# *"Let Your Characters Tell Their Story"*: A Dataset for Character-Centric Narrative Understanding

**Faeze Brahman**[1]     **Meng Huang**[2]     **Oyvind Tafjord**[3]
**Chao Zhao**[5]     **Mrinmaya Sachan**[4]     **Snigdha Chaturvedi**[5]

[1]University of California, Santa Cruz, [2]University of Chicago
[3]Allen Institute for AI, [4]ETH Zurich, [5]UNC Chapel Hill

`fbrahman@ucsc.edu`, `huangme@uchicago.edu`

## Abstract

When reading a literary piece, readers often make inferences about various characters' roles, personalities, relationships, intents, actions, etc. While humans can readily draw upon their past experiences to build such a character-centric view of the narrative, *understanding* characters in narratives can be a challenging task for machines. To encourage research in this field of character-centric narrative understanding, we present LiSCU – a new dataset of literary pieces and their summaries paired with descriptions of characters that appear in them. We also introduce two new tasks on LiSCU: *Character Identification* and *Character Description Generation*. Our experiments with several pre-trained language models adapted for these tasks demonstrate that there is a need for better models of narrative comprehension.[1]

## 1 Introduction

Previous works in literary analysis have discussed that the development of the plot and the main character(s) are among the most important components that contribute to a good piece of fiction (Kennedy and Gioia, 1983; Card, 1999). In particular, *character(s)* are central to narratives since their motivations, traits, and actions determine the flow of the plot. Hence, understanding and critically analyzing characters is an important facet of literary scholarship.

In Computational Narratives, prior work has exploited the potential of character-centric natural language understanding (Chambers, 2013; Chaturvedi et al., 2017; Chu et al., 2018; Zhang et al., 2019). However, these works are limited to only understanding certain aspects of characters and do not do an in-depth and systematic study.

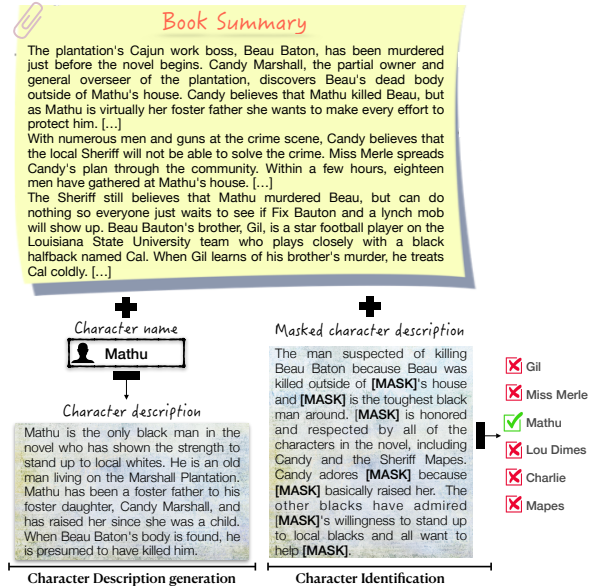To facilitate character-centric narrative understanding, we present LiSCU – a new dataset in



Figure 1: An illustration of the proposed dataset and the two tasks: *Character Description Generation* and *Character Identification*.

English, of literary pieces and their summaries paired with descriptions of characters that appear in them. These descriptions analyze the narrative from the perspective of the character highlighting their salient attributes, their role and contribution to the development of the narrative's plot.

Using this dataset, we devise two new tasks: (1) a *Character Identification* task to identify the character's name from an anonymized character description given the literature summary; and (2) a *Character Description Generation* task to generate the description for a given character of a literature summary. Our primary task, *Character Description Generation*, is related, but not identical to summarization. There are two main differences. Summarization typically has a one-to-one correspondence between documents and summaries, and focuses on copying (either extractively or abstractively) important content from the documents to create the summaries. On the other hand, character

---

descriptions are analysis, not merely summaries, of narratives from the character's point of view. They are created by abstracting out the low-level content of the narrative instead of simply identifying and paraphrasing important details. They describe events, roles, relationships, and salient attributes of the character that can be inferred from the narrative and might not be directly stated in the text. In particular, if the narrative describes several events where a character helps the protagonist, the character description will not simply mention all those events, but will instead describe the character as a helpful person (attribute) and a good friend of the protagonist (role). For example, in Fig. 1, "Mathu is virtually her foster father." in summary is expressed as "Candy adores Mathu because he basically raised her." in the character description. Thus, the *Character Description Generation* task, provides a unique opportunity for NLP systems to learn to *abstract* and model long-range dependency instead of simply *extracting* information.

Apart from this novel *abstraction* task, the dataset also poses another challenge for NLP systems by requiring them to process long documents. The average number of tokens in our summaries are 1022 which is beyond the comfort level of most existing systems. Understanding long narratives and modeling long contexts are new frontiers for NLP research (Roy et al., 2021; Fan et al., 2021) and LiSCU pushes us in this direction. To further facilitate research in this direction, we also release a small dataset where the goal is to read the entire literary piece and generate character descriptions.

We explore the ability of the modern neural models on both tasks. We demonstrate through experiments that although existing models can identify characters reasonably well in masked descriptions, there is still a scope for improvement considering human accuracy on this task. Also, while existing models can generate fluent and logically self-consistent text, they are not always faithful to the literature summaries and fail to capture salient details about the characters. Our contributions are:

- A new dataset of literature summaries paired with character descriptions to enable character-centric narrative understanding.

- A comprehensive human study to assess the quality of the proposed dataset.

- Novel tasks: a classification and an abstractive generation task to better understand characters in the narrative plot.

- Experiments with several strong baselines and a thorough qualitative analysis.

## 2 Background

The field of computational narrative understanding studies how to algorithmically represent, understand, and generate stories. Early computational studies on narratives had focused on learning procedural scripts and event sequences (Schank and Abelson, 1977; Manshadi et al., 2008; Regneri et al., 2010), narrative chains or schemas (Chambers and Jurafsky, 2008, 2009), and plot units (Goyal et al., 2010; McIntyre and Lapata, 2010; Elsner, 2012).

Computational linguists have also worked on character-centric modelling of narratives (Chambers, 2013). The character-centric perspective aims to understand characters – their personas, roles, goals, relationships, emotions, etc. Previous works have proposed methods to detect characters and infer latent personas in movie plot summaries and fictional novels (Bamman et al., 2013, 2014; Vala et al., 2015; Flekova and Gurevych, 2015), model inter-character relationships (Iyyer et al., 2016; Srivastava et al., 2016; Chaturvedi et al., 2017; Kim and Klinger, 2019), and emotions (Brahman and Chaturvedi, 2020). Earlier works have also considered constructing social networks of characters (Agarwal et al., 2014) from novels (Elson et al., 2010; Elsner, 2012) and films (Krishnan and Eisenstein, 2015).

Another line of work related to ours is on summarization of novels (Mihalcea and Ceylan, 2007). This work built a dataset of novel-summary pairs and used unsupervised summarization models such as *TextRank* (Mihalcea and Tarau, 2004) and *MEAD* (Radev, 2001). Instead of summarizing full novels, Ladhak et al. (2020) proposed a content-selection approach to create a gold-standard set of extractive summaries by aligning chapter sentences with abstractive summary sentences.

In a more related work, Zhang et al. (2019) collected a dataset of fictional stories along with author-written summaries. They proposed an extractive ranking and a classification approach to select a subset of salient attributes from a list of candidate attributes (extracted from the story) that describe a character's personality. While this work presented a collection of personality-related phrases as a potential summary for the actual novel, our dataset contains literature summaries and char-

acter descriptions, and we aim to generate natural language texts that analyze the narrative from the perspective of the characters. Such an analysis is more in-depth than a collection of phrases.

## 3 The LiSCU Dataset

We now describe our **Li**terature **S**ummary and **C**haracter **U**nderstanding (LiSCU) dataset. LiSCU is a dataset of literature summaries paired with descriptions of characters that appear in the summaries. Fig. 1 shows an example of our dataset.

Next, we describe the data collection pipeline for LiSCU (§3.1), followed by details on the reproducibility of the data collection process (§3.2).

### 3.1 Data Collection and Filtering

We collected LiSCU from various online study guides such as shmoop,[2] SparkNotes,[3] CliffsNotes,[4] and LitCharts.[5] These sources contain educational material to help students study for their literature classes. These study guides include summaries of various literary pieces as well as descriptions of characters that appear in them. These literature summaries and character descriptions were written by literary experts, typically teachers, and are of high pedagogical quality.

We used Scrapy,[6] a free and open-source web-crawling framework to crawl these study guides. Our initial crawl resulted in a set of $1,774$ literature summaries and $25,525$ character descriptions. These included all characters mentioned in the literary pieces. However, not all characters, especially those that played a minor role in the literary piece, appeared in the corresponding literature summaries. Since our task involves making inferences about characters from the literature summaries, we filtered out the characters which do not appear in the summaries or their names or the descriptions had very little overlap with the literature summaries. This is done to mitigate the reference divergence issue (Kryscinski et al., 2019; Maynez et al., 2020) and ensure that the literature summary has enough information about the character to generate the description. For this, we define the "information overlap" between two pieces of text $\mathcal{A}$ and $\mathcal{B}$, $IO(\mathcal{B}||\mathcal{A})$, as the ratio of the length of the

| | |
|---|---|
| # unique books | 1,220 |
| # literature summaries | 1,708 |
| # characters | 9,499 |
| # characters with accompanying full book | 2,052 |
| # unique books with full-text | 204 |
| avg. # characters per summary | 5.56 |
| min. # characters per summary | 1 |
| max. # characters per summary | 38 |
| avg. summary length (in tokens) | 1,022.32 |
| avg. # sentences in summary | 48.82 |
| avg. character description length (in tokens) | 184.57 |
| avg. # sentences in description | 8.56 |
| # characters in Train set | 7,600 |
| # characters in Test set | 957 |
| # characters in Validation set | 942 |

Table 1: Statistics of the LiSCU dataset.

longest overlapping word sub-sequence between $\mathcal{A}$ and $\mathcal{B}$, over the length of $\mathcal{A}$.[7] Note that this information overlap measure is not symmetric and intuitively measures how much information about $\mathcal{A}$ is present in $\mathcal{B}$. We used the information overlap measure to filter our dataset as follows. If the information overlap of the literature summary with the character name, $IO(\text{literature summary} || \text{character name})$, is less than $0.6$, then we consider that the character is not prominently mentioned in the literature summary and we remove that character from our dataset. Similarly, if the information overlap between the character description and the literature summary, $IO(\text{literature summary} || \text{character description})$, is less than $0.2$, then we consider the character description generation less feasible and we remove that data point from our dataset.[8]

However, during these filtering steps, we did not want to remove the most important characters of the narrative. The online study guides list characters in decreasing order of their importance in the literary piece. For example, narrators, protagonists, antagonists, etc., are always described first. Leveraging this ordering, we always retained the top 3 characters of the literary piece in our dataset.

After the filtering process, our final dataset consists of $1,708$ literature summaries and $9,499$ character descriptions in total. This set was split into train ($80\%$), test ($10\%$), and validation ($10\%$) sets.

---

[2] https://www.shmoop.com/study-guides/literature

[3] https://www.sparknotes.com/lit/

[4] https://www.cliffsnotes.com/literature

[5] https://www.litcharts.com

[6] https://scrapy.org/

[7] Technically this is the same as Rouge-L precision

[8] These thresholds were chosen by experimenting with different values and manually analyzing the quality of (a subset of) the data.

The data splits were created to avoid any data-leakages – each literary piece and all of its character descriptions were consistently part of only one of the train, test and validation sets. Table 1 shows the statistics of the final dataset. The dataset also contains the full-text of the books for $2,052$ of the character descriptions.

## 3.2 Dataset Reproducibility

LiSCU is drawn from various study guides on the web. While we do not have the rights to directly redistribute this dataset, to allow other researchers to replicate the LiSCU dataset and compare to our work, we provide a simple script that will allow others to recreate LiSCU from a particular time-stamped version of these study guides on *Wayback Machine*, a time-stamped digital archive of the web. Our script ensures that others will be able to recreate the same train, test and validation splits.

## 4 LiSCU Task Definitions

We introduce two new tasks on the LiSCU dataset:
- *Character Identification*
- *Character Description Generation*

### 4.1 Character Identification

The *Character Identification* task requires models to identify the character in an anonymized character description. Given a summary $S$, a candidate list of characters that appear in the literature summary $C = \{c_1, c_2, ..., c_k\}$, and an anonymized character description $D^{c*}_{masked}$, the goal in this task is to identify the name of the character $c^*$ described in the anonymized character description. We anonymize character descriptions by masking out all mentions of the character $c^*$ in the original description $D^{c*}$.

### 4.2 Character Description Generation

The *character description generation* task tests the ability of NLP models to critically analyze the narrative from the perspective of characters and generate coherent and insightful character descriptions. Formally, given a literature summary, $S$, and a character name, $c$, the goal in this task is to generate the character's description, $D^c$. Generating the character description necessitates understanding and analyzing every salient information about the character in the literature summary.

### 4.3 Human Assessment of LiSCU

In order to verify the tractability of these two tasks as well as assessing the quality of the collected



Figure 2: Human assessment of the feasibility of the character description generation task.

LiSCU dataset, we conducted a set of human evaluations on Amazon Mechanical Turk. We run our human assessment on the full test set of LiSCU.

**Assessing the Character Identification task:** In the first human assessment, we showed annotators the literature summaries, anonymized character descriptions , and a list of character names (plus one randomly sampled character from the literary piece). The descriptions were anonymized by replacing all mentions of the corresponding character names with blanks.[9] For each anonymized character description, we asked 3 judges to identify which character it is describing by choosing from the list of choices. The judges also had the option of saying that they are unable to identify the character given the literature summary and the anonymized character description.

**Assessing the Character Description Generation task:** In the second human assessment, the judges are shown the same summary along with the original de-anonymized character descriptions. For each character description, 3 judges were asked to evaluate the quality of the description by answering the following two questions:

1. **Fact coverage:** Specify how much of the information about the specific character in the corresponding "character description" is present in the summary (either explicitly or implicitly). Answer choices included: a) *almost all of the information*, b) *most of the information*, c) *some of the information*, d) *little or none of the information*, and e) *character does not appear in the summary at all*.

2. **Task difficulty:** Given the summary, how easy is it to write the character description on a Likert scale of 0-4 (0 being too difficult, 4 being too

---

[9] We identified mentions of a character in the summary by using a coreference system (Joshi et al., 2019b,a) as well as by matching the first name or the full name of the character.
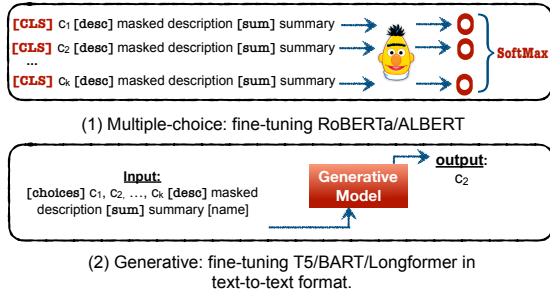
(1) Multiple-choice: fine-tuning RoBERTa/ALBERT



(2) Generative: fine-tuning T5/BART/Longformer in text-to-text format.

Figure 3: Approaches for *Character Identification*.

| Model | Description Setup | Accuracy (%) |
|---|---|---|
| *Random Guess* | - | 18.70 |
| RoBERTa-Large (Liu et al., 2019) | partial | 77.84 |
| ALBERT-XXL (Lan et al., 2020) | partial | **83.33** |
| T5-11B (Raffel et al., 2020a) | partial | <u>80.16</u> |
| BART-Large (Lewis et al., 2019) | partial | 74.89 |
| Longformer (Beltagy et al., 2020) | partial | 71.10 |
| BART-Large (Lewis et al., 2019) | full | 78.58 |
| Longformer (Beltagy et al., 2020) | full | 74.78 |
| *Human Performance* | - | 91.80 |

Table 2: Accuracy for the *Character Identification*. The 'partial' description setup used a truncated description (50 words) to allow including more of the summary.

easy)? If in the previous question the judges found that some of the information in the character description was not present in the summary, they are asked to disregard that while answering this question. In other words, they only need to consider the information in the character description which is explicitly or implicitly mentioned in the summary.

We recruited 200 crowd-workers who were located in the US, UK, or CA, and had a 98% approval rate for at least 5,000 previous annotations. We collected each annotation from 3 workers and use majority vote in our assessments. In the Appendix A, we describe several steps we took to alleviate limitations of using crowd-sourcing and ensure high quality annotations. Screenshots of our AMT experiments are provided in the Appendix.

For the first assessment on identifying characters, the human accuracy was 91.80% (Fleiss' Kappa (Landis and Koch, 1977) $\kappa = 0.79$), indicating the feasibility of the task.

For the second assessment of fact coverage and task difficulty, we summarize the result in Fig. 2. The top chart ('Fact Coverage') shows that around 75% of the of the literature summaries contain reasonable amount of information about the character represented in the corresponding character description. The bottom chart ('Task Difficulty') shows that more than 90% of the times, the human judges considered the task of writing the character descriptions from the literature summaries not too difficult.[10]

These results verify the feasibility of understanding and drawing reasonable inferences about characters in the literature summaries from the LiSCU

dataset. Next, we describe models and establish baseline performances on the two proposed tasks.

## 5 Character Identification

We present two approaches to address this task: (1) solving it as a multiple-choice classification problem, and (2) using a generative classifier that generates, instead of identifying, the character name, as shown in Fig. 3.

In the multiple-choice approach, we use the standard setup introduced in BERT (Devlin et al., 2019) where the text from $c_i$, $D_{masked}^{c*}$ and $S$ (with custom prefix tokens) are concatenated as input, and the [CLS] token is projected to a final logit. We apply a Softmax function to the logits to obtain the scores for each $c_i$. For training practicalities, we limit the number of choices to 4 during training (using the earliest window of choices which include the correct one). During inference, we can generate the logits for all the answer choices since they are independent before the final Softmax.

To establish a baseline performance, we experiment with finetuning RoBERTa (Liu et al., 2019), and ALBERT (Lan et al., 2020) which have been shown to perform well in several classification tasks. However, both these models cannot process inputs longer than 512 tokens and the concatenated inputs are generally much longer. So we also tried Longformer (Beltagy et al., 2020), a BERT-like model with an attention mechanism designed to scale linearly with sequence length, thus allowing the model to encode longer documents. However, despite trying various hyperparameters, Longformer was not able to match the scores in our experiments.

Our second approach, a generative classifier, is

---

[10]There is a natural label bias in the annotations: most of the responses fell into few categories. In this case, standard inter-annotator agreement statistics are not reliable (the well-known paradoxes of kappa (Feinstein and Cicchetti, 1990)). Thus, we simply report a pairwise agreement (i.e., how often do two judges agree on the answer for the same question) of 0.71 and 0.64 for 'fact coverage' and 'task difficulty', respectively.

inspired by Raffel et al. (2020b) who studied transfer learning by converting NLP problems into a text-to-text format. The generative classifier addresses the character identification problem by directly generating the character name $\hat{c}$, given all character names (answer choices), the masked character description, and the summary (see Fig. 3). During inference, we compute the model's probability of each of the answer choices, and output the one with the highest probability.

We use this procedure to train several strong baselines built on top of the following pre-trained transformer-based models: BART (Lewis et al., 2019), T5 (Raffel et al., 2020b), and Longformer (Beltagy et al., 2020).

**Implementation Details.** The RoBERTa and ALBERT multiple-choice classifiers were trained for 6 epochs, initial learning rate 1e-5 (ADAM optimizer), batch size 16. The generative classifier using BART was trained for 5 epochs, initial learning rate 5e-6, batch size 8. We used the Transformer package (Wolf et al., 2019) for training. The T5 model was trained for 12 epochs on a TPU using the default parameters from the T5 repository (learning rate 1e-3 with AdaFactor, batch size 8).[11] We truncate the summaries (and descriptions) to satisfy model-specific maximum input length.

**Results.** Table 2 shows the accuracies of different baselines. The highest accuracy is achieved by ALBERT-XXL (83.33%) followed by T5-11B (80.16%). Although both ALBERT and T5 were given partial character descriptions, their specific pre-training loss and larger number of parameters (for T5-11B) lead to superior performance over other baselines. We observe that there is still a significant difference between the human performance (91.80%) and the best model performance (83.33%) on the character identification task, warranting future work on this direction.

# 6 Character Description Generation

We present several strong baselines for generating character descriptions by fine-tuning pre-trained transformer-based language models (LM) (Vaswani et al., 2017). We study two types of models: (1) a standard left-to-right LM, namely GPT2-L (Radford et al., 2019) which is trained with LM objective to predict the next word; and (2) two encoder-

| Model | BLEU | ROUGE-1 | ROUGE-2 | ROUGE-L | BERT-F1 |
|---|---|---|---|---|---|
| **Length Truncated Input** | | | | | |
| GPT2-L | 0.67 | 19.25 | 3.50 | 17.51 | 77.71 |
| BART-L | **1.38** | **24.93** | **5.42** | **21.99** | 84.54 |
| Longformer | 1.05 | 21.47 | 4.66 | 19.37 | 84.64 |
| **Coref Truncated Input** | | | | | |
| GPT2-L | 0.58 | 18.69 | 3.15 | 16.91 | 78.46 |
| BART-L | 0.96 | 21.33 | 4.66 | 19.04 | 84.26 |
| Longformer | 0.98 | 21.18 | 4.40 | 19.13 | 84.59 |
| **Full Length Input** | | | | | |
| Longformer | 1.14 | 21.79 | 4.88 | 19.60 | **84.72** |

Table 3: Automatic evaluation results for *Character Description Generation*. BART-L achieved the best BLEU and ROUGE scores while Longformer performed best on BERTScore.

decoder models, namely BART[12] (Lewis et al., 2019) and Longformer (Beltagy et al., 2020)[13] which initialize the state of the Transformer by reading the input, and learn to generate the output.

One of the challenges of the proposed task is the length of the summaries, which might exceed the maximum allowable length for most existing pre-trained models. To overcome this, we either: (1) simply truncate the literature summary at the end, or (2) only keep sentences from the literature summary that have a mention of the character of interest. For the latter, we use a coreference resolution model, SpanBERT (Joshi et al., 2019b,a), to identify character mentions within a summary. This results in a modified dataset of character-specific literature summaries paired with character descriptions. In addition to these two approaches, we also fine-tune Longformer (Beltagy et al., 2020) with original full-length literature summary. Longformer leverages an efficient encoding mechanism to avoid the quadratic memory growth and has been previously explored for NLU tasks (encoder-only). We integrate this approach into the pre-trained encoder-decoder BART model to encode inputs longer than its maximum token limit. All the models take `[name]` $c$ `[sum]` $S$ `[desc]` as input and generate the character description $D^c$ as output.

**Experiment with Full Literary Pieces.** We also run an experiment on a subset of our data with accompanying full-text of the literary pieces.

---

[11]https://github.com/google-research/
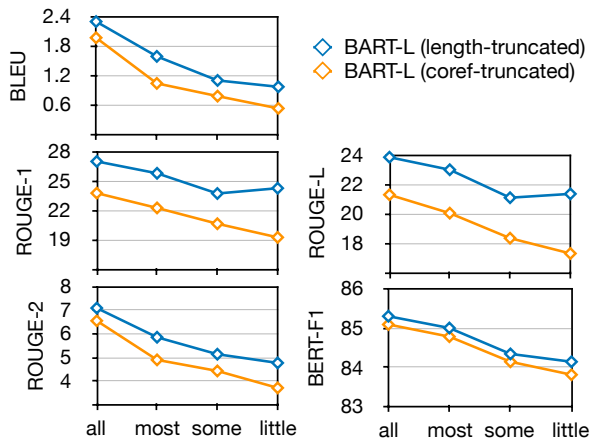text-to-text-transfer-transformer

Figure 4: Breakdown results for BART-L on subsets with annotated fact coverage as all/most/some/little. Results for other baselines are provided in Appendix.

| Model | BLEU | R-1/ R-2 /R-L | BERT-F1 |
|---|---|---|---|
| Longformer (w/ Books) | 0.73 | 17.61 /3.60 / 16.15 | 84.33 |
| Longformer (w/ Summaries) | **1.00** | **19.46 / 4.33 / 17.74** | **84.77** |

Table 4: Automatic evaluation results for models using full-text of books vs. literature summaries.

Since it is infeasible to use the full texts as input given the memory constraints of current models, we coarsely select spans of the full-text beginning $50$ tokens before, and $50$ tokens after the occurrence of character's name. We use a Longformer model where the input is simply the concatenation of the selected spans. Due to the small size of the this subset, we perform a 5-fold cross validation starting from a pre-trained model fine-tuned on summary-description pairs.[14]

**Implementation Details.** We use the Transformer library (Wolf et al., 2019). Each baseline was trained for 5 epochs with effective batch size of 8, and initial learning rate of 5e-6. We use the maximum input length of 1024 for GPT2, and 2048 for BART[15] and the variant of Longformer with truncated input. For experiment with original books, we use $16,384$ which is the maximum allowable input length for Longformer. During inference, we use beam search decoding with 5 beams.

## 6.1 Automatic Evaluation

Following previous works, we use several standard, widely used automatic evaluation metrics. We use **BLEU-4** (Papineni et al., 2002) that measures overlap of $n$-gram up to $n = 4$, **ROUGE-$n$** ($n$=1, 2), and **ROUGE-L** F-1 scores (Lin, 2004)[16] . However, recent works (Novikova et al., 2017; Wang et al., 2018) have raised concerns on the usage of

these metrics as they fail to capture paraphrases and conceptual information. To overcome these issues, we additionally include a model-based metric, **BERTScore** (Zhang et al., 2020), which measures the cosine similarity between contextualized embeddings of the gold and generated outputs.[17]

The result of the automatic evaluation is presented in Table 3. According to the table, BART-L consistently achieves the best performance across BLEU and ROUGE scores. However, Longformer achieves a slightly better BERTScore. Both BART and Longformer outperform GPT2 in general. This can be in part because BART and Longformer can handle longer context, and are initially pre-trained on a combination of books and Wikipedia data and further fine-tuned on summarization tasks, while GPT2 is pre-trained on WebText only.[18]

Models perform relatively better in the length truncation setups than in the coreference truncation. We posit that this is because a lot of the key points about major characters are likely to appear earlier in the book summary (favoring length truncation). Also, there might be errors introduced by the coreference resolution model itself.

In order to have a better insight into the models' performance with respect to varying level of task feasibility, in Fig. 4, we additionally report the breakdown of the results for BART-L on separate subsets with "almost all", "most", "some", "little or none" of the information about the character (refer to *Fact Coverage* in §4.3). As expected, we observe a consistent decline in the performance with lower amount of fact coverage. Results for other baselines are reported in Table 8 of the Appendix.

In Table 4, we compare the models when using selected spans from the original literary piece as the input vs. literature summaries as the input. We observe a decline in performance when we used the full text. This reveals that even though the literary pieces contain all the character information, this

---

[14]Pre-training data do not contain instances of this subset.

[15]BART originally accepts inputs of maximum 1024 BPE-tokens. We extend this to 2048 by adjusting its positional embeddings.

[16]Note that we did not include perplexity score as it is not comparable across LM-based and encoder-decoder models.

[17]We use the code at https://github.com/Tiiiger/bert_score

[18]While these models could have had access to the original book text, they do not have access to the character descriptions (our outputs) during pre-training. So, this information should not principally change any of our empirical conclusions.
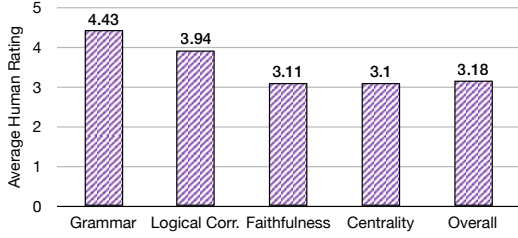
Figure 5: Human evaluation of generated character descriptions. While the descriptions are grammatically correct and logically coherent, they often misrepresent or miss important details about the character.

| Aspects | (1) | (2) | (3) | (4) | (5) |
|---|---|---|---|---|---|
| **Grammar** | 0.00 | 3.67 | 8.67 | 28.67 | 59.00 |
| **Logical Corr.** | 1.67 | 9.00 | 19.00 | 33.67 | 36.67 |
| **Faithfulness** | 12.67 | 23.67 | 21.00 | 24.67 | 18.00 |
| **Centrality** | 15.33 | 17.00 | 27.00 | 23.67 | 17.00 |
| **Overall** | 11.00 | 19.33 | 26.67 | 26.67 | 16.33 |

Table 5: Percentage of different ratings from human evaluation of generated descriptions (1=worst, 5=best).

| Error Type | Percentage |
|---|---|
| **Events** | 46.00 |
| **Role** | 24.33 |
| **Relationships** | 25.00 |
| **Personal characteristics** | 12.33 |
| **Behavioral characteristics** | 22.33 |
| **No major error** | 27.67 |

Table 6: Error Analysis: proportion of generated descriptions with different error types.

information is scattered which makes it harder for the model to identify important facts about the character. Using full texts also requires encoders which are better at understanding dialog, first-person narratives and different writing styles of the authors. We invite the community to consider this challenging but important problem.

## 6.2 Human Evaluation

To better evaluate the quality of the generated character descriptions, we conduct a human evaluation on 100 test pairs of literature summaries and character descriptions generated by the BART-L model on Amazon Mechanical Turk.[19] Given a literature summary and multiple generated character descriptions (shown one by one), the workers were asked to rate each generated description on a Likert scale of $1-5$ (1 being the worst, and 5 being the best) according to the following criteria: (1) **Grammatical correctness** to indicate if the generated description is grammatically correct, (2) **Logical correctness** to indicate whether the generated description is logically meaningful and coherent, (3) **Faithfulness** of the generated description with respect to the given summary (a faithful character description will not mention facts which are irrelevant to the character and/or not stated in the summary), (4) **Centrality** to evaluate whether the description captures important details and key facts about the character, and finally (5) the **Overall score** considering all the four criteria listed above. We provide a screenshot of the experiment in Fig. 7 of the Appendix.

Fig. 5 presents the results of this human evaluation. We observe that the generated descriptions show a reasonable level of grammatical (4.43) and logical correctness (3.94). However, they lack behind when it comes to faithfulness (3.11) and cen-

trality (3.10). We also report the distribution of ratings in Table 5. These results indicate that solving this task requires designing better models of character-centric analysis of narrative.

## 6.3 Qualitative Analysis

Here, we do a qualitative analysis for the *Character Description Generation* task. In our human evaluation of the generated character descriptions (§6.2), we additionally provided a questionnaire to collect in-depth feedback from crowd-workers on the type of errors the BART-L model made when generating character descriptions from the given literature summaries. The questionnaire asked "What details about the character does the given character description miss or describe inaccurately. Note that the description is supposed to describe only the important details and not necessarily all of them." The workers were asked to select all the applicable choices among the following error types:

1. **Events:** The character description misses or misrepresents some main event(s) that the character is involved in.
2. **Role:** The character's role in the narrative (e.g., protagonist, antagonist, etc.) is important but is not included or misrepresented in the character description.
3. **Relationships:** The character's relationship with other characters is important (e.g., the protagonist's wife) but is not included or mis-

---

| Generated Character Descriptions | Comments |
|---|---|
| **Book title:** The Three Sisters<br>**Character name:** Vershinin<br>**Generated Character Description:** Vershinin is one of the most interesting characters in Three Sisters. He's a <span style="color:green">young officer in the Russian army</span>, and <span style="color:red">he's got a lot of time on his hands.</span> But that doesn't stop him from having an affair <span style="color:red">with Irina</span>, the youngest sister, who falls in love with him because she's <span style="color:green">attracted to his philosophical ideas about the meaning of life.</span> | According to the summary "Masha", the middle sister, is the one who has an affair with the officer Vershinin, and not the youngest sister, Irina. |
| **Book title:** The Sun Also Rises<br>**Character name:** Pedro Romero<br>**Generated Character Description:** Pedro Romero is <span style="color:green">a young Spanish bull-fighter who is in love with Lady Brett Ashley, a beautiful British divorcee.</span> He's a talented and <span style="color:red">talented young man</span>, and, despite the fact that he is only <span style="color:red">twenty-five years old</span>, he has <span style="color:red">the potential to be one of the best bullfighters in Pamplona.</span> <span style="color:green">Jake feels terrible for introducing him to Brett, fearing that it has corrupted him.</span> | The description captures most of the important details about the character. |

Table 7: Examples of generated descriptions. Words in <span style="color:red">red</span> correspond to hallucinated or missing content, and words in <span style="color:green">green</span> correspond to faithful information. The input literature summaries are provided in the Appendix.

represented in the character description.

4. **Personal characteristics:** The character's personal characteristics (e.g., age, ethnicity, personality, etc.) are important for the narrative but are not included or misrepresented in the character description.

5. **Behavioral characteristics:** The character's motivation, desires, and behavior are important but are not included or misrepresented in the character description.

6. **No major error:** None of the above. The character description captures most of the important details about the character.

We also provided an optional text box for them to type in other details that are missing or misrepresented but not listed above.

The result of this analysis is shown in Table 6. We can see that the generated descriptions make fewer mistakes in capturing personality-related attributes (12.33%) and more mistakes in representing important events involving the characters (46%). They also sometimes omit or misrepresent roles (24%), relationships (25%), and behavioral characteristics (22%) of the characters. This indicates factors that future systems should consider improving upon when addressing this task.

We provide qualitative examples of the generated character descriptions along with the errors they made (as pointed out by the turkers) in Table 7. More examples with input literature summaries are provided in Tables 9 to 12 of the Appendix.

## 7 Conclusion

Understanding and critically analyzing fictional characters is an important element of understanding a literary piece. Human readers build a mental model of characters, understand what they look like, their role in the literary piece, and assess their psychology, motivations, and consequences of their behavior. However, building such a deep understanding of fictional characters in narratives is hard for machine reading systems. To encourage progress in character-centric understanding of narratives, we present `LiSCU`, a dataset of literature summaries paired with descriptions of characters that appear in them. We use `LiSCU` to propose two tasks that explore the ability of the modern neural models to understand the narrative from the perspective of characters. Performing human assessments on the model outputs show that there is still a lot of room for improvement on these tasks.

## Acknowledgments

## Broader Impacts and Ethics Statement

**Bias in Narrative Texts:** `LiSCU` is based on novels which often reflect societal norms and biases of their times. Such a dataset can be used to understand societal bias as well as design Natural Language Understanding models that can be more

aware of and possibly even avoid such biases. With this motivation, we analyzed the issue of gender bias in `LiSCU`.

First, we inferred the gender of the characters in our dataset using the pronouns used to refer to them. We could not infer the gender of some of the characters because of errors in the coreference system or lack of enough mentions, and we filtered them out for this analysis. We found that there are significantly more male characters than female characters in our dataset. Specifically, 66% of the characters are male. This suggests that systems that do not account for this bias might end up having more training data (and hence yield better performance) on descriptions of male characters than of female characters.

Second, we also investigated the scope of gender bias in the summaries. We computed the average number of mentions of male and female characters (in the summaries). We found that on average male and female characters are mentioned 32.1 and 31.7 times, respectively. This indicates that even though there are fewer female characters in the literary pieces of our dataset, the ones that are present play a significant role in the development of the narrative. Possibly because of their importance in the narrative, they are mentioned as many times as male characters in the summary (which describes the main developments and not all details from the literary piece).

Third, we investigated if the literary experts who composed the descriptions were biased in their analysis. For this, we compute the length of character descriptions of various characters. We found that there is no significant difference between male and female characters in this aspect. Specifically, the average number of tokens in the description of a male character was 203, and that of a female character was 200. Also, the average number of sentences in the description of a male character was 9.4 and that of a female character was 9.3. This also aligns with our observation in the previous experiment where we found that female characters, though fewer, play important roles in the narrative, and so their descriptions are not any shorter than descriptions of male characters. Overall, this analysis suggests that descriptions are not biased in their treatment of male and female characters.

In any language generation setting, such as ours, there is the possibility of (potentially harmful) social biases that can be introduced in the training data. As we did not specifically control or regularize our model to remove the possibility of such biases, we would urge downstream users to undertake the necessary quality-assurance testing to evaluate the extent to which such biases might be present and impacting their trained system and to make modifications to their model and procedures accordingly.

**Human participation in our study :** We conducted 2 human evaluations on Amazon Mechanical Turk. To ensure the annotators were fairly compensated, we did several rounds of test runs and estimated the average time to finish one HIT. Workers were paid $12/hr based on the HIT timings. We did not ask any personal, sensitive or identifying information from the annotators.

## References

Apoorv Agarwal, Sriramkumar Balasubramanian, Anup Kotalwar, Jiehan Zheng, and Owen Rambow. 2014. Frame semantic tree kernels for social network extraction from text. In *Proceedings of the 14th Conference of the European Chapter of the Association for Computational Linguistics*, pages 211–219, Gothenburg, Sweden. Association for Computational Linguistics.

David Bamman, Brendan O'Connor, and Noah A. Smith. 2013. Learning latent personas of film characters. In *Proceedings of the 51st Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 352–361, Sofia, Bulgaria. Association for Computational Linguistics.

David Bamman, Ted Underwood, and Noah A. Smith. 2014. A Bayesian mixed effects model of literary character. In *Proceedings of the 52nd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 370–379, Baltimore, Maryland. Association for Computational Linguistics.

Iz Beltagy, Matthew E. Peters, and Arman Cohan. 2020. Longformer: The long-document transformer. *ArXiv*, abs/2004.05150.

Faeze Brahman and Snigdha Chaturvedi. 2020. Modeling protagonist emotions for emotion-aware storytelling. In *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP)*, pages 5277–5294, Online. Association for Computational Linguistics.

Orson Scott Card. 1999. *Elements of Fiction Writing - Characters & Viewpoint*. Writer's Digest Books.

Nathanael Chambers. 2013. Event schema induction with a probabilistic entity-driven model. In *Proceedings of the 2013 Conference on Empirical Methods*

*in Natural Language Processing*, pages 1797–1807, Seattle, Washington, USA. Association for Computational Linguistics.

Nathanael Chambers and Dan Jurafsky. 2008. Unsupervised learning of narrative event chains. In *Proceedings of ACL-08: HLT*, pages 789–797, Columbus, Ohio. Association for Computational Linguistics.

Nathanael Chambers and Dan Jurafsky. 2009. Unsupervised learning of narrative schemas and their participants. In *Proceedings of the Joint Conference of the 47th Annual Meeting of the ACL and the 4th International Joint Conference on Natural Language Processing of the AFNLP*, pages 602–610, Suntec, Singapore. Association for Computational Linguistics.

Snigdha Chaturvedi, Mohit Iyyer, and Hal Daumé III. 2017. Unsupervised learning of evolving relationships between literary characters. In *Proceedings of the Thirty-First AAAI Conference on Artificial Intelligence,*, pages 3159–3165.

Eric Chu, Prashanth Vijayaraghavan, and Deb Roy. 2018. Learning personas from dialogue with attentive memory networks. In *Proceedings of the 2018 Conference on Empirical Methods in Natural Language Processing*, pages 2638–2646, Brussels, Belgium. Association for Computational Linguistics.

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

Micha Elsner. 2012. Character-based kernels for novelistic plot structure. In *Proceedings of the 13th Conference of the European Chapter of the Association for Computational Linguistics*, pages 634–644, Avignon, France. Association for Computational Linguistics.

David K. Elson, Nicholas Dames, and Kathleen R. McKeown. 2010. Extracting social networks from literary fiction. In *ACL 2010, Proceedings of the 48th Annual Meeting of the Association for Computational Linguistics, July 11-16, 2010, Uppsala, Sweden*, pages 138–147. The Association for Computer Linguistics.

Angela Fan, Thibaut Lavril, Edouard Grave, Armand Joulin, and Sainbayar Sukhbaatar. 2021. Addressing some limitations of transformers with feedback memory.

Alvan R Feinstein and Domenic V Cicchetti. 1990. High agreement but low kappa: I. the problems of two paradoxes. *Journal of clinical epidemiology*, 43(6):543–549.

Lucie Flekova and Iryna Gurevych. 2015. Personality profiling of fictional characters using sense-level links between lexical resources. In *Proceedings of the 2015 Conference on Empirical Methods in Natural Language Processing*, pages 1805–1816.

Amit Goyal, Ellen Riloff, and Hal Daumé III. 2010. Automatically producing plot unit representations for narrative text. In *Proceedings of the 2010 Conference on Empirical Methods in Natural Language Processing*, pages 77–86, Cambridge, MA. Association for Computational Linguistics.

Mohit Iyyer, Anupam Guha, Snigdha Chaturvedi, Jordan Boyd-Graber, and Hal Daumé III. 2016. Feuding families and former friends: Unsupervised learning for dynamic fictional relationships. In *Proceedings of the 2016 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*, pages 1534–1544.

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, Luke Zettlemoyer, and Daniel Weld. 2019b. BERT for coreference resolution: Baselines and analysis. In *Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing (EMNLP-IJCNLP)*, pages 5803–5808, Hong Kong, China. Association for Computational Linguistics.

X. J. Kennedy and D. Gioia. 1983. *Literature: An introduc- tion to fiction. Poetry, Drama, and writing.*

Evgeny Kim and Roman Klinger. 2019. Frowning frodo, wincing leia, and a seriously great friendship: Learning to classify emotional relationships of fictional characters. In *Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, NAACL-HLT 2019, Minneapolis, MN, USA, June 2-7, 2019, Volume 1 (Long and Short Papers)*, pages 647–653. Association for Computational Linguistics.

Vinodh Krishnan and Jacob Eisenstein. 2015. "you're mr. lebowski, I'm the dude": Inducing address term formality in signed social networks. In *Proceedings of the 2015 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*, pages 1616–1626, Denver, Colorado. Association for Computational Linguistics.

Wojciech Kryscinski, Nitish Shirish Keskar, Bryan McCann, Caiming Xiong, and Richard Socher. 2019. Neural text summarization: A critical evaluation. In *Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing (EMNLP-IJCNLP)*, pages 540–551, Hong Kong, China. Association for Computational Linguistics.

Faisal Ladhak, Bryan Li, Yaser Al-Onaizan, and Kathleen McKeown. 2020. Exploring content selection

in summarization of novel chapters. In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, pages 5043–5054, Online. Association for Computational Linguistics.

Zhenzhong Lan, Mingda Chen, Sebastian Goodman, Kevin Gimpel, Piyush Sharma, and Radu Soricut. 2020. Albert: A lite bert for self-supervised learning of language representations. In *International Conference on Learning Representations*.

J Richard Landis and Gary G Koch. 1977. The measurement of observer agreement for categorical data. *biometrics*.

Mike Lewis, Yinhan Liu, Naman Goyal, Marjan Ghazvininejad, Abdelrahman Mohamed, Omer Levy, Ves Stoyanov, and Luke Zettlemoyer. 2019. Bart: Denoising sequence-to-sequence pre-training for natural language generation, translation, and comprehension. *arXiv*.

Chin-Yew Lin. 2004. ROUGE: A package for automatic evaluation of summaries. In *Text Summarization Branches Out*, pages 74–81, Barcelona, Spain. Association for Computational Linguistics.

Yinhan Liu, Myle Ott, Naman Goyal, Jingfei Du, Mandar Joshi, Danqi Chen, Omer Levy, Mike Lewis, Luke Zettlemoyer, and Veselin Stoyanov. 2019. Roberta: A robustly optimized BERT pretraining approach. *CoRR*, abs/1907.11692.

Mehdi Manshadi, Reid Swanson, and Andrew S. Gordon. 2008. Learning a probabilistic model of event sequences from internet weblog stories. In *Proceedings of the Twenty-First International Florida Artificial Intelligence Research Society Conference, May 15-17, 2008, Coconut Grove, Florida, USA*, pages 159–164. AAAI Press.

Joshua Maynez, Shashi Narayan, Bernd Bohnet, and Ryan McDonald. 2020. On faithfulness and factuality in abstractive summarization. In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, pages 1906–1919, Online. Association for Computational Linguistics.

Neil McIntyre and Mirella Lapata. 2010. Plot induction and evolutionary search for story generation. In *Proceedings of the 48th Annual Meeting of the Association for Computational Linguistics*, pages 1562–1572, Uppsala, Sweden. Association for Computational Linguistics.

Rada Mihalcea and Hakan Ceylan. 2007. Explorations in automatic book summarization. In *Proceedings of the 2007 Joint Conference on Empirical Methods in Natural Language Processing and Computational Natural Language Learning (EMNLP-CoNLL)*, pages 380–389, Prague, Czech Republic. Association for Computational Linguistics.

Rada Mihalcea and Paul Tarau. 2004. TextRank: Bringing order into text. In *Proceedings of the 2004 Conference on Empirical Methods in Natural Language Processing*, pages 404–411, Barcelona, Spain. Association for Computational Linguistics.

Jekaterina Novikova, Ondřej Dušek, Amanda Cercas Curry, and Verena Rieser. 2017. Why we need new evaluation metrics for NLG. In *Proceedings of the 2017 Conference on Empirical Methods in Natural Language Processing*, pages 2241–2252, Copenhagen, Denmark. Association for Computational Linguistics.

Kishore Papineni, Salim Roukos, Todd Ward, and Wei-Jing Zhu. 2002. Bleu: a method for automatic evaluation of machine translation. In *Proceedings of the 40th Annual Meeting of the Association for Computational Linguistics*, pages 311–318, Philadelphia, Pennsylvania, USA. Association for Computational Linguistics.

Dragomir R. Radev. 2001. Experiments in single and multidocument summarization using mead. In *In First Document Understanding Conference*.

Alec Radford, Jeff Wu, Rewon Child, David Luan, Dario Amodei, and Ilya Sutskever. 2019. Language models are unsupervised multitask learners. *OpenAI Blog*, 1:8.

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

Colin Raffel, Noam Shazeer, Adam Roberts, Katherine Lee, Sharan Narang, Michael Matena, Yanqi Zhou, Wei Li, and Peter J. Liu. 2020b. Exploring the limits of transfer learning with a unified text-to-text transformer. *Journal of Machine Learning Research*, 21(140):1–67.

Michaela Regneri, Alexander Koller, and Manfred Pinkal. 2010. Learning script knowledge with web experiments. In *Proceedings of the 48th Annual Meeting of the Association for Computational Linguistics*, pages 979–988, Uppsala, Sweden. Association for Computational Linguistics.

Aurko Roy, Mohammad Saffar, Ashish Vaswani, and David Grangier. 2021. Efficient content-based sparse attention with routing transformers. *Trans. Assoc. Comput. Linguistics*, 9:53–68.

R. Schank and R. Abelson. 1977. *Scripts, plans, goals and understanding: An inquiry into human knowledge structures*. Lawrence Erlbaum Associates, Hillsdale, NJ.

Shashank Srivastava, Snigdha Chaturvedi, and Tom M. Mitchell. 2016. Inferring interpersonal relations in narrative summaries. In *Proceedings of the Thirtieth AAAI Conference on Artificial Intelligence*, pages 2807–2813.

Hardik Vala, David Jurgens, Andrew Piper, and Derek Ruths. 2015. Mr. bennet, his coachman, and the archbishop walk into a bar but only one of them gets recognized: On the difficulty of detecting characters in literary texts. In *Proceedings of the 2015 Conference on Empirical Methods in Natural Language Processing*, pages 769–774.

Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Ł ukasz Kaiser, and Illia Polosukhin. 2017. Attention is all you need. In *Advances in Neural Information Processing Systems*, volume 30. Curran Associates, Inc.

Xin Wang, Wenhu Chen, Yuan-Fang Wang, and William Yang Wang. 2018. No metrics are perfect: Adversarial reward learning for visual storytelling. In *Proceedings of the 56th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 899–909, Melbourne, Australia. Association for Computational Linguistics.

Thomas Wolf, Lysandre Debut, Victor Sanh, Julien Chaumond, Clement Delangue, Anthony Moi, Pierric Cistac, Tim Rault, R'emi Louf, Morgan Funtowicz, and Jamie Brew. 2019. Huggingface's transformers: State-of-the-art natural language processing. *arXiv*.

Tianyi Zhang, Varsha Kishore, Felix Wu, Kilian Q. Weinberger, and Yoav Artzi. 2020. Bertscore: Evaluating text generation with bert. In *International Conference on Learning Representations*.

Weiwei Zhang, Jackie Chi Kit Cheung, and Joel Oren. 2019. Generating character descriptions for automatic summarization of fiction. In *The Thirty-Third AAAI Conference on Artificial Intelligence, AAAI 2019*, pages 7476–7483. AAAI Press.

## A    Collecting Annotations from Crowd Workers

To alleviate the limitations of crowd-sourcing and ensure high quality of annotations, we took several steps. First, we conducted a pilot annotation exercise where we (authors) assessed the feasibility of the proposed task on a subset (250 instances) of the data. This pilot annotation helped us set up the task on AMT in a way that would make the task feasible for turkers (e.g. by asking clear concise questions). Second, we designed our setup to avoid annotator fatigue by asking them to read the summary context once and answer questions about all characters in that summary. Third, we ran a few experiments on AMT (before annotating the entire set) where we also included a 'comment' section for turkers to allow them to bring up issues or ambiguities in our setup. We then manually analyzed the results and modified the tasks based on the comments. Finally, after annotating the entire set, we computed inter-annotator agreement as a way to ensure trust in the annotation quality. We found reasonable agreements between annotators as reported in Footnote 10 of the paper. We would also like to mention that we received several comments from the annotators that they found the task very interesting and enjoyable.

| Model | BLEU | ROUGE-1 | ROUGE-2 | ROUGE-L | BERT-F1 |
|---|---|---|---|---|---|
| **Length Truncated Input** | | | | | |
| GPT2-L | 0.90/0.71/0.60/0.59 | 20.33/19.94/18.81/19.39 | 4.34/3.60/3.36/3.32 | 18.50/17.91/17.14/17.84 | 76.23/80.11/76.53/75.81 |
| Longformer | 1.86/1.06/0.92/0.70 | 24.16/21.78/20.55/20.20 | 6.80/4.62/4.33/3.98 | 22.05/19.65/18.71/18.26 | 85.60/84.92/84.46/84.24 |
| **Coreference Truncated Input** | | | | | |
| GPT2-L | 0.82/0.63/0.53/0.58 | 19.80/19.24/17.96/18.49 | 3.86/3.24/3.06/3.04 | 17.79/17.39/16.46/16.62 | 76.16/79.52/78.33/80.23 |
| Longformer | 1.78/1.09/0.77/0.65 | 23.32/22.23/20.23/19.90 | 6.04/4.80/3.96/3.57 | 21.41/20.22/18.16/17.70 | 85.43/85.07/84.47/84.12 |
| **Full Length Input** | | | | | |
| Longformer | 2.15/1.31/1.04/0.65 | 24.47/22.56/20.90/20.63 | 6.98/5.37/4.63/3.84 | 21.91/20.40/18.85/18.48 | 85.66/85.11/84.57/84.35 |

Table 8: Breakdown results on subsets of test set with annotated fact coverage as all/most/some/little.

| Generated Description | Comments |
|---|---|
| **Book title:** The Three Sisters<br>**Character name:** Vershinin<br>**Summary:** Three Sisters mainly follows the story of–wait for it–three sisters: Olga, Masha, and Irina Prozorov. They live with their brother, Andrey, in a big house on the edge of a small Russian town. The townspeople are kinda backward and boring compared to their educated and culture-lovin' family, so this set of sibs is not too fond of the town to begin with. Believe it or not, the only halfway interesting people around are the guys in the military. Basically, the Prozorov kids are worldly, well-educated army brats. And being in the army in Tsarist Russia pretty much meant you were in with the aristocracy and, once you got through the fighting stuff, probably developed a taste for the finer things in life. So ever since the family moved from Moscow eleven years prior (with their father, now dead), the sisters have obsessed over the dream of moving back to the big city. But guess what? It's not happening. Olga, the eldest, is a spinster schoolteacher and eventually becomes a headmistress living with her elderly maid. Masha, the middle sister, is married to another schoolteacher, Kulygin, whom she despises for his small-mindedness. She has an affair with the officer Vershinin because he's given to just the sort of philosophizing that really starts her engines. And guess what happens: the affair ends in heartbreak. Poor Irina, the youngest, has fanciful ideas about the value of work, but soon realizes that, in reality, work sucks the life out of her. She's also in love with the idea of love, but doesn't get to experience it. Finally she comes around to saying "yes" to Baron Tuzenbach, a friendly but ugly man who's been after her for years. On the day they're leaving to get married, he gets shot in a duel. Bummer. Andrey, the brother, gives up his intellectual dreams to pursue a town girl, Natasha. They marry, have kids, and little by little she takes over the estate. At the end of the play, the upper-class Prozorovs are pretty much evicted from their own house, while Natasha, a symbol of the working class, is on the rise. Allegory, much? For anyone not already reaching for the history books, this was just a few years before the Bolshevik Revolution of 1917–and you better believe smart Ruskies like Chekhov could already sense some storms on the horizon.<br>**Generated Character Description:** Vershinin is one of the most interesting characters in Three Sisters. He's a young officer in the Russian army, and he's got a lot of time on his hands. But that doesn't stop him from having an affair with Irina, the youngest sister, who falls in love with him because she's attracted to his philosophical ideas about the meaning of life. | According to the summary "Masha", the middle sister, is the one who has an affair with the officer Vershinin, and not the youngest sister, Irina. |

Table 9: Qualitative example 1 for the generated descriptions. Words in red correspond to hallucinated or missing content, words in green correspond to faithful information.

| Generated Description | Comments |
|---|---|
| **Book title:** The Sun Also Rises<br><br>**Character name:** Pedro Romero<br><br>**Summary**: Jake Barnes and his expatriate friends live in the topsy-turvy, hedonistic (sensual and self-indulgent) world of post-World War I Paris. There, they occasionally work, but spend most of their time partying, drinking, and arguing. From Jake's perspective, we meet the cast of characters that populates his story: the most important among them are Robert Cohn, a weak-willed, down-on-his-luck Princeton grad and unsuccessful writer, and Lady Brett Ashley, an exciting, beautiful, and unpredictable British divorcee. Although Jake and Brett are actually in love, they aren't together, presumably because a mysterious war wound has rendered Jake impotent. Cohn falls in love with Brett (as everyone does) and, despite the fact that she's not terribly impressed with him, she secretly goes on a trip with him to the Spanish resort town of San Sebastian. Cohn is infatuated with Brett–he's completely smitten. We're talking truly, madly, deeply in smit. Unfortunately for Cohn (and for everyone, for that matter), Brett is engaged to a wealthy, charming, and utterly inept drunkard named Mike. Jake's whimsical friend Bill returns to Paris from a trip and a plan is born: everyone agrees to decamp to Spain for some fishing and the running of the bulls in Pamplona. On their brief fishing trip, Bill and Jake have a splendid time communing with nature and with each other, but the relaxation quickly comes to an end. They return to civilization and meet up with Brett, Mike, and Cohn in Pamplona for a weeklong orgy of bullfights, alcohol, and high drama. Jake has a true passion (aficion) for bullfighting, but everyone else is simply there to have a good time. Brett begins a scandalous affair with a passionate and talented young bull-fighter, Pedro Romero. Jake feels terrible for many reasons– among them is the fear that he has corrupted Romero in some way by introducing him to Brett. Cohn's thwarted infatuation with Brett leads to arguments with everyone and, finally, he beats the unfortunate Romero to a bloody pulp. As the fiesta winds down, everyone leaves Pamplona in various states of anxiety, depression and frustration.Jake heads to San Sebastian, where he intends to decompress alone for a while. Unfortunately, desperate telegrams from Brett arrive immediately. He goes to her in Madrid, where she is alone, having sent Romero away. For the first time, we see Brett truly vulnerable, afraid, and guilty. The future looks just as bleak–Jake and Brett agree again that, even though they love each other, they can't be together.<br><br>**Generated Character Description:** Pedro Romero is <span style="color:green">a young Spanish bull-fighter who is in love with Lady Brett Ashley, a beautiful British divorcee.</span> He's a talented and <span style="color:green">talented young man</span>, and, despite the fact that he is only <span style="color:red">twenty-five years old</span>, he has <span style="color:red">the potential to be one of the best bullfighters in Pamplona.</span> <span style="color:green">Jake feels terrible for introducing him to Brett, fearing that it has corrupted him.</span> | The description captures most of the important details about the character. |

Table 10: Qualitative example 2 for the generated descriptions. Words in <span style="color:red">red</span> correspond to hallucinated or missing content, and words in <span style="color:green">green</span> correspond to faithful information.

| Generated Description | Comments |
|---|---|
| **Book title:** The Waves<br>**Character name:** Neville<br><br>**Summary:** The story begins by introducing us to the novel's six (yup, you read that right) narrators, Bernard, Neville, Louis, Jinny, Susan, and Rhoda, who meet as children in a nursery. During this phase of the novel, we learn a lot about the characters' personalities and their relationships to each other. After looking on as our new friends get embroiled in some kid-level dramas (e.g., trouble in math class and unrequited crushes), the six children head off to their respective boarding schools. At that time, the boys meet Percival, whom everyone seems to revere (and Neville falls in love with). The protagonists then all graduate and proceed into their adult careers (with a stop at university along the way, for some). At some point in there, Percival becomes friends with the girls as well, though we're not sure when that actually occurs. The narrators' paths diverge quite a bit as the novel progresses. After enduring a stint in a Swiss school, Susan returns to her beloved hometown, gets married, and starts having babies. Meanwhile, Bernard apparently moves to Waterloo (that's not entirely clear, but Woolf drops some clues to that effect), and we're not entirely sure what he does there, other than shave and make up stories about pedestrians. Meanwhile, Rhoda, Louis, Neville, and Jinny go to live in London, and their life paths are all over the map: Louis works for a shipping company, Jinny is a socialite, and Neville is a classics professor (we don't learn Rhoda's profession). Midway through the book, the friends meet up for dinner in London to see Percival off before he leaves to work in India, at which point Bernard announces that he is getting married. An unspecified amount of time later, Percival falls from his horse in India and dies, and our narrators are devastated. The death occurs just as Bernard's son is born, creating some serious cognitive dissonance for poor Bernard. The characters then struggle with aging and reflect upon the progress of their lives and adulthood. Toward the end, they all meet up at Hampton Court and experience the aches and pains that come when old friends reconnect after a long time has passed. The novel ends with Bernard talking to an apparent stranger, attempting to sum up the lives of the six narrators (i.e., the events of the novel) and work out his philosophies of language and lifeand that's why the last chapter is fifty pages long. He reflects upon his lifelong struggle to turn his phrasemaking into something meaningful and, despite the major ups and downs he's experienced along the way, he resolves to keep on trying.<br><br>**Generated Character Description:** Neville is one of the six narrators. <span style="color:green">He's a classics professor</span> at <span style="color:red">Oxford University</span> and, <span style="color:red">like the other five</span>, he's also a classicist. <u>We don't learn a whole lot about</u> Neville, other than the fact that <span style="color:green">he falls in love with Percival</span>, which makes him a bit of an oddball. In fact, we're not even sure what he does in his spare time away from Oxford. | The model hallucinated some details such as being a professor at Oxford University while the summary only mentions: He lives in London and is a classics professor. Also, phrases such as "We don't learn a whole lot about" appears quite frequently in the generated descriptions. |

Table 11: Qualitative example 3 for the generated descriptions. Words in <span style="color:red">red</span> correspond to hallucinated or missing content, words in <span style="color:green">green</span> correspond to faithful information, and <u>underline</u> corresponds to generic repetitive content.

| Generated Description | Comments |
|---|---|
| **Book title:** Travels with Charley | |
| **Character name:** Charley the Dog | |
| **Summary:** Because he's feeling pretty out of touch with his own country–and he's considered a great American author and all that–John Steinbeck decides to take a road trip around the U.S. to check it out and get a sense of where Americans and their hometowns are at in 1960. To get all prepped, he commissions a souped-up truck with a little house on the back that he can live in when he isn't crashing at hotels. He calls the truck "Rocinante" after Don Quixote's horse–clever, huh? When he's all set (and after a small run-in with a hurricane just before he was supposed to leave), he and Charley (his French poodle) hit the road. He starts out by driving over into Connecticut from his home in Long Island (with some assists from ferries, natch) and then heads north into New England. Along the way, he meets a pretty colorful group of characters and learns about their ways of life and their perspectives on the country and its politics. Also, he kind of takes the temperature of regional "temperaments" along the way. Then he comes back down out of New England and heads west, crossing through New York. He tries to cut through Canada, but he gets into a kerfuffle at the border because Charley doesn't have his proof of rabies vaccination, so he has to turn around. Steinbeck then passes through the Midwest, continuing to offer his reflections and thoughts about the people and places he encounters along the way. When he gets to Chicago, he puts Charley in a kennel and enjoys a couple of days with his wife, who flew out to meet him. He doesn't give us details of their time together, though. After that brief interlude, he heads further west into Minnesota and Wisconsin. He hits bad traffic and gets lost around the Twin Cities, and he's charmed by Wisconsin and its dells. He also visits Sauk Centre, the birthplace of author Sinclair Lewis. Then he heads toward Fargo, North Dakota, which apparently had been the subject of his boyhood fantasies. We picture Hawaii when we're fantasizing about faraway places, but okay... He heads through North Dakota and the Bad Lands, warming up to that area quite a bit when he sees it at sunset (otherwise, it's pretty scary, apparently). When he goes through Montana, he loves it. However, he soon has to head back south so he can see Yellowstone... though he doesn't stay long, since Charley freaks out about the bears. Then, he heads to Washington. Along the way, Charley ends up getting sick with a bladder problem, and Steinbeck tries to get him some help from a vet in Spokane. The doc isn't actually super-helpful, though. Still, they get the problem resolved enough that they can continue driving. Heading to Seattle, Steinbeck finds the place changed a lot from the last time he was there (and not for the better, sadly). Then he heads into Oregon, where his tire blows out. Through the kindness of a dude at the one service station open in that area on a Sunday, though, he gets some new heavy-duty tires and is on his way pretty quickly. Then he heads south through California, visiting Salinas (his hometown) on the way. While there, he argues with his sisters about politics and reconnects with old friends, but apparently it's kind of awkward, and he's reminded of that old adage that you can't go home again. After that, he starts circling back toward home, driving through the Mojave Desert and through the Southwest and Texas. While he's in Texas, he meets up with his wife once again and stays with some friends (oh, and actually takes Charley to a competent vet in Amarillo). We get a lot of detail about Steinbeck's impressions of Texas. After that little break, Steinbeck goes to Louisiana so he can see for himself what's going on there in terms of racial tensions (which appear to be pretty high). There was a well-known protest going on: two African American children had been admitted to a New Orleans school, and lots of people had a huge problem with that fact. So, some women known as the Cheerleaders basically stood outside every morning to scream at these poor children as they headed into school. Steinbeck went to get a firsthand look at this spectacle and was, as you might imagine, disgusted. He then headed back home, finally deciding he was "done" with the journey somewhere in Virginia. | The description never mention anything about the Charley's health struggles as appeared in the summary we have: "Along the way, Charley ends up getting sick with a bladder problem, and Steinbeck tries to get him some help from a vet in Spokane. The doc isn't actually super-helpful, though." |
| **Generated Character Description:** When we first meet Charley, he's a French poodle. He's Steinbeck's pet and his companion for most of the road trip, and we get to see him grow up a bit. In fact, we see a lot of him in the book, from the time he first meets him until he gets to Chicago. We also get a glimpse of his personality as a character in his own right, when he tells us, "He was a good-natured, loving dog" (1.2.1). | |

Table 12: Qualitative example 4 for the generated descriptions. words in green correspond to faithful information, and underline corresponds to generic repetitive content.

**Instructions** (Click to collapse)

## Task Description

In this page, you will be given a book summary along with short descriptions about 4 characters of the book. For each character, you need to do the following 2 tasks:

- Task 1 - In the short character description, all names and mentions of the character (see Notes for more details) are anonymized (replaced with blanks). You need to read the summary and the anonymized character description, then **identify which character the description is about** by choosing from a list of 5 characters (1 extra character choice).

- Task 2 - You will see the full version of the character description this time (character's name and mentions are not anonymized). Read the character description again and **evaluate the quality of the character description** by answering the following 2 questions:
  - How much of the information about the specific character in the corresponding character description is present in the book summary **either explicitly or implicitly**? Note that "implicit" information can be inferred from the book summary but not directly stated. For example, *"Candy adores Mathu because he basically raised her."* in character description is implicitly mentioned as *"Mathu is virtually her foster father."* in the book summary.
  - Given the book summary, how easy is it to write the character description? If in the previous question you found that some of the information in the character description was not present in the summary, please disregard that while answering this question. In other words, while answering this question please <u>only</u> consider the information in the character description which is <u>explicitly or implicitly mentioned in the summary</u>.

For each character, <u>Task 2 will appear after you submit the answer for Task 1.</u>

### Notes:

- Make sure you read **BOTH** the book sumamry and character descriptions carefully, since it will increase the accuracy of the response.
- Character's name and mentions include the following cases:
  - Full Name of Character(s)
  - First Name of Character(s)
  - Nickname/Alias of Character(s), *e.g., Captain America*
  - Pronouns, *e.g., He, She*
  - Possessive Pronouns, *e.g., His, Her*
  - Noun Phrases that show the relationship of the character(s) with others, *e.g., Jeff's Wife*

### Summary Example

> Three Weeks with My Brother is two stories in one. On the surface, it tells of a trip around the world that Nicholas Sparks takes with his brother Micah. Three Weeks with My Brother begins on the day that Nicholas Sparks receives the flier in the mail and ends with the brothers returning home. In between, Nicholas Sparks recounts the sights, sounds, and spectacles of various countries and continents, taking the reader on the journey with him. But Three Weeks with My Brother is more than just a travelogue. The text is the memoir of a successful author who seemingly is living the American Dream. Yet unbeknownst to most of his readers, he has also lived the American Tragedy. The story follows two brothers and their journey to becoming the best husbands, fathers, sons, brothers, and friends that they can be. Three Weeks with My Brother is a no-holds-barred memoir that shares the good, the bad, and the ugly that made Nicholas Sparks the man he is today.

### Task 1 Example

**Character Description (anonymized)**:

> Co-author and narrative voice of the memoir. Three Weeks with My Brother is both the story of a trip _____ takes with _____ brother, Micah, as well as the story of his immediate family as _____ grows up, marries, and has children of _____ own.

Based on the summary and the character description, **answer the following question**:

> Which character among the following does the description refer to?
> A. ○ Micah Sparks
> B. ○ Dana Sparks
> C. ◉ Nicholas Sparks
> D. ○ Bob
> E. ○ unable to identify character

### Task 2 Example

**Character Description (full version)**:

> Co-author and narrative voice of the memoir. Three Weeks with My Brother is both the story of a trip he takes with his brother, Micah, as well as the story of his immediate family as he grows up, marries, and has children of his own.

Based on the summary and the character description, **answer the following question**:

> How much of the information about the specific character in the corresponding character description is present in the summary (either explicitly or implicitly)?
> A. ◉ almost all information
> B. ○ most of the information
> C. ○ some of the information
> D. ○ little or none
> E. ○ character does not appear in the summary at all
>
> **Explanation**: The correct answer is "almost all information" since in the given character description, it is mentioned "the story of his immediate family as he grows up, <u>marries</u>, and has <u>children of his own</u>" and same information is implicitly present in the book summary when it says: "The story follows two brothers and their journey to <u>becoming the best husbands, fathers</u>, sons, brothers, and etc.
>
> Given the summary, how easy is it to write the character description? If in the previous question you found that some of the information in the character description was not present in the summary, please disregard that while answering this question. In other words, please only consider the information in the character description which is explicitly or implicitly mentioned in the summary.
>
> ○ Too Difficult  ○ Somewhat Difficult  ○ Medium  ○ Somewhat Easy  ◉ Too Easy

Figure 6: An illustration of human assessment on AMT.

**Instructions** (Click to collapse)

## Task Description

In this page, you will be given a book summary along with short descriptions about 4 characters of the book. For each character, you need to answer 6 questions that help us to measure how good these descriptions are. The first 4 questions ask you to measure the following 4 quality metrics **on a scale of 1 to 5 with 1 being the lowest quality and 5 being the highest quality**:

1. Grammatical Correctness (Q1) - A grammatical character description will mostly follow correct English grammar
2. Logical Correctness (Q2) - A logically coherent character description might be grammatically incorrect but will be able to convey the intended meaning as a whole and will make sense
3. Faithfulness (Q3) - A faithful character description will not mention facts, which are irrelevant to the character and/or not stated in the summary
4. Coverage of Key Facts (Q4a) - Whether the character description captures all the key facts about the character in the summary.

Question 4b asks you to provide what details/facts mentioned in the summary are missing from the descriptions. **Please select all that apply.**

Question 5 asks you to measure the overall quality of the descriptions by considering all 4 quality metrics describe above **on a scale of 1 to 5 with 1 being the lowest quality and 5 being the highest quality**.

### Notes:

- Make sure you read **BOTH** the sumamry and descriptions carefully, since it will increase the accuracy of the response.

### Example

**Summary**:

Three Weeks with My Brother is two stories in one. On the surface, it tells of a trip around the world that Nicholas Sparks takes with his brother Micah. Three Weeks with My Brother begins on the day that Nicholas Sparks receives the flier in the mail and ends with the brothers returning home. In between, Nicholas Sparks recounts the sights, sounds, and spectacles of various countries and continents, taking the reader on the journey with him. But Three Weeks with My Brother is more than just a travelogue. The text is the memoir of a successful author who seemingly is living the American Dream. Yet unbeknownst to most of his readers, he has also lived the American Tragedy. The story follows two brothers and their journey to becoming the best husbands, fathers, sons, brothers, and friends that they can be. Three Weeks with My Brother is a no-holds-barred memoir that shares the good, the bad, and the ugly that made Nicholas Sparks the man he is today.

**Description**:

**Nicholas Sparks**
Co-author and narrative voice of the memoir. Three Weeks with My Brother is both the story of a trip he takes with his brother, Micah, as well as the story of his immediate family as he grows up, marries, and has children of his own.

Based on the summary and the description, **answer the following question**:

**Q1: Is the character description grammatical?**
Please answer it on a scale of 1 to 5 with 1 being "not grammatical at all" and 5 being "very grammatical". A grammatical character description will mostly follow correct English grammar. Note that text can be grammatically correct but nonsensical. For example, "Where did the spoon take off?" should be rated as 5 for this question because it is grammatically correct.

1　2　3　4　**5** ⦿

**Q2: Is the character description logically meaningful and coherent?**
Please answer it on a scale of 1 to 5 with 1 being "not logically meaningful or coherent at all" and 5 being "very logically meaningful and coherent". A logically coherent character description might be grammatically incorrect but will be able to convey the intended meaning as a whole and will make sense. For example, "John and Mary goes to the market." should be rated as 5 for this question because it conveys the intended meaning.

1　2　3　4　**5** ⦿

**Q3: Is the character description faithful to the given summary?**
Please answer it on a scale of 1 to 5 with 1 being "not faithful at all" and 5 being "very faithful". A faithful character description will not mention facts, which are irrelevant to the character and/or not stated in the summary.

1　2　3　**4** ⦿　5

**Q4a: Does the character description capture all the key facts about this character in the summary?**
Please answer it on a scale of 1 to 5 with 1 being "none of the key facts" and 5 being "all of the key facts". These could include (but not restricted to) main events that the character is involved in, the character's role in the narrative, and his/her relationship with other key characters.

1　2　3　**4** ⦿　5

**Q4b: What details about the character does the above character description miss or describe incorrectly? (Select all that apply)**
Note that the description is supposed to describe only the important details and not necessarily all of them.

☐ The character description misses/misrepresents some main event(s) that the character is involved in

☐ The character's role in the narrative (e.g., protagonist, antagonist, etc.) is important but is not included or misrepresented in the character description

☐ The character's relationship with other characters in the narrative is important (e.g., the protagonist's wife) but is not included or misrepresented in the character description

☐ The character's personal characteristics (e.g. age, ethnicity, size, personality type, etc.) are important for the narrative but are not included or misrepresented in the character description

☐ The character's motivation, desires, needs, and behavior are important but are not included or misrepresented in the character description

☑ None of the above. The character description captures most of the important details about the character

☐ Other key points that are missing or misrepresented but not listed above

**Q5: Considering all four criteria listed above (question 1 - 4) how would you rate the overall quality of this character's description?**
Please answer it on a scale of 1 to 5 with 1 being "not good at all" and 5 being "very good".

1　2　3　**4** ⦿　5

Figure 7: An illustration of human evaluation for generated character description.