addition, we annotated the fine-grained information based on the results of the complete coarse-grained annotations. Two annotators annotated the same 150 dialogues and discussed them several times for consistency. Finally, the kappa value of all annotators for fine-grained annotations at this stage is 0.9032. These kappa values confirm that our annotators already have sufficient annotation capabilities for independent annotation, as well as the high quality of our corpus.

**Independent-annotation Stage** We ensured the quality of each annotator's annotation and judging criteria before starting the second phase of annotation. For both granularity annotations, we randomly assign the dialogues drawn from each domain to each annotator for independent annotation. At this stage, we annotate 1208 dialogues for coarse-grained annotations and 1158 dialogues for fine-grained annotations.

**Semi-automatic Rechecking Stage** Finally, we use a semi-automatic rechecking process to ensure that the corpus is still of high quality. On the one hand, we automatically format the dialogues with annotations to detect formatting

**Table 2.** Category and proportion of the corpus.

| Category | Train | Val. | Test | Sum. |
|---|---|---|---|---|
| Health(8%) | 85 | 11 | 11 | 107 |
| Education(16%) | 167 | 22 | 21 | 210 |
| Technology(17%) | 176 | 22 | 22 | 220 |
| Sports(33%) | 347 | 45 | 46 | 438 |
| Games(8%) | 86 | 11 | 11 | 108 |
| Entertainment(17%) | 180 | 23 | 22 | 225 |
| Total | 1041 | 134 | 133 | 1308 |

**Table 3.** Details of CNTD.

| | Min. | Max. | Avg. |
|---|---|---|---|
| Dialogue Turns | 20 | 26 | 20.1 |
| Utterance Words | 1 | 141 | 21.0 |
| Dialogue Words | 194 | 888 | 421.7 |
| Dialogue Topics | 2 | 9 | 5.2 |
| Topic Turns | 1 | 17 | 4.2 |

problems caused by manual annotation. On the other hand, we automatically match the related news to each dialogue and check that the topic attributes are consistent with the dialogue to rule out any possible errors.

### 3.5 Annotation Results

Due to the limited time, we randomly select 1308 dialogues from the Natural-Conv dataset and annotate them with four annotators. Finally, we construct a Chinese natural topic dialogues corpus containing 26K dialogue turns.

As shown in Table 2, we randomly split them into 1041 train, 134 validation, and 133 test dialogues respectively, according to the percentage of different categories. In addition, we show the details of CNTD in Table 3, which shows that our corpus has enough topics and long turns which is suitable for dialogue topic detection. Finally, there are the statistics of our fine-grained labels, as shown in Table 4.

We count the number of dialogues with different numbers of topics, as shown in Fig. 2. On another side, we count the distribution of topic shift signals in dialogues, shown in Fig. 3. We can see there are a total of 21 turns and three peaks of topic shift signals, which occur in $2^{nd}$, $4^{th}$, and $18^{th}$ turns, respectively. The reason is that the dialogue in our corpus usually starts with a greeting and

**Table 4.** Statistics for fine-grained labels.

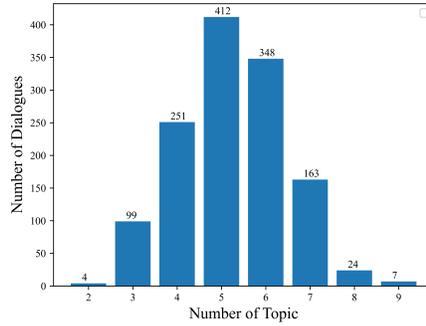| Fine-grained labels | Count |
|---|---|
| Commenting on the previous context | 15091 |
| Question answering | 3505 |
| Developing the dialogue to sub-topics | 857 |
| Introducing a relevant but different topic | 3106 |
| Completely changing the topic | 2439 |

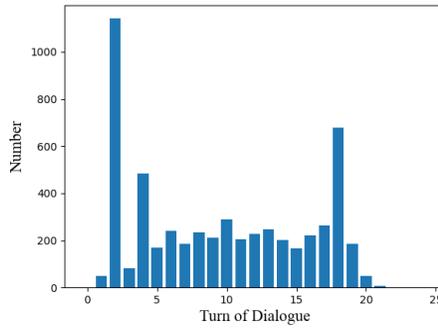**Fig. 2.** Number of dialogues with different numbers of topics.



**Fig. 3.** Topic shift distribution of CNTD.

ends with a farewell, which leads to more topic shifts at the beginning and end of the dialogues. In addition, the NaturalConv corpus gives a piece of news as the base document of the dialogue, so there are more frequent transitions from news to derived topics, leading to the third highest peak in $4^{th}$ turn. However, we think this is consistent with a natural dialogue scenario because people often talk about recent news after daily greetings.

## 4   Model

The framework of our model is shown in Fig. 4. We propose a teacher-student framework based on Hierarchical Contrastive Learning, which contains two parts: knowledge distillation and hierarchical contrastive learning which consists of two different contrastive learning.

### 4.1   Knowledge Distillation

Existing studies cannot effectively predict topic shifts due to the lack of future information. To address this problem, we introduced a teacher-student frame-
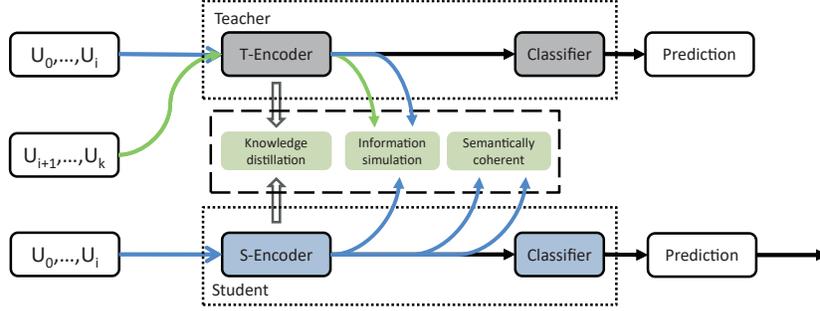
**Fig. 4.** Architecture of the model. The figure contains a teacher encoder, which can obtain the context and the response, and a student encoder, which has only context as input, where $u_{i+1}$ denotes the response in the current dialogue sample and $k$ denotes the length of the dialogue. In addition, the student is prompted to learn from the teacher through knowledge distillation and hierarchical contrastive learning which contains information simulation loss, and semantic coherence-aware loss.

work for dialogue topic shift detection. The student side learns an implicit way of topic prediction from the teacher side through knowledge distillation.

In this framework, we employ a pre-trained model as an encoder on both sides to obtain semantic representations of the dialogues. Besides, we share the encoder weights on the teacher when encoding the context and the response information. For instance, the representation of $[CLS]$ of the last hidden layer state is taken as a dialogue-level representation $H$ as follows.

$$H_{C,i}^T = E_T \left( X_{C,i} \right) \tag{1}$$

$$H_{P,i}^S = E_S \left( X_{P,i} \right) \tag{2}$$

where $H \in R^{N \times d}, i \in I = \{1, 2, ..., N\}$ denotes the index of samples in the batch, $C \in \{P, F, A\}$ where $P$ represents the context, $F$ represents the response, $A$ represents the full dialog. $E$ represents the encoder, $T$ represents the teacher, and $S$ represents the student.

On the teacher side, we connect the dialogue-level representations and then feed them into the linear layer to obtain the final detection results while on the student as follows.

$$Z_i^T = W_T[H_{P,i}; H_{F,i}; H_{A,i}] \tag{3}$$

$$Z_i^S = W_S H_{P,i} \tag{4}$$

where $W_T$ and $W_S$ denote the linear layers on the teacher and student sides, respectively.

In addition to calculating the cross-entropy loss of the final detection results, we establish the mean squared error loss between each hidden layer of the teacher and student encoders.

## 4.2   Hierarchical Contrastive Learning

We consider that the response information learned by knowledge distillation solely is insufficient. And the unbalanced proportion of shift or not in the dialogue corpus makes the model perform poorly in distinguishing different shift cases. Therefore, we propose hierarchical contrast learning, consisting of information simulation loss at the high level and semantic coherence-aware loss at the low level, respectively. Both losses are based on contrast learning, but the former is to strengthen the learning of features and the latter is to alleviate the imbalance of labels.

For information simulation loss, this loss enables active learning for each context representation. And this effectiveness has been demonstrated by several works [1,22,23], the incorporation of global information permits local information representation with some predictive information. This helps our encoder obtain a representation with more reliable predictive information.

To alleviate the unbalanced proportion of labels in the dialogues, we propose a semantic coherence-aware loss based on supervised contrast learning (SCL). The main concept of SCL[24,25] is to regard the samples of different categories as positive and negative samples from each other to address the issue of significant quantitative imbalance in the dataset[26]. This loss effectively alleviates the imbalance of shift cases and helps the model further distinguish between different shift cases.

**High-level Information Simulation Loss** We build the information simulation loss on both sides so that the context representations can be mapped to the same high-dimensional space. Thus, it is easier to learn the response information to improve final detection results. The following equation can be used to describe this loss.

$$L_{ISL}^{M} = \sum_{i \in I} \log \frac{exp\left(H_{P,i}^{M} \cdot H_{A,i}^{T}\right)}{\sum_{j \in A(i)} exp\left(H_{P,i}^{M} \cdot H_{A,j}^{T}\right)} \tag{5}$$

where $A(i) = I - \{i\}$ denotes the samples in the current batch other than itself, $P$ denotes the context, $A$ denotes the full dialog, and $M \in \{T, S\}$ where $T$ denotes the teacher, $S$ denotes the student.

**Low-level Semantic Conherent-aware Loss** For a batch with $N$ training samples, a copy of the dialogue's last hidden state $H$ is made to obtain $\overline{H}$ that is considered as the positive, and its gradient is detached. This results in $2N$ samples, then the semantic coherence-aware loss of all samples in a batch can be expressed as follows.

$$U = [H; \overline{H}] \tag{6}$$

$$L_{SCL} = \sum_{i \in I} \frac{-1}{|P(i)|} \sum_{p \in P(i)} \log \frac{exp\left(U_i \cdot U_p\right)}{\sum_{a \in A(i)} exp\left(U_i \cdot U_a\right)} \tag{7}$$

where $U \in R^{2N \times d}, i \in I = \{1, 2, ..., 2N\}$ denotes the index of samples in a batch, and $P(i) = I_{j=i} - \{i\}$ denotes samples of the same category as $i$ but not itself.

### 4.3   Model Training

We train our model in two steps. The teacher is trained first, and its loss consists of two parts: cross-entropy loss between predictions and manual annotation labels, named $L_{NCE}$. And the information simulation loss for learning response representation is named $L_{ISL}^{T}$. The overall training loss is as follows.

$$L^{T} = L_{NCE} + L_{ISL}^{T} \tag{8}$$

$$L^{S} = L_{NCE} + L_{KD} + L_{ISL}^{S} + L_{SCL} \tag{9}$$

Then, we train the student, which consists of four parts: cross-entropy loss between predictions and manual annotation labels named $L_{NCE}$, knowledge distillation between each hidden layer of both sides named $L_{KD}$, information simulation loss for the dialogue context named $L_{ISL}^{S}$, and the semantic coherence-aware loss for different shift cases named $L_{SCL}$. The above equation represents its overall training loss. As it proves to be arduous to fine-tune the weights assigned to the four losses, we ultimately opt for equal weighting across all of them.

## 5   Experiments and Analysis

### 5.1   Experimental Settings

Based on the train/validation/test dataset of CNTD we partitioned in Table 2 and previous work on TIAGE [7], we extract (context, response) pairs from each dialogue as input and the label of response as a target for the response-unknown task. In our experiments, every utterance except the first utterance of the dialogue can be considered as a response. As for evaluation, we report Precision (P), Recall (R), and Micro-F1 scores.

We use BERT as an encoder and fine-tune it during training. For both the TIAGE and CNTD corpus, all pre-trained model parameters are set to default values. We conduct our experiments on NVIDIA GeForce GTX 1080 Ti and NVIDIA GeForce GTX 3090 with batch sizes of 2 and 6 for both CNTD and TIAGE, with the initial learning rates of 2e-5. And we set the epochs of training to 20, and the dropout to 0.5.

For the pre-trained models in the experiment, we apply BERT-base-Chinese and MT5-base to obtain the semantic representation of the dialogues in CNTD, and we apply BERT-base-uncased and T5-base to obtain the semantic representation of the dialogues in TIAGE.

### 5.2   Experimental Results

Dialogue topic shift detection is a new task and there is no complex model available, besides a simple T5 [7] that can be considered as the SOTA model. Since we employ BERT as our encoder and the T5 model is used in TIAGE, we use the pre-trained models of T5 [9] and BERT [27] as baselines. For BERT, we

**Table 5.** Comparison of our model and three baselines on CNTD.

| Model | **P** | **R** | **F1** |
| --- | --- | --- | --- |
| T5 | 27.1 | 46.8 | 34.3 |
| BERT | 55.5 | 43.8 | 48.9 |
| TS | 52.7 | 47.4 | 49.9 |
| Ours | **56.0** | **52.0** | **53.9** |

connect the utterances in the context and separate the last utterance with $[SEP]$. For T5, we also connect utterances in the context and classify the undecidable predicted results to the 'not a topic shift' category.

Table 5 shows the performance comparison between our model and the baselines, in which TS denotes our teacher-student model without the hierarchical comparative learning (HCL) and Ours denotes our final model, i.e., the addition of SCL on the student side based on the addition of ISL on both the teacher and student sides.

It can be found that on CNTD, our model achieves a good improvement and improves both precision and recall in comparison with the baselines. Although T5 does not perform poorly on recall, its precision is inadequate in comparison with BERT, and it is clear that T5 is not effective in predicting topics. In contrast, TS improved by 1.0 in Micro-F1 in comparison with BERT, which confirms that the teacher-student framework is effective in introducing response information. As well, Ours improved by 4.0 in micro-F1 in comparison with TS, and also showed significant improvement in P and R, which fully demonstrates that our HCL can improve the model's ability to discriminate between different topic situations. In particular, our model improves on CNTD by 5.0 in comparison with the best baseline BERT, which shows the effectiveness of our proposed model.

### 5.3   Ablation Study

To verify the effectiveness of the components used in our model, we conduct ablation studies on CTND, and the experimental results are shown in Table 6.

If we remove ISL on the teacher side ($-ISL_S$) or the student side ($-ISL_T$), the performance of the model decreased by 1.5 and 1.3 on the Micro-F1 value, respectively, with the largest decrease after removing the ISL on the student side. Although $-ISL_T$ has the highest precision in predicting topics and lower error probability than $Ours$ and $-ISL_S$. However, it can be seen that adding ISL at both the teacher and student sides can better improve the correct prediction rate. Moreover, if we remove ISL both on the teacher and student side ($-ISL_{TS}$), it achieves a similar performance on Micro-F1, in comparison with $-ISL_S$ and $-ISL_T$. However, it achieves the highest precision (58.8%).

If we remove SCL (-SCL) or HCL (-HCL) from our model, the Micro-F1 value of the models -SCL and -HCL drop from 53.9 to 52.4 (-1.5) and 49.9 (-4.0), respectively. These results show that our Semantic Conherent-aware Loss(SCL), and Hierarchical Contrastive Learning(HCL) are effective for this task, especially HCL.

**Table 6.** Results of our model and its variants on CNTD where ISL, SCL, and HCL refer to the information simulation Loss, semantic conherent-aware Loss, and hierarchical contrastive Learning, respectively, and T, S, and TS refer to teacher side, student side and both of them.

|  | **P** | **R** | **F1** |
|---|---|---|---|
| Ours | 56.0 | 52.0 | **53.9** |
| $-ISL_S$ | 53.6 | 50.8 | 52.2 |
| $-ISL_T$ | 56.1 | 49.3 | 52.6 |
| $-ISL_{TS}$ | **58.8** | 47.0 | 52.3 |
| $-SCL$ | 51.6 | **53.1** | 52.4 |
| $-HCL$ | 52.7 | 47.4 | 49.9 |

**Table 7.** Performance on dialogues with the different number of topics.

| Topic Number | Our Model(F1) | BERT(F1) |
|---|---|---|
| 2 | 100 | 66.7 |
| 3 | 49.0 | 42.6 |
| 4 | 64.2 | 53.7 |
| 5 | 58.5 | 46.3 |
| 6 | 50.3 | 50.7 |
| 7 | 44.4 | 47.5 |
| 8 | 44.4 | 40.0 |
| 9 | 76.9 | 54.5 |

### 5.4   Analysis on Different Angles of Performance

In addition, we explore the performance of the dialogues with different numbers of topics to analyze our model in comparison with BERT, as shown in Table 7. It can be found that our model has a better performance than BERT on dialogues with fewer topics. Our model gets at least a 6% improvement in topic shift prediction on dialogues with 2 to 5 topics and obtains above-average performance. And when the number of topics increases to 9, the performance improves because the conversation length is still about 20 and the topics shift more significantly.

In Table 8, we also investigate the recall of the topic shift detection for various topic turns. Our model is improved for varying degrees across topic turns, with the most significant improvements in turns 7-9. Even in long topic shift cases, our model can obtain an effective boost. However, the performance of our model inevitably decreases compared to short topic shift cases. When there are fewer topic turns, the topic shift situation is simpler, so it is easier to determine. When the length of turns becomes longer and the situation becomes complicated, the topic of long turns has more information so it is easier to identify.

### 5.5   Results on English TIAGE

As shown in Table 9, it can be found that our model also achieves a good improvement on English TIAGE. Although our model is not the best on precision,

**Table 8.** Performances of topic shift with different turn lengths.

| Topic Turns | Our Model(Recall) | BERT(Recall) |
|:-----------:|:-----------------:|:------------:|
| 1-3   | 56.6 | 53.8 |
| 4-6   | 30.3 | 21.4 |
| 7-9   | 28.8 | 11.9 |
| 10-12 | 40.0 | 32.0 |
| 13-17 | 40.0 | 30.0 |

**Table 9.** Results on TIAGE.

|      | TIAGE | | |
|:----:|:--------:|:----:|:----:|
|      | **P**    | **R**  | **F1** |
| T5   | **34.0** | 17.0   | 22.0   |
| BERT | 28.1     | 17.9   | 21.7   |
| TS   | 26.9     | 20.1   | 22.9   |
| Ours | 27.4     | **28.3** | **27.8** |

we obtain the best performance on both recall and Micro-F1 values, especially on micro-F1 with a 5.8% improvement over T5. This proves that our model achieves the best performance both in English and Chinese.

### 5.6   Case Study and Error Analysis

We also conducted a case study. The prediction made by our model, the BERT model on the instance, and the manual labels are shown in Table 10. Compared with the BERT model, it is obvious that our model can accurately anticipate the change of topic in the instances corresponding to the utterances "Yes, that's right.","It is, indeed, should pay attention to it." etc., belonging to the question-answering scenario. However, if you respond "Well, the policy has been implemented in place this time." and "And now we are promoting the development of children's creative and practical skills." etc. belonging to the commenting on the previous context scenario, our model or BERT cannot accurately predict the topic shift in this scenario. This shows that detecting the topic shifts in natural dialogue is still challenging.

We further analyze the errors of the prediction produced in our experiments. Specifically, we analyzed the example to explore whether the error in the results of this example is prevalent in other dialogues. From Table 10, we can find that the wrong predictions at $14^{th}$ and $18^{th}$ turn. We predict "The teaching equipment must be updated, right?" as 'not a topic shift' and "Well, thanks to the government!" as 'topic shift'.

We counted the appearance of many errors, and the errors are mainly divided into two categories. One is for the "Introducing a relevant but different topic" type of utterance. It was predicted that no topic shift occurred due to the lack of information about the future of the conversation. The other is the "commenting on the previous context" category. Since this type of response does not affect the integrity of the previous topic, it is mostly predicted to be a topic shift.

**Table 10.** The results of BERT, Ours, and Human of different turns where "1" indicates that a topic shift has occurred and "0" indicates the opposite. We omit the lines with all 0.

| Turns | BERT | Ours | Human |
|:-----:|:----:|:----:|:-----:|
| 3     | 0    | 1    | 1     |
| 5     | 0    | 1    | 1     |
| 11    | 0    | 1    | 1     |
| 13    | 1    | 0    | 0     |
| 14    | 0    | 0    | 1     |
| 16    | 1    | 0    | 0     |
| 17    | 0    | 1    | 1     |
| 18    | 1    | 1    | 0     |
| 19    | 0    | 1    | 1     |

## 6   Conclusion

Based on the NaturalConv dataset, we create the CNTD dataset with manual annotations, which fill the gap in the Chinese natural dialogues topic corpus. And we propose a teacher-student model based on hierarchical contrastive learning to solve the lack of response information. We introduced response information through a teacher-student framework and constructed information simulation learning in high-level teachers and students and semantic conherent-aware learning in low-level students. The experiment results demonstrate that our model can perform better in dialogue with few topics. However, detecting the long turns topics or the dialogues with more topics remains a complex problem. Our future work will focus on how to better use response information and news information to detect topic shifts in real-time.

## Acknowledgements

## References

1. Shuyang Dai, Guoyin Wang, Sunghyun Park, and Sungjin Lee. Dialogue response generation via contrastive latent representation learning. In *Proceedings of the 3rd Workshop on Natural Language Processing for Conversational AI*, pages 189–197, 2021.
2. Jiaqi Li, Ming Liu, Zihao Zheng, Heng Zhang, Bing Qin, Min-Yen Kan, and Ting Liu. Dadgraph: A discourse-aware dialogue graph neural network for multiparty

dialogue machine reading comprehension. In *2021 International Joint Conference on Neural Networks (IJCNN)*, pages 1–8. IEEE, 2021.

3. Yiyang Li and Hai Zhao. Self-and pseudo-self-supervised prediction of speaker and key-utterance for multi-party dialogue reading comprehension. In *Findings of the Association for Computational Linguistics: EMNLP 2021*, pages 2053–2063, 2021.

4. Asma Ghandeharioun, Judy Hanwen Shen, Natasha Jaques, Craig Ferguson, Noah Jones, Agata Lapedriza, and Rosalind Picard. Approximating interactive human evaluation with self-play for open-domain dialog systems. *Advances in Neural Information Processing Systems*, 32, 2019.

5. Arash Einolghozati, Sonal Gupta, Mrinal Mohit, and Rushin Shah. Improving robustness of task oriented dialog systems. *arXiv preprint arXiv:1911.05153*, 2019.

6. Bing Liu, Gokhan Tür, Dilek Hakkani-Tür, Pararth Shah, and Larry Heck. Dialogue learning with human teaching and feedback in end-to-end trainable task-oriented dialogue systems. In *Proceedings of NAACL-HLT*, pages 2060–2069, 2018.

7. Huiyuan Xie, Zhenghao Liu, Chenyan Xiong, Zhiyuan Liu, and Ann Copestake. Tiage: A benchmark for topic-shift aware dialog modeling. In *Findings of the Association for Computational Linguistics: EMNLP 2021*, pages 1684–1690, 2021.

8. Yi Xu, Hai Zhao, and Zhuosheng Zhang. Topic-aware multi-turn dialogue modeling. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 35, pages 14176–14184, 2021.

9. Colin Raffel, Noam Shazeer, Adam Roberts, Katherine Lee, Sharan Narang, Michael Matena, Yanqi Zhou, Wei Li, Peter J Liu, et al. Exploring the limits of transfer learning with a unified text-to-text transformer. *J. Mach. Learn. Res.*, 21(140):1–67, 2020.

10. Zhuohan Li, Zi Lin, Di He, Fei Tian, Tao Qin, Liwei Wang, and Tie-Yan Liu. Hint-based training for non-autoregressive machine translation. In *Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing (EMNLP-IJCNLP)*, pages 5708–5713, 2019.

11. Saizheng Zhang, Emily Dinan, Jack Urbanek, Arthur Szlam, Douwe Kiela, and Jason Weston. Personalizing dialogue agents: I have a dog, do you have pets too? In *Proceedings of the 56th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 2204–2213, 2018.

12. Paweł Budzianowski, Tsung-Hsien Wen, Bo-Hsiang Tseng, Iñigo Casanueva, Stefan Ultes, Osman Ramadan, and Milica Gasic. Multiwoz-a large-scale multi-domain wizard-of-oz dataset for task-oriented dialogue modelling. In *Proceedings of the 2018 Conference on Empirical Methods in Natural Language Processing*, pages 5016–5026, 2018.

13. Mihail Eric, Lakshmi Krishnan, Francois Charette, and Christopher D Manning. Key-value retrieval networks for task-oriented dialogue. In *Proceedings of the 18th Annual SIGdial Meeting on Discourse and Dialogue*, pages 37–49, 2017.

14. Jacob Eisenstein and Regina Barzilay. Bayesian unsupervised topic segmentation. In *Proceedings of the 2008 Conference on Empirical Methods in Natural Language Processing*, pages 334–343, 2008.

15. Lan Du, Wray Buntine, and Mark Johnson. Topic segmentation with a structured topic model. In *Proceedings of the 2013 conference of the North American chapter of the Association for Computational Linguistics: Human language technologies*, pages 190–200, 2013.

16. Omri Koshorek, Adir Cohen, Noam Mor, Michael Rotman, and Jonathan Berant. Text segmentation as a supervised learning task. In *Proceedings of the 2018 Con-*

ference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 2 (Short Papers), pages 469–473, 2018.

17. Pinkesh Badjatiya, Litton J Kurisinkel, Manish Gupta, and Vasudeva Varma. Attention-based neural text segmentation. In European Conference on Information Retrieval, pages 180–193. Springer, 2018.

18. Sebastian Arnold, Rudolf Schneider, Philippe Cudré-Mauroux, Felix A Gers, and Alexander Löser. Sector: A neural model for coherent topic segmentation and classification. Transactions of the Association for Computational Linguistics, 7:169–184, 2019.

19. Yingcheng Sun and Kenneth Loparo. Topic shift detection in online discussions using structural context. In 2019 IEEE 43rd Annual Computer Software and Applications Conference (COMPSAC), volume 1, pages 948–949. IEEE, 2019.

20. Xiaoyang Wang, Chen Li, Jianqiao Zhao, and Dong Yu. Naturalconv: A chinese dialogue dataset towards multi-turn topic-driven conversation. In Proceedings of the AAAI Conference on Artificial Intelligence, volume 35, pages 14006–14014, 2021.

21. Wenquan Wu, Zhen Guo, Xiangyang Zhou, Hua Wu, Xiyuan Zhang, Rongzhong Lian, and Haifeng Wang. Proactive human-machine conversation with explicit conversation goal. In Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics, pages 3794–3804, 2019.

22. Aaron van den Oord, Yazhe Li, and Oriol Vinyals. Representation learning with contrastive predictive coding. arXiv preprint arXiv:1807.03748, 2018.

23. Shaoxiong Feng, Xuancheng Ren, Hongshen Chen, Bin Sun, Kan Li, and Xu Sun. Regularizing dialogue generation by imitating implicit scenarios. In Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP), pages 6592–6604, 2020.

24. Shimin Li, Hang Yan, and Xipeng Qiu. Contrast and generation make bart a good dialogue emotion recognizer. In Proceedings of the AAAI Conference on Artificial Intelligence, volume 36, pages 11002–11010, 2022.

25. Beliz Gunel, Jingfei Du, Alexis Conneau, and Veselin Stoyanov. Supervised contrastive learning for pre-trained language model fine-tuning. In International Conference on Learning Representations.

26. Yanran Li, Hui Su, Xiaoyu Shen, Wenjie Li, Ziqiang Cao, and Shuzi Niu. Dailydialog: A manually labelled multi-turn dialogue dataset. In Proceedings of the Eighth International Joint Conference on Natural Language Processing (Volume 1: Long Papers), pages 986–995, 2017.

27. Jacob Devlin Ming-Wei Chang Kenton and Lee Kristina Toutanova. Bert: Pre-training of deep bidirectional transformers for language understanding. In Proceedings of NAACL-HLT, pages 4171–4186, 2019.