# NovelCR: A Large-Scale Bilingual Dataset Tailored for Long-Span Coreference Resolution

**Meihan Tong**[1], **Shuai Wang**[2]
[1]Tsinghua University, Beijing, China
[2]JOYY Inc, China
tongmeihan@gmail.com, shuaiwa16@gmail.com

## Abstract

Coreference resolution (CR) links pronouns and noun phrases to their referent entities, serving as a key step in deep text understanding. Presently available CR datasets are either small in scale or restrict coreference resolution to a limited text span. In this paper, we present NovelCR, a large-scale bilingual benchmark designed for long-span coreference resolution. NovelCR features extensive annotations, including 148k mentions in NovelCR-en and 311k mentions in NovelCR-zh. Moreover, the dataset is notably rich in long-span coreference pairs, with 85% of pairs in NovelCR-en and 83% in NovelCR-zh spanning across three or more sentences. Experiments on NovelCR reveal a large gap between state-of-the-art baselines and human performance, highlighting that NovelCR remains an open issue. The NovelCR dataset is publicly available at: https://github.com/tongmeihan1995/NovelCR.

## 1 Introduction

Coreference resolution (CR) aims to identify mentions and their referent entities from text. For instance, given the sentence *"Recently, Apple sued Qualcomm, suing it for failing to cooperate by contracts"*, coreference resolution needs to distinguish that mention *it* here refers to entity *Qualcomm* instead of *Apple*. Coreference resolution is a core task in deep text analysis and acts as a prerequisite for multiple advanced natural language processing applications such as machine reading comprehension (Wu et al., 2020), information extraction (Zelenko et al., 2004), and multi-round dialogue construction (Yu et al., 2022).

However, existing coreference resolution datasets either suffer from small data scales or restrict coreference resolution within a limited text span. ACE2004 (Doddington et al., 2004) annotates coreferences from merely 451 documents. The data scales of WikiCoref (Ghaddar and Langlais, 2016), MUC-6 (muc, 1995), MUC-7

(Hirschman, 1997), STM-coref (Brack et al., 2021) are even smaller, containing 30, 60, 50, and 110 documents, respectively. Given their small scale, none of these coreference datasets can fairly assess the performance of modern neural coreference resolution models. The Winograd Schema Challenge (WSC) corpus (Levesque et al., 2012) restricts referential ambiguity to intra-sentential settings, so all anaphoric links are resolved within a single sentence. Other public benchmarks likewise exhibit only modest cross-sentence span lengths. The average sentence distance between antecedent and anaphor in CoNLL-2012 (Weischedel et al., 2011), ECB+ (Cybulska and Vossen, 2014), GAP (Webster et al., 2018a), LongtoNotes (Shridhar et al., 2022), and DWIE (Zaporojets et al., 2020) is just 2.9, 3.1, 2.3, 3.3, and 2.8 sentences, respectively. The predominance of short-span coreference pairs in these datasets substantially lowers the overall difficulty of the task these datasets pose.

In this paper, we introduce NovelCR, a large-scale bilingual benchmark designed to tackle long-span coreference resolution. Long-span coreference resolution requires models to reason over broad discourse context rather than local sentence cues. As Figure 1 shows, linking *the lady on the ground* back to *Jerebal* demands that the models track dialogue-role shifts across several sentences, ignore intervening mentions, and overcome an apparent gender mismatch — highlighting the depth of contextual understanding required for robust coreference resolution. To ensure that NovelCR captures such complex phenomena, we focus on annotating the coreference chains of novel characters, whose mentions are naturally dispersed across extended narrative spans. For example, characters like Jerebal, Quila, and Quil in Figure 1(a) are frequently referenced after multiple intervening sentences, making them ideal anchors for evaluating discourse-level resolution capabilities.

| Jerebai | Quila | Quil |

**Chapter 1**

Hearing the voice of the visitor, the lady on the ground finally moved. Her cracked lips quivered, asking, "Quila, how's Quil? Perhaps it was because she hadn't spoken for such a long time, but her voice sounded extremely hoarse, like the grinding of gravel on the floor.

Qulla frowned, with ever-growing abhorrence in her eyes. "Haaa--? My brother?" She hooked her lips into a smile full of ridicule and derision, "Jerebai are you still expecting him to come and save you? Do you know what day it is today? Today is the day that he marries my new sister-in-law! He is in love - do you really expect that you, a murderous demoness would even cross his mind?!" The man's sister cried.

He actually...

Jerebai heart felt as though it had been stabbed by a needle - and it wasn't an acute unbearable type of pain, but the type of pain that reverberates and lingers, even eking out traces of blood ever so slowly.

She should have known. After all, that person had not come to save her after such a long time…

Jerebai unconsciously held her abdomen. She once carried a child belonging to her and that man.

**Figure 1(a)**

于修逸很想知道什么意思，可秦亦封却什么都不说了，气得他直跳脚。

算了算了，要是不想说，谁也拿他没办法，于修逸长叹了口气，觉得这一刻的秦亦封很可怕。

Entity: 秦亦封
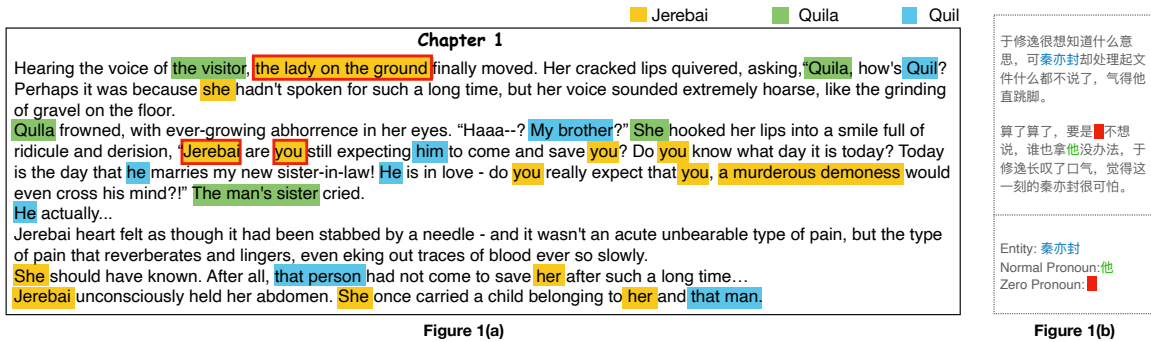Normal Pronoun:他
Zero Pronoun:

**Figure 1(b)**

Figure 1: Figure 1(a): An example of NovelCR. NovelCR resolves coreferences for novel characters Jerebal, Quila, and Quil. *Jerebal* and *the lady on the ground* form a long-span coreference pair, with the entity and mention separated by 6 sentences. *Jerebal* and *you* form a short-span coreference pair, with the entity and mention occurring in 1 sentence. Figure 1(b): An example of zero pronoun resolution in NovelCR-zh.

The construction process of NovelCR is as follows: we first obtain English and Chinese novels from online websites. Then, we leverage NER tools and prompt learning to collect candidate entities and mentions from novel chapters. We cover a wide range of mentions in our dataset, including pronouns (e.g., *she*, *her*, and *him*), proper and common noun phrases (e.g., *the visitor*, *the man's sister*, and *a murderous demones*) to reduce the likelihood of missing labels. Afterward, we utilize crowdsourcing to remove improper mentions and re-edit mention boundaries to ensure that all mentions adhere to the maximum span principle. Finally, annotators are required to answer multiple-choice questions to match mentions to their referent entities.

We highlight the three contributions of NovelCR: (1) Large scale. NovelCR contains a total of 460k mentions and 402k coreferences, making it significantly larger than existing CR datasets. (2) Abundant long-span coreference pairs. The ratio of coreference pairs scattered over 3 or more sentences reaches 85% in NovelCR-en and 83% in NovelCR-zh, giving NovelCR the highest percentage of long-span coreference pairs among current datasets. (3) Bilingual. NovelCR provides coreference annotations in both English and Chinese, supporting the exploration of cross-lingual learning within the dataset. Additionally, NovelCR-zh includes a substantial number of zero pronouns (as shown in Figure 1(b)), making the task of zero pronoun resolution an integral part of the dataset and further increasing its complexity and challenge.

We evaluate NovelCR against 12 state-of-the-art CR baselines. Experiments show that there is still a large gap between the CR baselines and human beings, revealing that NovelCR remains an unresolved challenge. Detailed experiments show that the cutting-edge CR model still performs poorly when dealing with long-span coreference resolution.

## 2 Related Work

MUC-6 (muc, 1995) and MUC-7 (Hirschman, 1997) are the two earlier proposed coreference resolution datasets, and the data scale is relatively small, with 60 and 50 documents, 30k and 25k tokens respectively. WikiCoref (Ghaddar and Langlais, 2016) contains merely 30 documents and 7955 mentions. STM-coref (Brack et al., 2021) annotates coreferences from no more than 110 research papers. GUM (Zeldes, 2017) and ARRAU (Uryupina et al., 2016) solve anaphora resolution from open source multi-layer corpus with barely 300 documents. ACE2004 (Doddington et al., 2004) is a widely adopted CR dataset that covers multiple domains, including news communications, broadcast programs, and online blogs. However, it contains a relatively small amount of data, with just 451 documents and 22,550 mentions. In contrast, the proposed dataset NovelCR features an extensive dataset, with 28k documents and 460k mentions, far exceeding existing CR datasets. Additionally, the proposed dataset NovelCR focuses on both English and Chinese coreference resolution, unlike PreCo (Chen et al., 2018), LongtoNotes (Shridhar et al., 2022) and LitBank (Bamman et al., 2020), which is a single-language dataset.

Winograd Schema Challenge (WSC) (Levesque et al., 2012) is a well-known CR benchmark proposed by Hector Levesque, consisting of 803 coref-

| Datasets | #Doc. | #Sent. | #Mnt. | #Coref. | #Dis/s | #Dis/tok | #LongCoref/s | #LongCoref/tok |
|---|---|---|---|---|---|---|---|---|
| ACE2004 | 451 | 18,530 | 22,550 | - | - | - | - | - |
| MUC-6 | 60 | 3,750 | - | - | - | - | - | - |
| MUC-7 | 50 | 3,197 | - | - | - | - | - | - |
| WikiCoref | 30 | 2,292 | 7,955 | 6,700 | 3.5 | 58.4 | 3,082 (46%) | 1,069 (16%) |
| WSC | - | 803 | 2,409 | 1,606 | 1.0 | - | 0 (0%) | 0 (0%) |
| GAP | - | 8,908 | 26,724 | 17,816 | 2.3 | 43.6 | 0 (0%) | 1,394 (8%) |
| STM-coref | 110 | 1,480 | 2,577 | 1,669 | 2.4 | 41.6 | 484 (29%) | 116 (7%) |
| CoNLL2012 | 3,493 | 112,941 | 56,371 | 43,560 | 2.9 | 48.9 | 14,810 (34%) | 89,23 (20%) |
| LongtoNotes | 2,415 | 112,941 | 38,640 | 32,715 | 3.3 | 54.2 | 12,104 (37%) | 7,824 (24%) |
| LitBank | 100 | 108,000 | 57,514 | 28,411 | 6.8 | 109.4 | 19,603 (69%) | 15,613 (55%) |
| ECB+ | 502 | 9,171 | 32,297 | 12,930 | 3.1 | 52.7 | 5,301 (41%) | 1,843 (14%) |
| DWIE | 802 | 13,628 | 43,373 | 20,243 | 2.8 | 49.9 | 7,085 (35%) | 2,731 (13%) |
| NovelCR-en(ours) | 9,462 | 289,285 | 148,529 | 128,847 | 8.2 | 114.5 | 109,520 (85%) | 81,137 (63%) |

Table 1: Statistics of English coreference resolution datasets. Doc.: novel chapters, Mnt.: mentions, Coref.: coreference pairs, Dis/s: average distance between coreference pairs, measured in sentence. Dis/s: average distance between coreference pairs, measured in LLM tokens. LongCoref/s: number of long-span coreference pairs, where the gap between pairs is no less than three sentences, LongCoref/tok: number of long-span coreference pairs, where the gap between pairs is no less than 100 LLM tokens (Llama3.1).

| Datasets | #Doc. | #Sent. | #Mnt. | #Coref. | #Dis/s | #Dis/tok | #LongCoref/s | #LongCoref/tok |
|---|---|---|---|---|---|---|---|---|
| ACE2004 | 646 | 14,233 | 28,135 | - | - | - | - | - |
| CoNLL2012 | 2,280 | 83,763 | 15,136 | 8,859 | 3.1 | 79.0 | 3,986 (45%) | 438 (5%) |
| CLUEWSC2020 | - | 1,648 | 4,944 | 3,296 | 1.0 | - | 0 (0%) | 0 (0%) |
| NovelCR-zh(ours) | 19,288 | 80,872 | 311,482 | 273,379 | 5.2 | 184.3 | 258,530 (83%) | 216,538 (70%) |

Table 2: Statistics of Chinese coreference resolution datasets.

erences. WSCR (Rahman and Ng, 2012), PDP (Davis et al., 2017), WINOBIAS (Zhao et al., 2018), and WinoGrande (Sakaguchi et al., 2021) are datasets evolved from the WSC. All the above datasets limit coreference resolution in a single sentence. Besides, most of the coreference pairs in CoNLL2012-en, CoNLL2012-zh (Weischedel et al., 2011), GAP (Webster et al., 2018a), ECB+ (Cybulska and Vossen, 2014), and DWIE (Zaporojets et al., 2020) appear within the scope of three sentences. The prevalence of short-span coreferences in these datasets makes them less challenging. Unlike them, our proposed dataset NovelCR contains a significant number of long-span coreferences, and this abundance of long-span coreferences necessitates a more robust semantic understanding model to handle NovelCR effectively.

## 3 Dataset Construction

In this section, we illustrate the dataset construction process. As shown in Figure 2, We construct NovelCR in three steps: novel chapter collection, mention detection, and coreference identification. Novel chapter collection aims to gather chapters from a wide range of genres sourced from online

novel websites. Mention detection leverages NER tools (Stanford CoreNLP tool for English, LTP tool for Chinese) and prompt learning to mine potential entities and mentions from novel chapters. Coreference identification uses crowdsourcing to distinguish coreference pairs in chapters by converting coreference resolution into multiple-choice questions.

### 3.1 Novel Chapter Collection

We select online novels as our data source. The underlined reason is that novels, unlike news articles, exhibit strong narrative coherence and are more likely to include long-span coreferences. Specifically, we crawl hundreds of popular English and Chinese novels from the online reading site WUX-IAWORLD [1], all of which are open source and free to access. The crawled novels encompass a wide range of genres such as cultivation, fantasy, comedy, suspense, romance, science fiction, etc. In total, we collected 1000 English novels (97,723 chapters) for NovelCR-en and 2000 Chinese novels (187,492 chapters) for NovelCR-zh. These novels were originally written in Chinese and translated
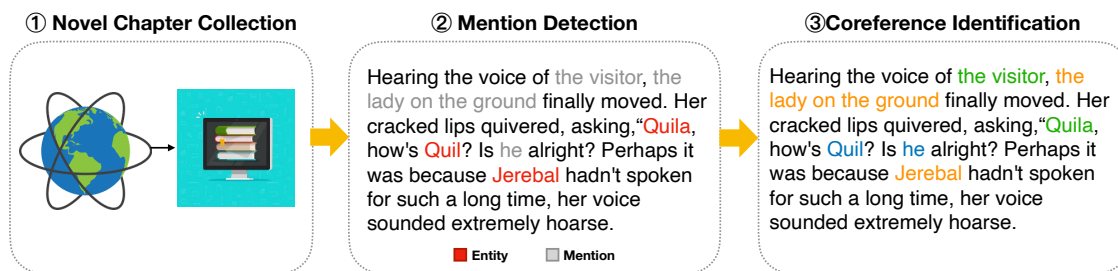
---

[1] https://www.wuxiaworld.com/

① **Novel Chapter Collection**     ② **Mention Detection**     ③**Coreference Identification**

Hearing the voice of the visitor, the lady on the ground finally moved. Her cracked lips quivered, asking,"Quila, how's Quil? Is he alright? Perhaps it was because Jerebal hadn't spoken for such a long time, her voice sounded extremely hoarse.

■ Entity    ▢ Mention

Hearing the voice of the visitor, the lady on the ground finally moved. Her cracked lips quivered, asking,"Quila, how's Quil? Is he alright? Perhaps it was because Jerebal hadn't spoken for such a long time, her voice sounded extremely hoarse.

Figure 2: Labeling Process of NovelCR.

into English by human experts. Due to the incomplete translation, the number of English novels is less than that of Chinese novels.

We filter out novel chapters with less than 256 tokens and more than 32,768 tokens to balance the document lengths. Additionally, we utilize NER tools (Stanford NLP for English and LTP for Chinese) to filter out chapters with less than 8 entities, ensuring abundant coreference annotations. After two rounds of filtering, we collect 9,462 novel chapters for NovelCR-en and 19,288 novel chapters for NovelCR-zh.

### 3.2 Mention Detection

This section aims to detect candidate mentions and entities from novel chapters. To reduce the burden on annotators, mention detection is divided into two steps. The first step is to use NER tools and prompt learning to mine candidate entities and mentions, and the second step is to employ annotators to do manual verification.

#### 3.2.1 Candidates Collection

To detect candidate entities, we employ Stanford CoreNLP [2] and LTP NER tool[3] to recognize named entities from English and Chinese chapters, respectively. Finally, we detected 42,849 and 98,571 person entities for NovelCR-en and NovelCR-zh respectively. We involved three students to conduct human evaluations to assess the quality of annotations. The average recall rates of NER on NovelCR-en and NovelCR-zh are 99.1% and 98.9%, respectively, demonstrating the effectiveness of the named entity tools.

Previous datasets usually employ POS tagging to detect pronoun mentions and semantic parsers to detect noun phrase mentions, yet these mention detection methods are pattern-dependent and the

mention recall rate is not that high. In this paper, we employ prompt-learning (Ouyang et al., 2022) to detect pronoun and noun phrase mentions. Empowered by ChatGPT, prompt learning has strong text comprehension capabilities and can identify a wider variety of mentions. Finely designed prompts are shown in Appendix C. We take the union of annotations of different prompts as the final annotation result. We engage three students to conduct human evaluations. As shown in Appendix B Table 9, compared to the traditional method (POS tagging + Semantic Parser), our proposed method (Prompt-Learning) improves the recall rate by 7.8% and 8.2% on NovelCR-en and NovelCR-zh, respectively, effectively reducing the risk of missing annotations.

We additionally train a sequence labeling model to handle Chinese zero pronoun resolution. We leverage OntoNotes (Weischedel et al., 2011) as our training corpus and adopt BERT as the backbone. The training goal is to insert a special token before the zero pronoun. For instance, given the sentence "She poured water until *it* was full", where *it* is omitted in Chinese, the output of the sequence labeling model is "She poured water until *[Zero Pronoun]* was full". The average recall rate in human evaluation is 87.4% on Chinese zero pronoun resolution.

#### 3.2.2 Manual Verification

In the section, we manually verify the entities and mentions obtained in Section 3.2. We invite a total of 136 Chinese college students to participate in our crowdsourcing annotation. The annotators of NovelCR-en are English-major students with TOEFL scores higher than 100 or IELTS scores higher than 7.5, and the annotators of NovelCR-zh are native Chinese speakers.

As shown in the guideline in Appendix A, annotators first remove invalid mentions, i.e., mentions

---

[2]https://github.com/stanfordnlp/CoreNLP
[3]https://www.ltp-cloud.com/intro_en

that do not refer to a person entity, such as *the bank* and *this beautiful knife*. In particular, *her* in *her split lips* is also considered an invalid mention as it functions as a modifier of *lips* rather than an independent personal pronoun. Only mentions verified by at least two annotators will be retained. By removing invalid mentions, we ensure the quality of mentions in the proposed dataset, but it may also cause missing annotations, which we will discuss in the limitations section.

After that, the annotators are required to refine the boundary of the mention. We adopt the principle of maximum span. For example, given the mention *a little child*, if the original annotation is *child*, the annotator needs to adjust the boundary to *a little child*. If two of the three annotators edit the boundary in the same way, we will accept the revision. Otherwise (0.8% of the time), we ask another experienced annotator to make the final decision. This experienced annotator should have annotated more than 50 chapters with an accuracy rate of more than 95%.

### 3.3 Coreferences Identification

In this section, we leverage crowdsourcing to identify coreferences from the novel chapters.

We reframe coreference identification as a multiple-choice question. Specifically, we first collect the entity set $E$ from the chapter and deduplicate it. Then, for each mention $m$ in the chapter, we ask the annotators to determine which entity in $E$ the mention $m$ refers to. Taking Figure 2 as an example, the entity set in the novel chapter is {*Quila*, *Quil*, *Jerebal*}. Given the mention *the visitor*, annotators need to determine which entity *the visitor* refers to, *Quila*, *Quil* or *Jerebal*. We adopt the answer *Quila* as the final coreference annotation.

Each mention undergoes labeling by three annotators, with the final result determined by the majority vote. If the three annotators cannot reach a consensus (4.5% of the cases), we engage another experienced annotator to make the final decision. The experienced annotator should have annotated more than 50 chapters with an accuracy rate of more than 95%. The guideline is shown in Appendix B. We remove the singleton mentions after finishing the annotation.

### 3.4 Annotation Quality & Remuneration

We use Cohen's kappa coefficient (Artstein and Poesio, 2008; McHugh, 2012) to measure the inter-

annotator agreement (IAA) of crowdsourced labeling. The IAA scores are respectively 96% and 92% for mention verification (Section 3.2.2) and coreference identification (Section 3.3), indicating very high labeling agreement.

We pay $0.1 per data point per annotator for mention verification and $0.3 per data point per annotator for coreference identification. According to our standards, the hourly wage of annotators is not less than 10 US dollars per hour, which exceeds the US minimum hourly wage of 7.25 US dollars per hour. We release NovelCR under the Creative Commons Attribution-NonCommercial License (CC BY-NC).

## 4 Data Analysis

### 4.1 Overall Statistic

We compare NovelCR-en and NovelCR-zh to existing representative English and Chinese coreference resolution datasets in Table 1 and Table 2 respectively.

From the tables, we can draw the following observations. First, our dataset is much larger than existing CR datasets. As shown in Table 1, NovelCR-en contains 9,462 documents, 289,285 sentences, 8.1M tokens, 148,529 mentions, and 128,837 coreference pairs. Even compared with the current large CR datasets CoNLL2012 and LongtoNotes, our dataset is still 2.9 and 2.6 times larger in terms of documents and 3.0 and 4.0 times larger in terms of the number of coreference pairs. This phenomenon is more pronounced in comparisons involving Chinese datasets. As shown in Table 2, NovelCR-zh contains 19,288 documents, 80,872 sentences, 21M tokens, 311,482 mentions, and 273,379 coreference pairs. The number of coreference pairs is 30.9 times that of CoNLL2012 and 82.9 times that of CLUEWSC2020.

In addition, our dataset contains abundant long-span coreference pairs. As shown in Table 1, the average distance between coreference pairs in NovelCR-en is 8.2 sentences, longer than that in LongtoNotes (3.3 sentences), ECB+ (3.1 sentences), and LitBank (6.8 sentences). NovelCR-en also has the largest proportion of coreference pairs spread over 3 or more sentences, reaching 85%, which is greater than LongtoNotes (37%), ECB+ (41%), and LitBank (69%). Given the varying lengths of sentences, we additionally quantify the number and proportion of long-coreference pairs in terms of LLM tokens. As shown in Table 1,
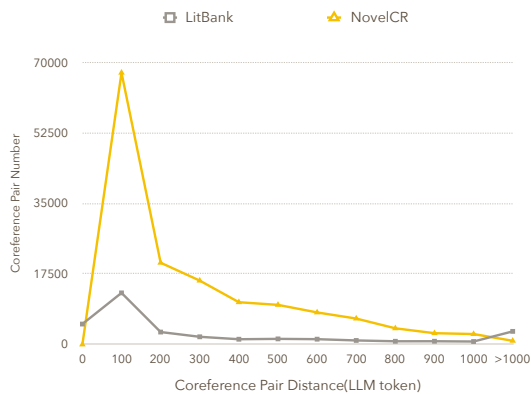
Figure 3: Dataset Comparision between LitBank and NovelCR(ours)

NovelCR-en exhibits the largest number of long-span coreference pairs with a gap of no less than 100 tokens (81,137), surpassing the STM-coref (116), CoNLL2012 (89,23), LongtoNotes(7,824), and LitBank (15,613). This highlights the challenge posed by NovelCR in long-span coreference resolution. This challenge becomes even more pronounced when comparing Chinese datasets. As shown in Table 2, the ratio of long-span coreference pairs in NovelCR-zh has reached 83%, far exceeding the proportion of long-span coreference pairs in existing Chinese CR datasets.

**Annotation Policy Comparison** Table 3 illustrates the different annotation policies adopted by existing CR datasets and our dataset.

As shown in Table 3, our annotation follows three main principles: (1) singletons are excluded from coreference chains, in line with OntoNotes; (2) indefinite and generic references are omitted, as they do not refer to specific, identifiable entities; and (3) copular constructions are annotated only when they involve a named entity. For example, in "*Barack Obama* is the president," since Barack Obama is a named entity, we include the coreference between Barack Obama and the president in our dataset.

## 4.2 Detailed Statistic

As shown in Table 1, LitBank contains the highest proportion of long-span coreference pairs among existing datasets. Moreover, both LitBank and our dataset (NovelCR) use novels as their data source. Therefore, in this section, we first present a detailed comparative analysis between LitBank and NovelCR, focusing specifically on the distribution of coreference pair lengths measured in LLM tokens.

As shown in Figure 3, the number of coreference

pairs with varying gaps in NovelCR (ours) far exceeds that in LitBank, highlighting the large-scale nature of NovelCR (ours). Besides, in terms of the number of long-span coreference pairs with a gap of at least 100 tokens, NovelCR (ours) also surpasses LitBank, aligning with its original intention of constructing a dataset for long-span coreference resolution. We observe that LitBank contains slightly more super long-span coreference pairs (with gaps greater than 1000 tokens), since LitBank annotates coreference across entire novels, while NovelCR (ours) annotates within individual chapters.

Then, we analyze the distribution of mention lengths in NovelCR-en. According to statistics, 54% mentions contain 1 word, most of which are entities and personal pronouns, such as *she* and *her*. 36% mentions consist of 2-5 tokens, and 10% mentions exceed 5 tokens, most of which were noun phrases of named entities, such as *that person*, and *the beloved woman in front of me*.

After that, we analyze the distribution of document lengths in NovelCR-en. Statistics reveal that 61% of documents consist of less than 10k tokens. 33% of documents are comprised of 10k-20k tokens, while 6% of documents extend beyond 20k tokens.

We also analyze gender bias within our dataset. Following (Karimi et al., 2016; Webster et al., 2018b), we use the Gender Guesser library4 [4] to determine the gender of the mentions. According to the statistics, 45.1% of mentions belong to *male* or *mostly male* names, 34.2% of mentions belong to *female* or *mostly female* names, and 20.7% were classified as *unknown*. The ratio between female and male candidates is estimated to be 58%, with male candidates predominating.

In addition, we count the number of zero pronouns in NovelCR-zh. In total, we annotate 84,738 zero pronouns, accounting for 27.2% of the annotated mentions in NovelCR-zh.

## 5 Experiment

### 5.1 Experimental Setup

**Data Split** Table 4 shows the detailed statistics of dataset splitting. We split NovelCR-en and NovelCR-zh in a ratio of 8:1:1 to form training, validation, and test sets.

**Hyperparameters** Our experiments are conducted on eight A100 GPUs, each with 80GB of

---

[4]https://pypi.org/project/gender-guesser/

5166

| Dataset | Singletons | Indefinite/Generic References | Copular Constructions |
|---|---|---|---|
| ACE2004 | Exclude | Exclude | Exclude |
| WikiCoref | Include | Include | Include |
| CoNLL2012 | Exclude | Exclude | Exclude |
| LitBank | Include | Include | Include |
| NovelCR (Ours) | Exclude | Exclude | Include |

Table 3: Comparison of annotation policies across coreference datasets.

| Method | NovelCR-en | | | NovelCR-zh | | |
|---|---|---|---|---|---|---|
| | Train | Dev | Test | Train | Dev | Test |
| #Doc. | 7K | 1.5K | 1.5K | 15K | 2K | 2K |
| #Men. | 118K | 15K | 15K | 247K | 31K | 32K |
| #Coref. | 107K | 10K | 11K | 212K | 33K | 28K |

Table 4: Data Split in NovelCR. #Doc., #Men. and #Coref. refer to the number of novel chapters, mentions, and coreference pairs.

GPU memory. We ran the experiment using the default hyperparameters in the baseline release codes. The training time for the baselines is about half an hour. Long chapters are split into non-overlapping segments of up to 2048 word-piece tokens. For human evaluation, we invited three students to annotate 200 documents randomly selected from NovelCR-en and NovelCR-zh and report the average accuracy of the three students as the final results.

**Metrics** We utilize precision, recall, and F1 to evaluate the performance of existing baselines on the proposed dataset. All the metrics are calculated in B3 (Bagga and Baldwin, 1998), MUC (Vilain et al., 1995), CEAFe (Luo, 2005), and CoNLL (Pradhan et al., 2011, 2012) (average of B3, MUC and CEAFe), respectively, to allow adequate comparison. We report the average result of five rounds.

## 5.2 Baseline

We introduce twelve baselines to validate the challenges of the NovelCR, including: **e2e-coref** (Lee et al., 2017) is an end-to-end coreference resolution model, which considers all spans as potential mentions and learns the probabilities of possible antecedents for each mention. **c2f-coref** (Lee et al., 2018) introduces a coarse-to-fine approach to accelerate coreference resolution, which allows for more aggressive span pruning without compromising accuracy. **CR-BERT** (Joshi et al., 2019b) applies BERT to coreference resolution, achieving significant improvements on the CoNLL2012 and GAP benchmarks. **SpanBERT** (Joshi et al., 2019a) upgrades BERT from word-level pre-training to span-

level pre-training via geometric masking to better cope with span-level coreference resolution. **WL-COREF** (Dobrovolskii, 2021) finds coreferences at the granularity of tokens rather than word spans, and then reconstructs the word spans to reduce the complexity of the coreference model. **Link-Append** (Bohnet et al., 2022) uses the seq2seq paradigm and transition matrix to jointly predict mentions and entities, which formulate coreference resolution as a generation task. **BookNLP** (Bamman, 2020) is an open-source NLP pipeline for analyzing long-form English texts (novels and historical documents), offering tools to extract rich linguistic and narrative information. **LongDoc** (Wu et al., 2021) proposes a neural coreference model that efficiently handles long documents by selectively ignoring irrelevant context using bounded memory mechanisms. **Maverick** (Martinelli et al., 2024) is an efficient and accurate coreference resolution system that outperforms much larger models while using significantly fewer parameters and achieving faster inference. **Fastcoref** (Otmazgin et al., 2022a) is a precise and user-friendly coreference resolution algorithm that is widely used. We employ LingMess implementation (Otmazgin et al., 2022b) in our experiments. **Llama3.1(CoT)** (Touvron et al., 2023) prompts the Llama3.1 to think step by step (Wei et al., 2022). **GPT-o1** (OpenAI, 2024) allows LLMs to think longer before they answer a prompt.

Among the baselines, e2e-coref, c2f-coref, CR-BERT, SpanBERT, WL-COREF, Link-Append, LongDoc, and Fastcoref are fine-tuned on NovelCR, whereas BookNLP, Maverick, Llama3.1(CoT), and GPT-o1 are used for inference only.

## 5.3 Overall Performance

Table 5 and Table 6 show the experimental results of NovelCR-en and NovelCR-zh, from which we have the following observations.

(1) Human beings have achieved good performance on NovelCR, achieving an F1 score of 91.4% on NovelCR-en and 90.5% on NovelCR-zh

| Methods | B3 | | | MUC | | | CEAFe | | | CoNLL |
|---|---|---|---|---|---|---|---|---|---|---|
| | P | R | F | P | R | F | P | R | F | F |
| e2e-coref | 59.4 | 57.1 | 58.2 | 62.3 | 59.4 | 60.8 | 58.8 | 61.2 | 60.0 | 59.7 |
| c2f-coref | 64.7 | 66.5 | 65.6 | 67.2 | 65.9 | 66.5 | 65.3 | 68.7 | 67.0 | 66.4 |
| CR-BERT | 74.3 | 71.8 | 73.0 | 75.5 | 73.9 | 74.7 | 74.7 | 72.5 | 73.6 | 73.8 |
| SpanBERT | 68.4 | 72.2 | 70.2 | 72.6 | 73.4 | 73.0 | 73.4 | 71.2 | 72.3 | 71.8 |
| WL-COREF | 73.1 | 71.6 | 72.3 | 73.3 | 72.8 | 73.0 | 70.6 | 74.9 | 72.7 | 72.7 |
| Link-Append | 63.4 | 62.7 | 63.0 | 65.5 | 68.1 | 66.8 | 67.8 | 64.2 | 66.0 | 65.3 |
| BookNLP | 64.7 | 68.3 | 66.5 | 71.7 | 69.8 | 70.7 | 64.5 | 65.1 | 64.8 | 67.3 |
| LongDoc | 71.2 | 68.5 | 69.8 | 73.3 | 72.4 | 72.8 | 72.6 | 68.0 | 70.2 | 70.9 |
| Maverick | 70.8 | 74.3 | 72.7 | 74.5 | 76.8 | 75.6 | 71.3 | 72.4 | 71.8 | 73.4 |
| Fastcoref | 76.8 | 78.3 | 77.5 | 79.4 | 76.6 | 78.0 | 78.6 | 76.5 | 77.5 | 77.7 |
| Llama3.1(CoT) | 80.9 | 79.0 | 79.9 | 80.6 | 80.1 | 80.3 | 77.9 | 80.7 | 79.3 | 79.3 |
| GPT-o1 | 86.3 | 82.4 | 84.3 | 87.6 | 82.1 | 84.8 | 86.5 | 80.8 | 83.6 | 84.2 |
| Human | **93.6** | **89.1** | **91.3** | **94.0** | **90.3** | **92.1** | **93.2** | **88.3** | **90.7** | **91.4** |

Table 5: Overall Performance on NovelCR-en (%).

| Methods | B3 | | | MUC | | | CEAFe | | | CoNLL |
|---|---|---|---|---|---|---|---|---|---|---|
| | P | R | F | P | R | F | P | R | F | F |
| e2e-coref | 53.2 | 62.3 | 57.4 | 59.3 | 58.6 | 58.9 | 59.4 | 56.8 | 58.1 | 58.1 |
| c2f-coref | 58.3 | 68.8 | 63.1 | 66.4 | 67.1 | 66.7 | 67.3 | 64.8 | 66.0 | 65.3 |
| CR-BERT | 62.7 | 70.8 | 66.5 | 68.6 | 67.2 | 67.9 | 63.9 | 69.1 | 66.4 | 67.0 |
| SpanBERT | 68.1 | 67.4 | 67.7 | 72.4 | 65.8 | 69.0 | 67.4 | 70.2 | 68.8 | 68.5 |
| WL-COREF | 60.7 | 63.3 | 62.0 | 68.5 | 63.7 | 66.0 | 64.7 | 62.2 | 63.4 | 63.8 |
| Link-Append | 58.9 | 67.2 | 62.8 | 65.4 | 67.1 | 66.2 | 63.0 | 66.7 | 64.8 | 64.6 |
| LongDoc | 67.6 | 68.7 | 68.1 | 70.4 | 68.0 | 69.2 | 66.8 | 68.3 | 67.5 | 68.3 |
| Fastcoref | 67.9 | 68.1 | 68.0 | 69.5 | 67.3 | 68.4 | 68.3 | 64.7 | 66.5 | 67.6 |
| LLaMa3.1(CoT) | 70.6 | 71.8 | 71.2 | 73.4 | 72.6 | 73.0 | 72.8 | 70.6 | 71.7 | 72.0 |
| GPT-o1 | 73.1 | 72.5 | 72.8 | 74.2 | 73.4 | 73.8 | 72.5 | 71.9 | 72.1 | 72.9 |
| Human | **96.3** | **85.1** | **90.4** | **94.3** | **86.8** | **90.4** | **95.4** | **86.2** | **90.6** | **90.5** |

Table 6: Overall Performance on NovelCR-zh (%).

using the CoNLL metric, demonstrating the high quality of NovelCR. (2) Current CR baselines still suffer from a performance gap compared to human beings, with the state-of-the-art model(GPT-o1) achieving 84.2% F1 score on NovelCR-en and 72.9% F1 score on NovelCR-zh, about 10% lower than the scores of human evaluations. Humans can not only utilize extensive world knowledge to infer coreference relationships, but also possess strong logical reasoning abilities, capable of handling complex scenarios such as indirect references and implicit information. Therefore, humans achieve better results than current CR models in this regard.

### 5.3.1 Short-Span or Long-Span

In this section, we provide a focused analysis of how existing CR models perform under different span lengths, with particular attention to the contrast between short-span and long-span coreference cases. Specifically, we categorize the coreference pairs in NovelCR-en into three groups based on the sentence-level distance between mentions: pairs occurring within fewer than three sentences (*<3*), between three and five sentences (*3–5*), and those

spanning more than five sentences (*>5*). To ensure a comprehensive and reliable assessment, we adopt Fastcoref, LLaMA 3.1, and GPT-o1 as representative baselines.

| Sent. | <3 | 3-5 | >5 |
|---|---|---|---|
| Fastcoref | 88.9 | 79.2 | 69.8 |
| Llama3.1 | 93.5 | 76.1 | 73.9 |
| GPT-o1 | 92.8 | 82.8 | 76.3 |

Table 7: Short-Span VS Long-Span(%).

From the results presented in Table 7, we observe a clear trend: as the distance between coreferent mentions increases, from <3 sentences, 3-5 sentences, to >5 sentences, the performance of existing SOTA models (Fastcoref, LLaMA3.1, and GPT-4o) degrades substantially. Specifically, Fastcoref's F1 score declines from 88.9% to 69.8%, LlaMA3.1's F1 score drops from 95.6% to 74.4%, and GPT-o1 exhibits a similar pattern, with the F1 score decreasing from 92.8% to 76.3%. The observed performance degradation underscores that long-span coreference resolution remains an unsolved problem, motivating the introduction of NovelCR to facilitate the development of improved solutions.

| Error Types | Examples |
|---|---|
| Closest Selection | Jerebal, are you still expecting him to save you? Today is the day that he gets married! He is in love – do you really expect that **you** would even cross his mind?!" Quila cried.<br>Predict: Quila  Golden: Jerebal |
| Gender Confusion | Dad, you should mind your own business, she said. Don't say that to father, a little boy said. See what a sweet daughter you've got, the man's wife said.<br>Predict: a little boy  Golden: a sweet daughter |
| Multiple Entities | Emma said "I am not the killer, and I think it was James that killed Mason". "I didn't do that. I saw Oliver last night. It must be him". "No **you** are lying. Oliver does not hate Mason, and we all know that.", Ava said.<br>Predict: Mason  Golden: James |

Table 8: Error Analysis in NovelCR.

## 5.4 Error Analysis

In this section, we analyze common errors in NovelCR.

One of the common errors is the nearest selection. Existing CR models often simply and rudely believe that a mention refers to its closest entity. For instance, in the first example in Table 8, existing CR models do not take context into account and mistakenly assume that the mention *you* refers to the closer entity *Quila*, rather than the farther but correct entity *Jerebal*.

Another common error in NovelCR is that existing CR models lack the common sense to discern the gender of the mention. For instance, in the second example in Table 8, existing CR models fail to understand that the pronoun of *she* should be female rather than male, which leads to the model incorrectly resolving *she* to *a little boy* instead of *a sweet daughter*.

The third common error in NovelCR is that existing CR models will be very confused if there are too many entities surrounding the mention in the text. For instance, in the third example in Table 8, there are numerous entities in the text, including *Emma, James, Mason, Oliver, Ava*. Faced with so many choices, it is difficult for existing CR models to understand that *you* here refers to *James* rather than *Emma, Mason, Oliver, Ava*.

## 6 Conclusion

We propose NovelCR, a large-scale bilingual benchmark designed for long-span coreference resolution. NovelCR features a substantial dataset size and contains numerous lengthy coreferences. Extensive experiments on NovelCR demonstrate that the performance of the state-of-the-art baselines cannot catch up with human beings, showing that NovelCR remains an unresolved challenge.

## 6.1 Limitations

While we have made significant strides in constructing a high-quality CR dataset, it is important to acknowledge the limitations that may affect the interpretation and generalizability of our work.

**Few Entity Types** As outlined in the introduction, we concentrate on resolving coreferences of characters in the novel. This is a double-edged choice. On one side, it enables NovelCR to contain abundant long-span coreferences. On the other side, it restricts NovelCR's entity type exclusively to persons, omitting locations, organizations, times, events, and others. The restricted entity type compromises NovelCR's diversity and constrains NovelCR's applicability across diverse natural language understanding contexts. Future endeavors could explore extracting more long-span coreferences for additional entity types from varied data sources.

**Missing Mention Annotation** As described in Section 3.2, our mention detection process follows a two-stage pipeline: we first apply automatic tools and models to pre-identify candidate mentions, followed by manual filtering of invalid spans. This approach improves annotation efficiency and precision but may overlook a small number of valid mentions. To assess the impact, we manually examined 200 documents each from NovelCR-en and NovelCR-zh, identifying a missing mention rate of 0.9% and 1.1%, respectively. Given the large scale of our dataset, this level of omission is minimal and acceptable.

## References

1995. *Sixth Message Understanding Conference (MUC-6): Proceedings of a Conference Held in Columbia, Maryland, November 6-8, 1995*.

Ron Artstein and Massimo Poesio. 2008. Inter-coder

agreement for computational linguistics. *Computational Linguistics*, 34(4):555–596.

Amit Bagga and Breck Baldwin. 1998. Algorithms for scoring coreference chains. In *The first international conference on language resources and evaluation workshop on linguistics coreference*, volume 1, pages 563–566. Citeseer.

David Bamman. 2020. Booknlp: Natural language processing for long-form text. In *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing: System Demonstrations*, pages 1–7. Association for Computational Linguistics.

David Bamman, Olivia Lewke, and Anya Mansoor. 2020. An annotated dataset of coreference in English literature. In *Proceedings of the Twelfth Language Resources and Evaluation Conference*, pages 44–54, Marseille, France. European Language Resources Association.

Bernd Bohnet, Chris Alberti, and Michael Collins. 2022. Coreference resolution through a seq2seq transition-based system.

Arthur Brack, Daniel Uwe Müller, Anett Hoppe, and Ralph Ewerth. 2021. Coreference resolution in research papers from multiple domains. *CoRR*, abs/2101.00884.

Hong Chen, Zhenhua Fan, Hao Lu, Alan L Yuille, and Shu Rong. 2018. Preco: A large-scale dataset in preschool vocabulary for coreference resolution. *arXiv preprint arXiv:1810.09807*.

Agata Cybulska and Piek Vossen. 2014. Using a sledgehammer to crack a nut? lexical diversity and event coreference resolution. In *Proceedings of the Ninth International Conference on Language Resources and Evaluation (LREC'14)*, pages 4545–4552, Reykjavik, Iceland. European Language Resources Association (ELRA).

Ernest Davis, Leora Morgenstern, and Charles L Ortiz. 2017. The first winograd schema challenge at ijcai-16. *AI Magazine*, 38(3):97–98.

Vladimir Dobrovolskii. 2021. Word-level coreference resolution. In *Proceedings of the 2021 Conference on Empirical Methods in Natural Language Processing*, pages 7670–7675, Online and Punta Cana, Dominican Republic. Association for Computational Linguistics.

George R Doddington, Alexis Mitchell, Mark A Przybocki, Lance A Ramshaw, Stephanie M Strassel, and Ralph M Weischedel. 2004. The automatic content extraction (ace) program-tasks, data, and evaluation. In *Lrec*, volume 2, pages 837–840. Lisbon.

Abbas Ghaddar and Philippe Langlais. 2016. Wikicoref: An english coreference-annotated corpus of wikipedia articles. In *Proceedings of the Tenth International Conference on Language Resources and Evaluation (LREC'16)*, pages 136–142.

Lynette Hirschman. 1997. Muc-7 coreference task definition, version 3.0. *Proceedings of MUC-7, 1997*.

Mandar Joshi, Danqi Chen, Yinhan Liu, Daniel S. Weld, Luke Zettlemoyer, and Omer Levy. 2019a. Spanbert: Improving pre-training by representing and predicting spans. *CoRR*, abs/1907.10529.

Mandar Joshi, Omer Levy, Daniel S. Weld, and Luke Zettlemoyer. 2019b. Bert for coreference resolution: Baselines and analysis.

Fariba Karimi, Claudia Wagner, Florian Lemmerich, Mohsen Jadidi, and Markus Strohmaier. 2016. Inferring gender from names on the web: A comparative evaluation of gender detection methods. In *Proceedings of the 25th International conference companion on World Wide Web*, pages 53–54.

Kenton Lee, Luheng He, Mike Lewis, and Luke Zettlemoyer. 2017. End-to-end neural coreference resolution. *arXiv preprint arXiv:1707.07045*.

Kenton Lee, Luheng He, and Luke Zettlemoyer. 2018. Higher-order coreference resolution with coarse-to-fine inference. *arXiv preprint arXiv:1804.05392*.

Hector Levesque, Ernest Davis, and Leora Morgenstern. 2012. The winograd schema challenge. In *Thirteenth international conference on the principles of knowledge representation and reasoning*.

Xiaoqiang Luo. 2005. On coreference resolution performance metrics. In *Proceedings of Human Language Technology Conference and Conference on Empirical Methods in Natural Language Processing*, pages 25–32, Vancouver, British Columbia, Canada. Association for Computational Linguistics.

Giuliano Martinelli, Edoardo Barba, and Roberto Navigli. 2024. Maverick: Efficient and accurate coreference resolution defying recent trends. *arXiv preprint arXiv:2407.21489*.

Mary L McHugh. 2012. Interrater reliability: the kappa statistic. *Biochemia medica*, 22(3):276–282.

OpenAI. 2024. Openai o1 system card. Accessed: 2025-05-31.

Shon Otmazgin, Arie Cattan, and Yoav Goldberg. 2022a. F-coref: Fast, accurate and easy to use coreference resolution.

Shon Otmazgin, Arie Cattan, and Yoav Goldberg. 2022b. Lingmess: Linguistically informed multi expert scorers for coreference resolution. *arXiv preprint arXiv:2205.12644*.

Long Ouyang, Jeff Wu, Xu Jiang, Diogo Almeida, Carroll L. Wainwright, Pamela Mishkin, Chong Zhang, Sandhini Agarwal, Katarina Slama, Alex Ray, John Schulman, Jacob Hilton, Fraser Kelton, Luke Miller, Maddie Simens, Amanda Askell, Peter Welinder, Paul Christiano, Jan Leike, and Ryan Lowe. 2022. Training language models to follow instructions with human feedback.

Sameer Pradhan, Alessandro Moschitti, Nianwen Xue, Olga Uryupina, and Yuchen Zhang. 2012. Conll-2012 shared task: Modeling multilingual unrestricted coreference in ontonotes. In *Joint conference on EMNLP and CoNLL-shared task*, pages 1–40.

Sameer Pradhan, Lance Ramshaw, Mitchell Marcus, Martha Palmer, Ralph Weischedel, and Nianwen Xue. 2011. CoNLL-2011 shared task: Modeling unrestricted coreference in OntoNotes. In *Proceedings of the Fifteenth Conference on Computational Natural Language Learning: Shared Task*, pages 1–27, Portland, Oregon, USA. Association for Computational Linguistics.

Altaf Rahman and Vincent Ng. 2012. Resolving complex cases of definite pronouns: the winograd schema challenge. In *Proceedings of the 2012 Joint Conference on Empirical Methods in Natural Language Processing and Computational Natural Language Learning*, pages 777–789.

Keisuke Sakaguchi, Ronan Le Bras, Chandra Bhagavatula, and Yejin Choi. 2021. Winogrande: An adversarial winograd schema challenge at scale. *Communications of the ACM*, 64(9):99–106.

Kumar Shridhar, Nicholas Monath, Raghuveer Thirukovalluru, Alessandro Stolfo, Manzil Zaheer, Andrew McCallum, and Mrinmaya Sachan. 2022. Longtonotes: Ontonotes with longer coreference chains. *arXiv preprint arXiv:2210.03650*.

Hugo Touvron, Thibaut Lavril, Gautier Izacard, Xavier Martinet, Marie-Anne Lachaux, Timothée Lacroix, Baptiste Rozière, Naman Goyal, Eric Hambro, Faisal Azhar, Aurelien Rodriguez, Armand Joulin, Edouard Grave, and Guillaume Lample. 2023. Llama: Open and efficient foundation language models.

Olga Uryupina, Ron Artstein, Antonella Bristot, Federica Cavicchio, Kepa Rodriguez, and Massimo Poesio. 2016. Arrau: Linguistically-motivated annotation of anaphoric descriptions. In *Proceedings of the Tenth International Conference on Language Resources and Evaluation (LREC'16)*, pages 2058–2062.

Marc Vilain, John Burger, John Aberdeen, Dennis Connolly, and Lynette Hirschman. 1995. A model-theoretic coreference scoring scheme. In *Sixth Message Understanding Conference (MUC-6): Proceedings of a Conference Held in Columbia, Maryland, November 6-8, 1995*.

Kellie Webster, Marta Recasens, Vera Axelrod, and Jason Baldridge. 2018a. Mind the gap: A balanced corpus of gendered ambiguou. In *Transactions of the ACL*, page to appear.

Kim Webster, Kristin Diemer, Nikki Honey, Samantha Mannix, Justine Mickle, Jenny Morgan, Alexandra Parkes, Violeta Politoff, Anastasia Powell, Julie Stubbs, et al. 2018b. *Australians' attitudes to violence against women and gender equality*. Australia's National Research Organisation for Women's Safety.

Jason Wei, Xuezhi Wang, Dale Schuurmans, Maarten Bosma, Ed H. Chi, Quoc Le, and Denny Zhou. 2022. Chain of thought prompting elicits reasoning in large language models. *CoRR*, abs/2201.11903.

Ralph Weischedel, Sameer Pradhan, Lance Ramshaw, Martha Palmer, Nianwen Xue, Mitchell Marcus, Ann Taylor, Craig Greenberg, Eduard Hovy, Robert Belvin, et al. 2011. Ontonotes release 4.0. *LDC2011T03, Philadelphia, Penn.: Linguistic Data Consortium*.

Mingzhu Wu, Nafise Sadat Moosavi, Dan Roth, and Iryna Gurevych. 2020. Coreference reasoning in machine reading comprehension. *CoRR*, abs/2012.15573.

Yufang Wu, Daniel Khashabi, and Dan Roth. 2021. Learning to ignore: Long document coreference with bounded memory neural networks. In *Proceedings of the 2021 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*, pages 2923–2933. Association for Computational Linguistics.

Xintong Yu, Hongming Zhang, Ruixin Hong, Yangqiu Song, and Changshui Zhang. 2022. Vd-pcr: Improving visual dialog with pronoun coreference resolution. *Pattern Recognition*, 125:108540.

Klim Zaporojets, Johannes Deleu, Chris Develder, and Thomas Demeester. 2020. DWIE: an entity-centric dataset for multi-task document-level information extraction. *CoRR*, abs/2009.12626.

Amir Zeldes. 2017. The GUM corpus: Creating multilayer resources in the classroom. *Language Resources and Evaluation*, 51(3):581–612.

Dmitry Zelenko, Chinatsu Aone, and Jason Tibbetts. 2004. Coreference resolution for information extraction. In *Proceedings of the Conference on Reference Resolution and Its Applications*, pages 24–31, Barcelona, Spain. Association for Computational Linguistics.

Jieyu Zhao, Tianlu Wang, Mark Yatskar, Vicente Ordonez, and Kai-Wei Chang. 2018. Gender bias in coreference resolution: Evaluation and debiasing methods. *arXiv preprint arXiv:1804.06876*.

## A Annotation Guideline of Mention Verification

In mention verification, the annotation instructions are outlined below.

Please read the novel chapter, and finish the two tasks: (1)delete invalid mentions, and (2)re-edit the mention boundaries. The second task can only be started after the first task is completed.

When deleting invalid mentions, you should remove mentions that do not refer to the person entities, such as *the bank* and *this beautiful knife*. Note

that dependent personal pronouns should also be deleted. For instance, *her* in *her split lips* is also an invalid mention since it functions as a modifier of *lips*. To delete invalid mentions, click the mention to highlight it and then click the *Delete* button.

When re-editing the boundary of the mentions, we follow the maximum span principle. This means that you should identify the longest string representing the mention. For instance, in the sentence *the sad man is looking for his wife*, you should annotate the mention as *the sad man* rather than just *man*. If the mention does not meet the maximum span criteria, you should drag the gray border to correct the boundaries of the mention. Please do nothing if no mistakes are found. When you have completed all annotations on a page, remember to click the *Submit* button to store the annotation results. We assure you that all annotations will be utilized solely for research purposes.

| Datasets | NovelCR-en | NovelCR-zh |
|---|---|---|
| | Recall | |
| POS+Semantic Parser | 91.3 | 90.7 |
| Prompt-Learning(ours) | 99.1 | 98.9 |

Table 9: Candidate Mention Detection Performance(%).

## B  Annotation Guideline of Coreference Identification

As shown in Figure 5, annotators need to match entities and mentions. The annotation instructions are as follows.

Please read the novel chapter and match each mention to the entity it refers to. We recommend reading the entire chapter before starting any annotations, as coreference resolution relies on a broad context span understanding. We already highlight mentions in grey and list the entity options at the top of the chapter. All you need to do is click the mention and then the entity it refers to to match them. If the mention doesn't refer to any entities, you can simply click on the *None* option. When you have completed all annotations on a page, remember to click the *Submit* button to store the annotation results. We promise that all annotations will be used for research purposes.

## C  Prompt for Mention Detection.

We leverage direct prompting, chain-of-thought (CoT) prompting, and ReAct prompting, respectively, to detect mentions from the novel chapter.
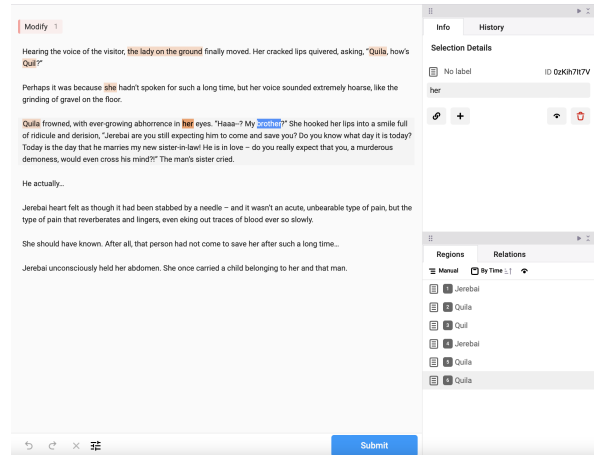


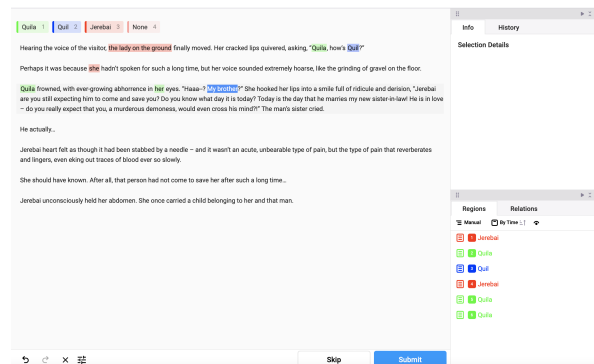Figure 4: Screenshot of Mention Verification.



Figure 5: Screenshot of Coreference Identification.

**Prompt 1 (direct prompting)**

*Question*:Please find all words or phrases that may refer to a person in the following passage:[novel chapter].

**Prompt 2 (CoT prompting)**

*Question*: Please find all words or phrases that may refer to a person in the following passage:[novel chapter]

*Thought*:The possible candidates include pronouns, human names, and noun phrases. pronouns could be he, she, him, her, their, and them. Noun phrases could be nouns like man, woman, girl, and boy with their adjectives. Human names can be discovered using the rules of different languages.

**Prompt 3 (ReAct prompting)**

*Tools*:NER(p) takes a passage as parameter and returns Named Entities that belong to human beings. PosTag(p) takes a passage as parameter and returns all pronouns and nouns phrases.

*Question*: Please find all words or phrases that may refer to a person in the following passage: [novel chapter]

*Thought*:The possible candidates include pro-

nouns, human names, and noun phrases. Human
names can be found by NER first.

    *Action*:NER
    *Observation*: [entities]
    *Thought*:Then noun phrases and pronouns can
be found by PosTag.
    *Action*:PosTag
    *Observation*: [mentions]