

Inconsistent dialogue responses and how to recover from them

Anonymous ACL submission

Abstract

001 One critical issue for chat systems is to stay
002 consistent about preferences, opinions, beliefs
003 and facts of itself, which has been shown a
004 difficult problem. In this work, we study meth-
005 ods to assess and bolster utterance consistency
006 of chat systems. A dataset is first developed
007 for studying the inconsistencies, where inconsis-
008 tent dialogue responses, explanations of the
009 inconsistencies, and recovery utterances are au-
010 thored by annotators. This covers the life span
011 of inconsistencies, namely introduction, under-
012 standing, and resolution. Building on this, we
013 introduce a set of tasks centered on dialogue
014 consistency, specifically focused on its detec-
015 tion and resolution. Our experimental findings
016 indicate that our dataset significantly helps the
017 progress in identifying and resolving conversa-
018 tional inconsistencies, and current popular
019 large language models like ChatGPT which are
020 good at resolving inconsistencies however still
021 struggle with detection.¹

022 1 Introduction

023 For years, inconsistencies in human-to-chatbot con-
024 versations have been evident (Dziri et al., 2019;
025 Rashkin et al., 2021; Ji et al., 2023), even in the era
026 of large language models (Mündler et al., 2023).
027 We categorize these inconsistencies as either extrin-
028 sic or intrinsic. *Extrinsic* inconsistencies (Rashkin
029 et al., 2021; Santhanam et al., 2021) arise when
030 there’s a discrepancy between a statement and an
031 external source of information, such as a knowl-
032 edge base. On the other hand, *intrinsic* inconsisten-
033 cies (Dziri et al., 2019; Nie et al., 2021; Zheng et al.,
034 2022) occur within the dialogue itself. These can
035 manifest in two ways: through an intra-utterance
036 contradiction (Zheng et al., 2022), where a single
037 sentence contains conflicting information, or
038 a history contradiction (Nie et al., 2021), where
039 a current statement conflicts with a previous one.

¹The dataset and codebase will be released at <https://url>.

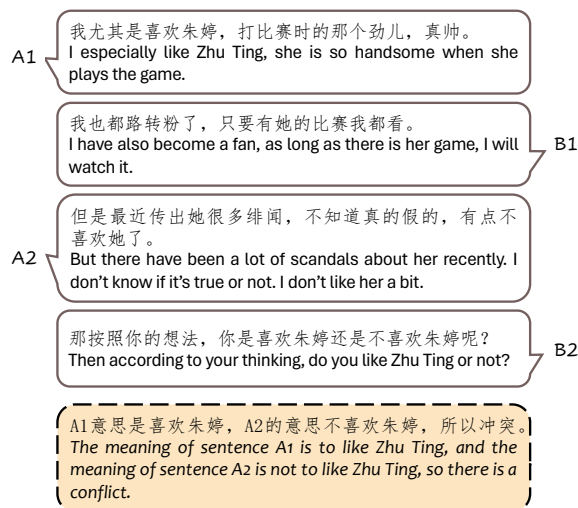


Figure 1: An instance in **CIDER** dataset. $\{A, B\}_x$ denotes the x -th utterance of one of the two speakers (A or B). An inconsistent utterance (A2 in this case), an explanation of the inconsistency (the dotted box), and a clarification response (B2 in this case) are written for each dialogue.

040 Our study particularly addresses history contradic-
041 tions, a persistent challenge in conversational mod-
042 els due to the nature of language modeling: models
043 could forget what they said due to intervening con-
044 text (Roller et al., 2021).

045 Researchers have been actively exploring how
046 to resolve inconsistencies between utterances gen-
047 erated by conversational models in recent years.
048 Li et al. (2020); Rashkin et al. (2021) has made
049 progress in this domain by enhancing the training
050 of these models, incorporating additional features
051 and objectives to bolster self-consistency. Further-
052 more, Lee et al. (2022); Su and Collier (2022) in-
053 troduced innovative decoding algorithms aimed at
054 fostering greater coherence in utterances. These
055 preemptive approaches are designed to mitigate
056 conversational inconsistencies by reducing the like-
057 lihood of generating responses that contradict pre-
058 vious dialogue. However, these approaches cannot

059 resolve the inconsistencies that do occur, possibly
060 from the user or from model errors. Therefore it’s
061 equally important to robustly address inconsisten-
062 cies that do arise. Various remedial techniques
063 have shown promise in other tasks, from grammar
064 error correction (Wu et al., 2023) and moderating
065 inappropriate dialogue content (Zhang et al., 2023),
066 to generating clarifying questions in information
067 retrieval (Zamani et al., 2020a) and conversational
068 question answering (Guo et al., 2021). However,
069 there seems to be a significant gap in the research
070 when it comes to directly addressing inconsisten-
071 cies that do arise between utterances.

072 In this work, we first propose a human-authored
073 dataset with 27,180 dialogues to study the in-
074 consistencies between utterances. At a high
075 level, the dataset, called **CIDER**, covers the
076 whole life span of inConsistencies, encompassing
077 their Introduction, understanding, and Resolution.
078 Specifically, for each dialogue, annotators are first
079 asked to write an utterance with inconsistent con-
080 tent regarding one utterance in the history to con-
081 tinue the conversation (A2 in Figure 1), and then
082 explain why the two utterances are inconsistent
083 with natural language (the dotted box in Figure 1),
084 and finally provide a clarification response to con-
085 tinue the dialogue to resolve the inconsistency² (B2
086 in Figure 1). Given its large collection of inconsis-
087 tent utterances paired with clarifying responses,
088 **CIDER** stands out as a valuable resource for re-
089 searching strategies to mitigate conversational in-
090 consistencies.

091 Utilizing the **CIDER** dataset, we conduct com-
092 prehensive experiments and analyses to study dia-
093 logue inconsistencies. Our findings underscore that
094 **CIDER** can facilitate the development of robust in-
095 consistency checkers compared to models trained
096 on comparable public datasets. In addition, our
097 research indicates that classic models like T5 and
098 BART face challenges in adeptly resolving inconsis-
099 tencies by providing clarifying responses. When
100 assessing the proficiency of large language models
101 (LLMs) in identifying and resolving conversational
102 inconsistencies, we discerned two key points: 1)
103 LLMs, when employed as inconsistency checkers,
104 still leave much to be desired in terms of perfor-
105 mance. 2) In contrast, as resolvers of inconsistency,
106 LLMs exhibit a higher success rate compared to
107 the fully supervised BART resolver.

²The dialogues and annotation in the dataset are in Chinese. We also offer an English version translated by ChatGPT to facilitate research.

2 Related work 108

Consistency checking. Natural Language Infer- 109
ence (NLI) (Hu et al., 2020; Saha et al., 2020) 110
is a task closely related to our work, where a 111
provided hypothesis is evaluated for its logical 112
consistency with a given premise, with both pre- 113
sented in natural language. Within the context of 114
dialogues, Welleck et al. (2019) framed the con- 115
sistency checking in dialogue as NLI and anno- 116
tated binary consistency labels between dialogue- 117
persona or persona-persona sentence pairs from the 118
Persona-Chat dataset (Zhang et al., 2018). Dziri 119
et al. (2019) employed NLI models to assess topic 120
coherence between a current response and the pre- 121
ceding dialogue history. Meanwhile, Shuster et al. 122
(2022) delved into the issue of role confusion, 123
where dialogue systems might inadvertently adopt 124
the identity of the other party involved, and pro- 125
posed a reranker trained with human judgments of 126
identity consistency. The most relevant works are 127
from (Nie et al., 2021) and (Zheng et al., 2022), 128
where they created datasets providing supervision 129
for contradiction detection between conversational 130
sentences. Our work distinguishes itself by provid- 131
ing more extensive annotations, including explana- 132
tions and clarification responses. 133

Consistency resolving in dialogue. To enhance 134
the self-consistency of conversational models, 135
Rashkin et al. (2021) employed controllable fea- 136
tures, steering models towards generating more 137
consistent responses. Lee et al. (2022) introduced 138
factual-nucleus sampling and factuality-enhanced 139
continued training to augment the reliability of lan- 140
guage models during both decoding and training 141
phases. Shuster et al. (2022) encouraged the con- 142
versational models to maintain an identity with the 143
help of a role-playing accuracy classifier. Li et al. 144
(2020) explored unlikelihood training (Welleck 145
et al., 2020) to curb inconsistencies in dialogue. 146
However, given computational constraints, contem- 147
porary conversational models tend to rely predomi- 148
nantly on recent dialogue history when formulating 149
responses. This predisposes them to produce con- 150
tent that may contradict earlier parts of the dialogue, 151
especially distant sections. Generating clarifica- 152
tion questions has emerged as a strategy to address 153
communication breakdowns in dialogues, such as 154
resolving ambiguities in a query during conversa- 155
tional information retrieval (Zamani et al., 2020b) 156
or clarifying ambiguous user questions in conver- 157

sational question answering (Guo et al., 2021) scenarios. In this research, we propose an approach to recover from conversational inconsistencies by generating clarification questions, with the support of the proposed dataset.

Large language models. Recent advancements in AI have been dominated by the rise of large language models, notably ChatGPT (Ouyang et al., 2022), GPT-4 (OpenAI, 2023) and others. They have shown that by scaling up language models, they can be equipped to tackle intricate tasks, such as question answering, machine translation, and numerical reasoning. In this study, leveraging the extensive annotations of our proposed dataset, **CIDER**, we assess these models’ proficiency in detecting and addressing conversational inconsistencies.

3 Data collection

The candidate conversations for annotation are sampled from two open-source conversation datasets: LCCC and NaturalConv. LCCC (Wang et al., 2020) is a large collection of short conversations from the Chinese social media platform Weibo. NaturalConv (Wang et al., 2021) is an annotator-written dataset containing conversations around news items on topics like film and sports. They are different in content and style. LCCC conversations tend to be short in number of turns, and more in the style of daily chitchat. NaturalConv conversations, in contrast, are two to five times longer and contain more serious discussions about events in sports, films, and other areas. 20,000 and 10,000 conversations are sampled from the LCCC and NaturalConv respectively for annotation. When sampling, conversations that are shorter than 4 turns or contain utterances shorter than 5 words are filtered out.

The sampled conversations are generally consistent, therefore the goal of data annotation is to create an alternative conversation that contains inconsistent utterances. To achieve this, we truncate the original conversation to create a common conversation context. For LCCC, the last utterance is truncated for inconsistent continuation writing; for NaturalConv, a random turn between 8 and $l - 4^3$ and the following turns are chosen for truncation, where l is the length of the conversation.

Finally, a specified source turn is sampled from the last turn or the turn before the last. This source

³The last turns of NaturalConv tend to be goodbyes, therefore we choose to truncate before such utterances.

	LCCC			NaturalConv		
	Train	Dev	Test	Train	Dev	Test
# of Convs	14,126	1,883	1,797	7,537	917	920
Ave. Cont. Len.	29.3	28.9	28.9	40.4	40.9	40.5
Ave. Exp. Len.	40.9	40.5	41.0	50.4	50.3	50.3
Ave. Res. Len.	16.2	16.1	16.1	20.3	20.1	20.0

Table 1: Some basic statistics of the annotated datasets. Ave. Cont. Len. is the average continuation length in number of Chinese characters; Ave. Exp. Len. is the average explanation length; Ave. Res. Len. is the average resolution question length. They correspond to the outcome from the three annotation tasks.

turn is designated to be the source of the inconsistency where the following inconsistent continuation needs to form inconsistency with the utterance from the same speaker in this turn.

4 Annotation guidelines

The annotation process has been divided into three different tasks: inconsistent continuation, inconsistency explanation, and inconsistency resolution, which are required to be performed to each candidate conversation by one annotator when given a candidate conversation and a specified source turn.

Inconsistent continuation. The annotator first tries to create a natural continuation of the conversation by providing a possible utterance to the candidate conversation, but forms an inconsistency with the specified source utterance (A2 in Figure 1 is the continuation, and A1 is the source.) The annotators are instructed to write the utterance with contradictory viewpoints, reasoning, and argumentation, instead of providing simple negation to the source utterance. For example, for the specified utterance *I went to the supermarket yesterday.*, the continuation meeting the annotation requirement is *I have been staying home for the past four days, not really wanting to go anywhere*, instead of *I didn’t go to the supermarket yesterday.*

Inconsistency explanation. After writing the continuation of the candidate conversation, the annotator is instructed to write down the rationale behind the created inconsistency (the dashed box in Figure 1). They are asked to follow this template when writing the rationale: *The specified utterance means X, but the continuation utterance means Y, which is in contradiction with X.*, where the utterance meanings should be explicit. In the example above, the explanation one may write is *The specified utterance indicates that I went out of my home*

	Pair-Check						Diag-Check					
	Train		Valid		Test		Train		Valid		Test	
	#Pos	#Neg	#Pos	#Neg	#Pos	#Neg	#Pos	#Neg	#Pos	#Neg	#Pos	#Neg
STANCE	1816	3959	195	446	346	644	1816	3959	195	446	346	644
OCNLI	14837	30601	1639	3409	900	2100	14837	30601	1639	3409	900	2100
CDConv	2623	4373	880	1452	848	1484	2623	4373	880	1452	848	1484
CIDER	21663	53012	2800	6692	2717	6569	21663	21663	2800	2800	2717	2717

Table 2: Dataset statistics for checking tasks.

yesterday, but the continuation utterance means that I didn't go out for many days including yesterday, which is in contradiction with the previous statement.

Inconsistency resolution. Finally, the annotator provides another utterance to expose and question the inconsistency from a different party than the continuation party (B2 in Figure 1). The annotator is asked to write the resolution question naturally with the main purpose being clarifying the situation instead of complaining. They are also asked to try varying how the clarification question is raised, because the most intuitive way is asking by providing a binary choice. The resolution question for the example above is *So were you home yesterday or did you go to the supermarket?*

Twelve examples collected from the two data sources and annotated by the authors were provided to the annotators along with the guidelines, which cover a number of common mistakes that the authors discovered in the trial annotation. The annotation project lasted two months, with six annotators⁴ participating in the project from a commercial annotation provider, who was chosen amongst three providers based on the performance in the trial annotation task. The items for annotation were segmented into batches, each with 3000 conversations. The annotated items are checked first by quality assurance specialists from the annotation provider by batch, and then spot-checked by the authors with the acceptance rate setting at 95%.⁵ Candidate conversations which are not possible to form inconsistencies, such as conversations containing mostly utterances of simple greeting or agreeing,

⁴The chosen provider created a qualification test based on the annotation guidelines for selecting annotators. The annotators with the highest agreement with the authors were then chosen as annotators. They then went through an online training session with the authors to align with the understanding of guidelines from the authors. They were paid twice the local average monthly salary for their contributions.

⁵The spot-check rate is 10%.

are dropped in the annotation process.

5 Data overview

After annotation, 17,806 conversations from LCCC and 9,374 conversations from NaturalConv have valid annotation. They are further split into train, dev and test sets, shown in Table 1. The average continuation and explanation lengths from LCCC conversations are substantially shorter than from NaturalConv, indicating the simple nature of social media conversations. The resolution question lengths are closer than the other lengths, showing that resolution questions tend to be less influenced by context and style.

6 Consistency checking

In this section, we experimentally verify whether the proposed **CIDER** could help the detection of inconsistency in conversation via two task settings: (1) checking the consistency between two sentences (*Pair-Check*); (2) checking the consistency between an utterance and its preceding context (*Diag-Check*). The (in)consistency checker is initialized as RoBERTa-base (Liu et al., 2019) with a linear binary classification head on the top. The input of the encoder for *Pair-Check* is formatted as "[CLS] {sentence 1} [SEP] {sentence 2} [SEP]" while for *Diag-Check*, "[CLS] {context} [SEP] {utterance} [SEP]", where the [CLS] and [SEP] are special tokens.

Baselines. We compare **CIDER** with several related datasets:

- CDConv (Zheng et al., 2022): a dataset with 12K dialogues for conversational contradiction detection. Compared to **CIDER**, CDConv covers another two types of contradiction: intra-sentence contradiction and role confusion. Each dialogue of CDConv contains two turns of utterances between a user and a bot

	<i>STANCE Test</i>			<i>OCNLI Test</i>			<i>CDConv Test (Turn)</i>			<i>CIDER Test (Turn)</i>		
	Pre.	Rec.	F1	Pre.	Rec.	F1	Pre.	Rec.	F1	Pre.	Rec.	F1
C_{STANCE}^{Turn}	72.8	60.4	<u>66.0</u> \uparrow 14.3	37.7	19.4	25.7	38.1	21.3	27.4	37.5	14.4	20.8
C_{OCNLI}^{Turn}	31.6	36.1	33.7	72.9	74.9	<u>73.9</u> \uparrow 10.2	51.3	37.3	43.2	35.7	37.4	36.5 \uparrow 1.4
C_{CDConv}^{Turn}	41.8	8.1	13.6	40.9	15.0	22.0	56.3	72.9	<u>63.5</u> \uparrow 14.7	29.8	42.8	35.1
C_{CIDER}^{Turn}	61.0	44.8	51.7 \uparrow 18.0	30.7	76.2	63.7 \uparrow 38.0	37.7	69.3	48.8 \uparrow 5.6	76.2	69.3	<u>72.6</u> \uparrow 36.1

(a) Performance of *Pair-Check* checkers.

	<i>STANCE Test</i>			<i>OCNLI Test</i>			<i>CDConv Test (Diag)</i>			<i>CIDER Test (Diag)</i>		
	Pre.	Rec.	F1	Pre.	Rec.	F1	Pre.	Rec.	F1	Pre.	Rec.	F1
C_{STANCE}^{Turn}	72.8	60.4	<u>66.0</u> \uparrow 20.4	37.7	19.4	25.7	25.9	4.5	7.6	48.4	21.8	30.0
C_{OCNLI}^{Turn}	31.6	36.1	33.7	72.9	74.9	<u>73.9</u> \uparrow 40.8	46.6	37.6	41.6 \uparrow 18.6	52.5	42.7	47.1 \uparrow 17.1
C_{CDConv}^{Diag}	54.5	8.7	15.0	31.5	16.2	21.4	62.5	60.8	<u>61.7</u> \uparrow 20.1	61.3	8.3	14.6
C_{CIDER}^{Diag}	38.8	55.2	45.6 \uparrow 11.9	33.7	32.4	33.1 \uparrow 7.4	52.7	14.7	23.0	89.4	91.6	<u>90.5</u> \uparrow 43.4

(b) Performance of *Diag-Check* checkers.

Table 3: Performance of the checking tasks. The checker trained on dataset Y for task *X-Check* is denoted as C_Y^X . The best result in each column is in bold. The best F1 score on each dataset is underscored and the points by which it exceeds the second best are shown by \uparrow . The transferring F1 scores on each dataset are in italics and the points by which they exceed the second best transferring score are shown by \uparrow . The performance of C_{STANCE}^{Turn} and C_{OCNLI}^{Turn} on *STANCE Test* and *OCNLI Test* in Table 3b is copied from Table 3a.

	<i>Merge</i>						<i>Pretrain</i>					
	<i>Pair-Check</i>			<i>Diag-Check</i>			<i>Pair-Check</i>			<i>Diag-Check</i>		
	Pre.	Rec.	F1	Pre.	Rec.	F1	Pre.	Rec.	F1	Pre.	Rec.	F1
C_{CIDER}	76.2	69.3	72.6	89.4	91.6	90.5	76.2	69.3	72.6	89.4	91.6	90.5
+CDConv	76.7	72.5	74.6 \uparrow 2.0	90.7	91.9	<u>91.3</u> \uparrow 0.8	76.4	71.1	73.7 \uparrow 1.1	88.4	91.4	89.9 \downarrow 0.6
+OCNLI	70.1	77.4	73.6 \uparrow 1.0	89.8	92.1	90.9 \uparrow 0.4	77.4	70.7	<u>73.9</u> \uparrow 1.3	88.6	93.1	<u>90.8</u> \uparrow 0.3
+STANCE	72.4	77.9	<u>75.1</u> \uparrow 2.5	88.2	92.9	90.5 \uparrow 0.0	76.2	70.3	73.2 \uparrow 0.6	87.3	92.7	89.9 \downarrow 0.6

Table 4: Performance of checkers leveraging extra data on the test set of **CIDER**. The best are in bold. The relative increasing (\uparrow) and decreasing (\downarrow) points are calculated based on the performance of C_{CIDER} .

and annotation of *consistent* or *inconsistent* between the replies of the bot.

- STANCE⁶: a dataset for stance classification of articles of debating topics from online forums, where sentence pairs against each other are marked as *inconsistent* and otherwise *consistent*.
- OCNLI (Hu et al., 2020): a large-scale natural language inference (NLI) dataset, consisting of about 56,000 annotated sentence pairs. We regard sentence pairs with *contradiction* label as *inconsistent* and others as *consistent*.

Implementation details. For **CIDER**, when creating *consistent* training instances of *Pair-Check*, we regard all the utterances in the context of the

⁶www.fudan-disc.com/sharedtask/AIDebater21/tracks.html

same speaker without *inconsistent* label as being consistent with the current response; and when creating the training instances of *Diag-Check*, we drop current response with inconsistency and regard the previous response as being consistent with the context. Table 2 shows the statistics of the datasets for these two checking tasks.

We adopt AdamW (Loshchilov and Hutter, 2019) to optimize models for 50 epochs with a learning rate of 1e-6 and a batch size of 16. We evaluate the model on the validation set at each epoch and keep the one with the best performance with an early stop patience of 3. All the results are averaged over three runs. Our experiments are run on two Nvidia V100 GPUs.

Results for *Pair-Check*. The performance of checkers trained on different datasets for *Pair-*

345 *Check* is demonstrated in Table 3a. For each
 346 checker, we show its performance on all the test
 347 sets of the evaluating datasets.

348 There is a substantial distribution difference be-
 349 tween the datasets with the checker trained on
 350 one dataset performing the best on the correspond-
 351 ing test set. $C_{\text{CIDER}}^{\text{Turn}}$ has the largest exceeding F1
 352 points over the second best, 36.1, indicating that
 353 the checker trained on other datasets is not good at
 354 detecting the consistency in the test set of **CIDER**
 355 and the training set of **CIDER** could provide useful
 356 supervision for it. Moreover, we compare the 0-
 357 shot transfer ability of checkers across the datasets.
 358 Results show that $C_{\text{CIDER}}^{\text{Turn}}$ has the best transfer re-
 359 sults on all the other three datasets, surpassing the
 360 second best by 18.0, 38.0, and 5.6 F1 points, re-
 361 spectively, demonstrating $C_{\text{CIDER}}^{\text{Turn}}$ covering many
 362 similar linguistic phenomena in other datasets. On
 363 the whole, **CIDER provides robust supervision to**
 364 **check whether a pair of sentences are consistent,**
 365 **regardless of they are in a dialogue or not.**

366 **Results for *Diag-Check*.** The performance of the
 367 checkers trained on different datasets for *Diag-*
 368 *Check* is demonstrated in Table 3b. The results
 369 of $C_{\text{CDConv}}^{\text{Diag}}$ and $C_{\text{CIDER}}^{\text{Diag}}$ indicates again the distribu-
 370 tion difference between **CIDER** and CDConv also
 371 being significant for *Diag-Check* task: **CIDER** do
 372 not cover role confusion and intra-sentence contra-
 373 diction these two types of inconsistency while be-
 374 ing much larger than CDConv. In addition, $C_{\text{CIDER}}^{\text{Diag}}$
 375 outperforms $C_{\text{CDConv}}^{\text{Diag}}$ on *STANCE Test* by 30.6 F1
 376 points and on *OCNLI Test* by 11.7 F1 points, which
 377 demonstrates better transferring ability of $C_{\text{CIDER}}^{\text{Diag}}$
 378 to non-conversational scenarios. Therefore, along
 379 with the transferring results in Table 3a, **CIDER**
 380 **offers more transferable patterns for checking**
 381 **consistency, and may be complementary to CD-**
 382 **Conv in the conversational scenarios.** We also
 383 notice that $C_{\text{OCNLI}}^{\text{Turn}}$ is superior to $C_{\text{CIDER}}^{\text{Diag}}$ on *CD-*
 384 *Conv Test (Diag)* and to $C_{\text{CDConv}}^{\text{Diag}}$ on *CIDER Test*
 385 (*Diag*), showing that the knowledge of inconsis-
 386 tency between sentences in OCNLI is also useful
 387 for the inconsistency checking in dialogue.

388 **Role of extra data.** We are interested in whether
 389 other datasets could improve the performance of
 390 C_{CIDER} . We leverage the training data of STANCE,
 391 OCNLI, and CDConv via two ways: 1) directly
 392 merging one of them into the training data of
 393 **CIDER (Merge)**; 2) pretraining the checker on one
 394 of them before training on **CIDER (Pretrain)**.

	<i>Pair-Check</i>			<i>Diag-Check</i>		
	Pre.	Rec.	F1	Pre.	Rec.	F1
C_{CIDER}	76.2	69.3	72.6	89.4	91.6	90.5
ChatGPT	42.0	79.0	54.8	57.2	84.9	68.4
GPT4	49.9	76.2	60.3	68.8	82.1	74.8

Table 5: Performance of LLMs on checking tasks.

395 The results are presented in Table 4. It’s evident
 396 that **incorporating additional data generally en-**
 397 **hances the overall performance of C_{CIDER} ,** The
 398 only exception is that only pretraining on OCNLI
 399 could improve the checker for *Diag-Check* task,
 400 which indicates better supervision signal from OC-
 401 NLI for checking the inconsistency of an utterance.
 402 Compared with pretraining on extra data, directly
 403 merging them is superior, which could be ascribed
 404 to the phenomenon of catastrophic forgetting (Kirk-
 405 patrick et al., 2017) of pretrained models. More-
 406 over, *Pair-Check* generally benefits from the extra
 407 datasets more than *Diag-Check* because most of the
 408 extra datasets are intrinsically designed for check-
 409 ing of sentence pairs and in large quality so models
 410 could learn generalized patterns from them.

411 **LLMs as consistency checker.** We investigated
 412 the potential of large language models (LLMs) to
 413 function as robust consistency checkers. We pre-
 414 examine five human-crafted prompts for each task
 415 using a small-scale test set (50 instances) and se-
 416 lect the best. The prompts applied for the checking
 417 tasks are illustrated in Figure 2. The evaluating
 418 LLMs are ChatGPT and GPT4. As shown in Ta-
 419 ble 5, LLM-based checkers significantly lag behind
 420 the fully supervised C_{CIDER} , indicating that there
 421 is still much room for improvement. Moreover, the
 422 higher performance of GPT4 over ChatGPT under-
 423 scores that larger LLMs possess a better capability
 424 to detect inconsistencies.

<i>Pair-Check</i>	Whether the following two sentences are semantically related and have semantic inconsistencies, please answer "yes" or "no". sentence 1: {sentence 1} sentence 2: {sentence 2}
<i>Diag-Check</i>	Please answer "yes" or "no" if the speaker of the last sentence in the following dialogue contradicts himself, and give an explanation. {dialogue}

Figure 2: Prompts of checking tasks.

Model	Pair-Resolve				Diag-Resolve			
	BLEU	R-1	R-2	R-L	BLEU	R-1	R-2	R-L
#1 T5	26.9	55.3	33.0	52.2	14.8	43.0	20.6	40.4
#2 BART	28.2 \uparrow 1.3	57.2 \uparrow 1.9	34.8 \uparrow 1.8	53.7 \uparrow 1.5	14.9 \uparrow 0.1	43.7 \uparrow 0.7	21.7 \uparrow 1.1	41.0 \uparrow 0.6
#3 T5 _{oracle}	46.2	71.5	53.0	68.3	46.7	71.7	53.2	68.3
#4 BART _{oracle}	49.4 \uparrow 3.2	74.4 \uparrow 2.9	56.2 \uparrow 3.2	70.7 \uparrow 2.4	47.4 \uparrow 0.7	72.4 \uparrow 0.7	53.9 \uparrow 0.7	68.7 \uparrow 0.4
#5 ChatGPT	14.3	45.2	22.2	41.4	5.3	29.8	9.9	26.9
#6 GPT4	10.8	42.7	20.2	38.0	4.1	28.0	9.8	24.2

Table 6: Performance of resolvers on the test set of **CIDER**. The relative increasing (\uparrow) points of BART (BART_{oracle}) are calculated based on the performance of T5 (T5_{oracle}).

7 Consistency resolution

Inconsistent responses of a conversational model could be detected by a consistency checker in advance, avoiding being exposed to users. However, inconsistent responses from a user can not be ignored by chat systems. The existence of inconsistent content may confuse the conversational model and induce undesired responses. Resolving the occurred inconsistency is necessary to maintain a smooth dialogue flow with clear semantics. The proposed **CIDER** dataset contributes to resolving the occurred inconsistency in a dialogue with *clarification responses*, which is a valuable source to train an inconsistency resolution model.

We choose the base version of two representative conditional generative models to initialize the resolver: BART (Lewis et al., 2020) and T5 (Raffel et al., 2020). They both follow an encoder-decoder structure and generate clarification responses in a sequence-to-sequence fashion: the conversational text with inconsistency is fed into the encoder and the clarification response is generated aggressively by the decoder. Like the checking experiments in section 6, we consider two task settings: (1) generating a clarification response for a pair of inconsistent utterances (*Pair-Resolve*); (2) generating a clarification response for a dialogue, of which the current response is inconsistent to the preceding context (*Diag-Resolve*). The input of the encoder for *Pair-Resolve* is formatted as "[CLS] {utterance 1} [SEP] {utterance 2} [SEP]" while for *Diag-Resolve*, "[CLS] {context} [SEP] {response} [SEP]".

Implementation details. We use the same optimization configuration of checkers to train the resolvers, except that a learning rate of $3e-4$ is used for T5. BART and T5 are loaded with pretrained parameters from Zhao et al. (2019) and Shao et al.

(2021), respectively. In decoding, we adopt Nucleus Sampling (Holtzman et al., 2020) with top-0.90 probability mass across the experiments.

Evaluation. We use BLEU (Papineni et al., 2002) and ROUGE (Lin, 2004), including ROUGE-1 (R-1), ROUGE-2 (R-2) and ROUGE-L (R-L), to measure the similarity between the generated text and the ground truth.

Results. According to rows #1 and #2 in Table 6, BART shows better performance in both *Pair-Resolve* and *Diag-Resolve* tasks than T5, indicating the pretrained parameters of BART are more suitable to inconsistency resolving. Meanwhile, the points of *Pair-Resolve* are higher than those of *Diag-Resolve*, which could be ascribed to *Diag-Resolve* being a more difficult task than *Pair-Resolve* because recognizing inconsistent contents between conversational context and a response is harder than between a pair of sentences. We also try appending *explanations* to the input of the encoder to aid the generation process. Specifically, the input becomes "[CLS] {utterance 1} [SEP] {utterance 2} [SEP] {explanation} [SEP]" for *Pair-Resolve* and "[CLS] {context} [SEP] {response} [SEP] {explanation} [SEP]" for *Diag-Resolve*. The models with *explanation* are denoted as T5_{oracle} and BART_{oracle}, whose performances are shown at rows #3 and #4 in Table 6. We could see that T5_{oracle} and BART_{oracle} surpass T5 and BART by a significant margin, showing that with *explanations* informing what inconsistency the input delivers, the models are able to produce clarification responses more semantically similar to the ground truth. Moreover, BART_{oracle} performs better than T5_{oracle} across all the metrics, demonstrating BART is better at exploiting *explanations* to resolve semantic inconsistency.

Model	#Succ. / #Total	
	Pair-Resolve	Diag-Resolve
BART	56 / 100	36 / 100
BART _{oracle}	91 / 100	82 / 100
ChatGPT	76 / 100	64 / 100
GPT4	92 / 100	79 / 100

Table 7: The number of successfully resolved instances.

Analysis. We go through 200 randomly selected instances (100 from *Pair-Resolve* and 100 from *Diag-Resolve*) of the best-performing BART resolver to 1) check whether the generated responses successfully clarify the inconsistent content and 2) explore the possible reasons that the clarification fails. The numbers of successful instances are presented in Table 7. We could see **both T5 and BART face challenges in inconsistency resolution** and there is still large room for improvement. The higher success count for *Pair-Resolve* compared to *Diag-Resolve* indicates again that resolving inconsistencies between a response and its context poses greater challenges. We summarise the main types of failed clarification as follows:

1. The resolver misses inconsistent content and just picks irrelevant semantic units to form a clarifying response. For instance, the user first says *I want to buy a cup of coffee because I'm so sleepy.* and then *Great, let's try Chinese tea!*. The resolver responds with *Are you on earth sleepy or not?* This error type is common in *Diag-Resolve* because long context contains irrelevant information that interferes with locating inconsistent content.

2. The resolver includes the inconsistent content in the response but fails to form a fluent, contextual coherent response. For example, the user first says *Are you free? I want you to do me a favor.* and then *I am busy now.* and the resolver replies with *Can you do a favor at all?*. In this case, the resolver misunderstands who is the subject of the action, thus providing a response incoherent to the context.

LLMs as consistency resolver. We examine the consistency resolution ability of LLMs by asking LLMs to form a clarification response for the two resolving tasks via the prompts shown in Figure 3 (one in-context example is included in the prompts to ensure a fixed output format).

We report automatic evaluation results in rows #5 and #6 of Table 6. On the selected instances in subsection **Analysis**, we conduct the same human evaluation of the generated clarification response

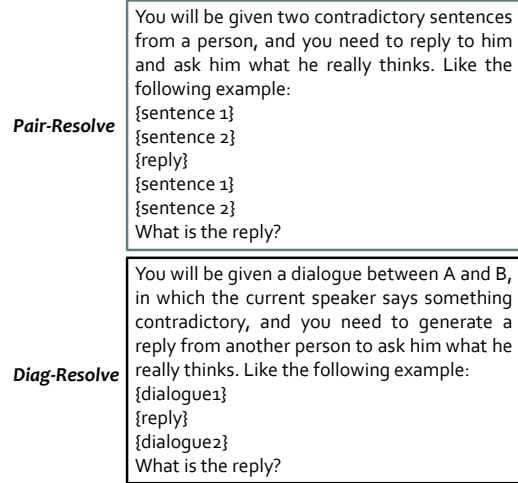


Figure 3: Prompts of resolving tasks.

of the LLMs and show the results in Table 7. Results indicate that: **while ChatGPT and GPT4, both cutting-edge LLMs, score lower in BLEU and ROUGE compared to T5 and BART, they excel in addressing inconsistencies in dialogue history**, whose performance rivals that of the oracle resolvers. The lower BLEU and ROUGE scores of LLMs can be attributed to their tendency to produce more varied and extensive sentences. To illustrate, consider the reference clarification sentence: *Do you really want to eat hot pot or barbecue?.* BART's response is, *Do you really want to eat hot pot or not?*, whereas GPT4 offers, *So, are you more attracted to hot pot, or does barbecue appeal to you more?.*

8 Conclusion

We present **CIDER**, a comprehensive dialogue dataset comprising 27,180 annotated dialogues to investigate conversational inconsistencies. The annotations of **CIDER** cover the whole life span of inconsistencies: the human-authored utterances with inconsistent content demonstrate the introduction of inconsistencies; the explanations help understand the inconsistencies; and the clarification responses exemplify how to resolve the inconsistencies. Through rigorous experiments and analysis, we show that **CIDER** significantly advance the detection and resolution of conversational inconsistencies, and large language models, ChatGPT and GPT4, exhibit commendable performance in identifying and resolving these conversational inconsistencies.

574 Limitation

575 Our work has following limitations:

- 576 • Our proposed dataset emphasizes contradic-
577 tions between utterances. For a truly effective
578 system that detects or resolves inconsistencies,
579 it is essential to incorporate resources that ad-
580 dress other types of inconsistencies, such as
581 intra-utterance or extrinsic discrepancies.
- 582 • We’ve currently evaluated the ability of LLMs
583 to function as independent resolvers under
584 specific prompts to generate clarification ques-
585 tions. The potential for these models to au-
586 tonomously identify and clarify inconsisten-
587 cies remains an intriguing avenue for future
588 exploration. Moreover, while our evaluation
589 of LLMs relies on the optimal prompts cho-
590 sen from several human-crafted options, a
591 more rigorous approach to prompt engineer-
592 ing could potentially yield superior outcomes.

593 Ethical consideration

594 Our dataset, along with the LCCC (Wang et al.,
595 2020) and NaturalConv (Wang et al., 2021) sources,
596 have been cleaned to ensure no breaches of privacy
597 (further details are available in their respective pa-
598 pers). All annotation guidelines (as detailed in
599 Section 4) have received approval from the ethics
600 review committee. We are confident that **CIDER**
601 will play a pivotal role in crafting more human-
602 friendly conversational models.

603 References

604 Nouha Dziri, Ehsan Kamalloo, Kory Mathewson, and
605 Osmar Zaiane. 2019. [Evaluating coherence in dia-](#)
606 [logue systems using entailment](#). In *Proceedings of*
607 *the 2019 Conference of the North American Chap-*
608 *ter of the Association for Computational Linguistics:*
609 *Human Language Technologies, Volume 1 (Long and*
610 *Short Papers)*, pages 3806–3812, Minneapolis, Min-
611 nesota. Association for Computational Linguistics.

612 Meiqi Guo, Mingda Zhang, Siva Reddy, and Malihe
613 Alikhani. 2021. [Abg-CoQA: Clarifying ambiguity in](#)
614 [conversational question answering](#).

615 Ari Holtzman, Jan Buys, Li Du, Maxwell Forbes, and
616 Yejin Choi. 2020. [The curious case of neural text](#)
617 [degeneration](#). In *8th International Conference on*
618 *Learning Representations, ICLR 2020, Addis Ababa,*
619 *Ethiopia, April 26-30, 2020*. OpenReview.net.

Hai Hu, Kyle Richardson, Liang Xu, Lu Li, Sandra
Kübler, and Lawrence Moss. 2020. [OCNLI: Orig-](#)
[inal Chinese Natural Language Inference](#). In *Find-*
ings of the Association for Computational Linguistics:
EMNLP 2020, pages 3512–3526, Online. Association
for Computational Linguistics.

Ziwei Ji, Nayeon Lee, Rita Frieske, Tiezheng Yu, Dan
Su, Yan Xu, Etsuko Ishii, Ye Jin Bang, Andrea
Madotto, and Pascale Fung. 2023. [Survey of halluci-](#)
[nation in natural language generation](#). *ACM Comput-*
ing Surveys, 55(12):1–38.

James Kirkpatrick, Razvan Pascanu, Neil Rabinowitz,
Joel Veness, Guillaume Desjardins, Andrei A Rusu,
Kieran Milan, John Quan, Tiago Ramalho, Ag-
nieszka Grabska-Barwinska, et al. 2017. [Over-](#)
[coming catastrophic forgetting in neural networks](#).
Proceedings of the national academy of sciences,
114(13):3521–3526.

Nayeon Lee, Wei Ping, Peng Xu, Mostofa Patwary, Pas-
cale N Fung, Mohammad Shoeybi, and Bryan Catan-
zaro. 2022. [Factuality enhanced language models](#)
[for open-ended text generation](#). *Advances in Neural*
Information Processing Systems, 35:34586–34599.

Mike Lewis, Yinhan Liu, Naman Goyal, Marjan
Ghazvininejad, Abdelrahman Mohamed, Omer Levy,
Veselin Stoyanov, and Luke Zettlemoyer. 2020.
[BART: Denoising sequence-to-sequence pre-training](#)
[for natural language generation, translation, and com-](#)
[prehension](#). In *Proceedings of the 58th Annual Meet-*
ing of the Association for Computational Linguistics,
pages 7871–7880, Online. Association for Computa-
tional Linguistics.

Margaret Li, Stephen Roller, Ilia Kulikov, Sean Welleck,
Y-Lan Boureau, Kyunghyun Cho, and Jason Weston.
2020. [Don’t say that! making inconsistent dialogue](#)
[unlikely with unlikelihood training](#). In *Proceedings*
of the 58th Annual Meeting of the Association for
Computational Linguistics, pages 4715–4728, Online.
Association for Computational Linguistics.

Chin-Yew Lin. 2004. [ROUGE: A package for auto-](#)
[matic evaluation of summaries](#). In *Text Summariza-*
tion Branches Out, pages 74–81, Barcelona, Spain.
Association for Computational Linguistics.

Yinhan Liu, Myle Ott, Naman Goyal, Jingfei Du, Man-
dar Joshi, Danqi Chen, Omer Levy, Mike Lewis,
Luke Zettlemoyer, and Veselin Stoyanov. 2019.
[Roberta: A robustly optimized bert pretraining ap-](#)
[proach](#). *ArXiv preprint*, abs/1907.11692.

Ilya Loshchilov and Frank Hutter. 2019. [Decoupled](#)
[weight decay regularization](#). In *7th International*
Conference on Learning Representations, ICLR 2019,
New Orleans, LA, USA, May 6-9, 2019. OpenRe-
view.net.

Niels Mündler, Jingxuan He, Slobodan Jenko, and Mar-
tin Vechev. 2023. [Self-contradictory hallucinations](#)
[of large language models: Evaluation, detection and](#)
[mitigation](#). *ArXiv preprint*, abs/2305.15852.

677	Yixin Nie, Mary Williamson, Mohit Bansal, Douwe	Yunfan Shao, Zhichao Geng, Yitao Liu, Junqi Dai,	735
678	Kiela, and Jason Weston. 2021. I like fish, espe-	Fei Yang, Li Zhe, Hujun Bao, and Xipeng Qiu.	736
679	cially dolphins: Addressing contradictions in dia-	2021. Cpt: A pre-trained unbalanced transformer	737
680	logue modeling . In <i>Proceedings of the 59th Annual</i>	for both chinese language understanding and genera-	738
681	<i>Meeting of the Association for Computational Lin-</i>	tion . <i>ArXiv preprint</i> , abs/2109.05729.	739
682	<i>guistics and the 11th International Joint Conference</i>		
683	<i>on Natural Language Processing (Volume 1: Long</i>	Kurt Shuster, Jack Urbanek, Arthur Szlam, and Jason	740
684	<i>Papers)</i> , pages 1699–1713, Online. Association for	Weston. 2022. Am I me or you? state-of-the-art dia-	741
685	Computational Linguistics.	logue models cannot maintain an identity . In <i>Find-</i>	742
		<i>ings of the Association for Computational Linguis-</i>	743
686	OpenAI. 2023. Gpt-4 technical report . <i>ArXiv preprint</i> ,	<i>tics: NAACL 2022</i> , pages 2367–2387, Seattle, United	744
687	abs/2303.08774.	States. Association for Computational Linguistics.	745
688	Long Ouyang, Jeffrey Wu, Xu Jiang, Diogo Almeida,	Yixuan Su and Nigel Collier. 2022. Contrastive search	746
689	Carroll Wainwright, Pamela Mishkin, Chong Zhang,	is what you need for neural text generation . <i>ArXiv</i>	747
690	Sandhini Agarwal, Katarina Slama, Alex Ray, et al.	<i>preprint</i> , abs/2210.14140.	748
691	2022. Training language models to follow instruc-		
692	tions with human feedback. <i>Advances in Neural</i>	Xiaoyang Wang, Chen Li, Jianqiao Zhao, and Dong	749
693	<i>Information Processing Systems</i> , 35:27730–27744.	Yu. 2021. Naturalconv: A chinese dialogue dataset	750
		towards multi-turn topic-driven conversation . In	751
694	Kishore Papineni, Salim Roukos, Todd Ward, and Wei-	<i>Thirty-Fifth AAAI Conference on Artificial Intelli-</i>	752
695	Jing Zhu. 2002. Bleu: a method for automatic evalua-	<i>gence, AAAI 2021, Thirty-Third Conference on In-</i>	753
696	tion of machine translation . In <i>Proceedings of the</i>	<i>novative Applications of Artificial Intelligence, IAAI</i>	754
697	<i>40th Annual Meeting of the Association for Computa-</i>	<i>2021, The Eleventh Symposium on Educational Ad-</i>	755
698	<i>tional Linguistics</i> , pages 311–318, Philadelphia,	<i>vances in Artificial Intelligence, EAAI 2021, Vir-</i>	756
699	Pennsylvania, USA. Association for Computational	<i>tual Event, February 2-9, 2021</i> , pages 14006–14014.	757
700	Linguistics.	AAAI Press.	758
701	Colin Raffel, Noam Shazeer, Adam Roberts, Katherine	Yida Wang, Pei Ke, Yinhe Zheng, Kaili Huang, Yong	759
702	Lee, Sharan Narang, Michael Matena, Yanqi Zhou,	Jiang, Xiaoyan Zhu, and Minlie Huang. 2020. A	760
703	Wei Li, and Peter J. Liu. 2020. Exploring the limits	large-scale chinese short-text conversation dataset .	761
704	of transfer learning with a unified text-to-text	In <i>Natural Language Processing and Chinese Com-</i>	762
705	transformer . <i>J. Mach. Learn. Res.</i> , 21:140:1–140:67.	<i>puting: 9th CCF International Conference, NLPCC</i>	763
		<i>2020, Zhengzhou, China, October 14–18, 2020, Pro-</i>	764
706	Hannah Rashkin, David Reitter, Gaurav Singh Tomar,	<i>ceedings, Part I 9</i> , pages 91–103. Springer.	765
707	and Dipanjan Das. 2021. Increasing faithfulness		
708	in knowledge-grounded dialogue with controllable	Sean Welleck, Iliia Kulikov, Stephen Roller, Emily Di-	766
709	features . In <i>Proceedings of the 59th Annual Meet-</i>	nan, Kyunghyun Cho, and Jason Weston. 2020. Neu-	767
710	<i>ing of the Association for Computational Linguistics</i>	ral text generation with unlikelihood training . In	768
711	<i>and the 11th International Joint Conference on Natu-</i>	<i>8th International Conference on Learning Represen-</i>	769
712	<i>ral Language Processing (Volume 1: Long Papers)</i> ,	<i>tations, ICLR 2020, Addis Ababa, Ethiopia, April</i>	770
713	pages 704–718, Online. Association for Computa-	<i>26-30, 2020</i> . OpenReview.net.	771
714	tional Linguistics.		
715	Stephen Roller, Emily Dinan, Naman Goyal, Da Ju,	Sean Welleck, Jason Weston, Arthur Szlam, and	772
716	Mary Williamson, Yinhan Liu, Jing Xu, Myle Ott,	Kyunghyun Cho. 2019. Dialogue natural language	773
717	Eric Michael Smith, Y-Lan Boureau, and Jason We-	inference . In <i>Proceedings of the 57th Annual Meet-</i>	774
718	ston. 2021. Recipes for building an open-domain	<i>ing of the Association for Computational Linguistics</i> ,	775
719	chatbot . In <i>Proceedings of the 16th Conference of</i>	pages 3731–3741, Florence, Italy. Association for	776
720	<i>the European Chapter of the Association for Computa-</i>	Computational Linguistics.	777
721	<i>tional Linguistics: Main Volume</i> , pages 300–325,		
722	Online. Association for Computational Linguistics.	Haoran Wu, Wenxuan Wang, Yuxuan Wan, Wenxiang	778
		Jiao, and Michael Lyu. 2023. Chatgpt or grammarly?	779
723	Swarnadeep Saha, Yixin Nie, and Mohit Bansal. 2020.	evaluating chatgpt on grammatical error correction	780
724	ConjNLI: Natural language inference over conjunc-	benchmark . <i>ArXiv preprint</i> , abs/2303.13648.	781
725	tive sentences . In <i>Proceedings of the 2020 Confer-</i>		
726	<i>ence on Empirical Methods in Natural Language</i>	Hamed Zamani, Susan Dumais, Nick Craswell, Paul	782
727	<i>Processing (EMNLP)</i> , pages 8240–8252, Online. As-	Bennett, and Gord Lueck. 2020a. Generating clar-	783
728	sociation for Computational Linguistics.	ifying questions for information retrieval. In <i>Pro-</i>	784
		<i>ceedings of The Web Conference 2020, WWW ’20</i> ,	785
729	Sashank Santhanam, Behnam Hedayatnia, Spandana	pages 418–428, New York, NY, USA. Association	786
730	Gella, Aishwarya Padmakumar, Seokhwan Kim,	for Computing Machinery.	787
731	Yang Liu, and Dilek Hakkani-Tur. 2021. Rome was		
732	built in 1776: A case study on factual correctness	Hamed Zamani, Susan Dumais, Nick Craswell, Paul	788
733	in knowledge-grounded response generation . <i>ArXiv</i>	Bennett, and Gord Lueck. 2020b. Generating clar-	789
734	<i>preprint</i> , abs/2110.05456.	ifying questions for information retrieval. In <i>Pro-</i>	790
		<i>ceedings of The Web Conference 2020, WWW ’20</i> ,	791

- 792 pages 418–428, New York, NY, USA. Association
793 for Computing Machinery.
- 794 Mian Zhang, Lifeng Jin, Linfeng Song, Haitao Mi, Wen-
795 liang Chen, and Dong Yu. 2023. Safeconv: Explain-
796 ing and correcting conversational unsafe behavior. In
797 *Proceedings of the 61st Annual Meeting of the As-
798 sociation for Computational Linguistics (Volume 1:
799 Long Papers)*, pages 22–35.
- 800 Saizheng Zhang, Emily Dinan, Jack Urbanek, Arthur
801 Szlam, Douwe Kiela, and Jason Weston. 2018. [Per-
802 sonalizing dialogue agents: I have a dog, do you
803 have pets too?](#) In *Proceedings of the 56th Annual
804 Meeting of the Association for Computational Lin-
805 guistics (Volume 1: Long Papers)*, pages 2204–2213,
806 Melbourne, Australia. Association for Computational
807 Linguistics.
- 808 Zhe Zhao, Hui Chen, Jinbin Zhang, Xin Zhao, Tao Liu,
809 Wei Lu, Xi Chen, Haotang Deng, Qi Ju, and Xiaoy-
810 ong Du. 2019. [UER: An open-source toolkit for pre-
811 training models.](#) In *Proceedings of the 2019 Confer-
812 ence on Empirical Methods in Natural Language Pro-
813 cessing and the 9th International Joint Conference
814 on Natural Language Processing (EMNLP-IJCNLP):
815 System Demonstrations*, pages 241–246, Hong Kong,
816 China. Association for Computational Linguistics.
- 817 Chujie Zheng, Jinfeng Zhou, Yinhe Zheng, Libiao Peng,
818 Zhen Guo, Wenquan Wu, Zheng-Yu Niu, Hua Wu,
819 and Minlie Huang. 2022. [CDConv: A benchmark
820 for contradiction detection in Chinese conversations.](#)
821 In *Proceedings of the 2022 Conference on Empirical
822 Methods in Natural Language Processing*, pages 18–
823 29, Abu Dhabi, United Arab Emirates. Association
824 for Computational Linguistics.