

	English	avg. len.	vocab size	density	Chinese	avg. len.	vocab size	density
human	<b>All</b>	142.50	<b>79157</b>	<b>2.33</b>	<b>All</b>	102.27	<b>75483</b>	<b>5.75</b>
ChatGPT		<b>198.14</b>	66622	1.41		<b>115.3</b>	45168	3.05
human	<i>reddit_eli5</i>	134.21	<b>55098</b>	<b>2.46</b>	<i>nlpcc_dbqa</i>	24.44	10621	<b>25.43</b>
ChatGPT		<b>194.84</b>	44926	1.38		<b>78.21</b>	<b>11971</b>	8.96
human	<i>open_qa</i>	35.09	9606	<b>23.06</b>	<i>open_qa</i>	93.68	<b>40328</b>	<b>13.13</b>
ChatGPT		<b>131.68</b>	<b>16251</b>	10.40		<b>150.66</b>	26451	5.35
human	<i>wiki_csai</i>	<b>229.34</b>	<b>15859</b>	<b>8.21</b>	<i>baike</i>	<b>112.25</b>	<b>28966</b>	<b>5.59</b>
ChatGPT		208.33	9741	5.55		77.19	14041	3.94
human	<i>medicine</i>	92.98	<b>11847</b>	<b>10.42</b>	<i>medicine</i>	92.34	<b>9855</b>	<b>9.94</b>
ChatGPT		<b>209.61</b>	7694	3.00		<b>165.41</b>	7211	4.06
human	<i>finance</i>	202.07	<b>25500</b>	<b>3.21</b>	<i>finance</i>	80.76	2759	<b>5.05</b>
ChatGPT		<b>226.01</b>	21411	2.41		<b>120.84</b>	<b>4043</b>	4.94
human	-	-	-	-	<i>psychology</i>	<b>254.82</b>	<b>16160</b>	<b>5.77</b>
ChatGPT		-	-	-		164.53	5897	3.26
human	-	-	-	-	<i>law</i>	28.77	2093	<b>19.55</b>
ChatGPT		-	-	-		<b>143.76</b>	<b>3857</b>	7.21

Table 3: Average answer length, vocabulary size and density comparisons on our corpus.

## 4 Linguistic Analysis

In this section, we analyze the linguistic features of both humans’ and ChatGPT’s answers, and try to find some statistical evidence for the characteristics concluded in Section 3.

### 4.1 Vocabulary Features

In this part, we analyze the vocabulary features of our collected corpus. We are interested in how humans and ChatGPT differ in the choice of words when answering the same set of questions.

Since the number of human/ChatGPT answers is unbalanced, we randomly sample one answer from humans and one answer from ChatGPT during our statistical process. We calculated the following features: **average length** ( $L$ ), which is the average number of words in each question; **vocab size** ( $V$ ), the number of unique words used in all answers; we also propose another feature called **density** ( $D$ ), which is calculated by  $D = 100 \times V / (L \times N)$  where  $N$  is the number of answers. Density measures how *crowded* different words are used in the text. For example, if we write some articles that add up to 1000 words, but only 100 different words are used, then the density is  $100 \times 100 / 1000 = 10$ . The higher the density is, the more different words are used in the same length of text.

In Table 3, we report the vocabulary features for both English and Chinese corpus. Looking at both features of *average length* and *vocab size*, we can see that: **compared to ChatGPT, human answers are relatively shorter, but a larger vocabulary is used.** This phenomenon is particularly obvious in the Chinese *open\_qa* split and the *medical* splits in both languages, where the average length of ChatGPT is nearly twice longer than that of humans, but the vocab size is significantly smaller.

This phenomenon is also reflected by the *density* factor. The word density of humans is greater than ChatGPT’s in **every split**, which further reveals that **humans use a more diverse vocabulary in their expressions.**

### 4.2 Part-of-Speech & Dependency Analysis

In this part, we compare the occurrences of different part-of-speech (POS) tags and the characteristics of the dependency relations.

#### 4.2.1 Part-of-Speech

Figure 1 illustrates the comparisons between humans and ChatGPT in terms of POS usage. In HC3-English, ChatGPT uses **more** NOUN, VERB, DET, ADJ, AUX, CONJ and PART words, while using less ADV and PUNCT words.

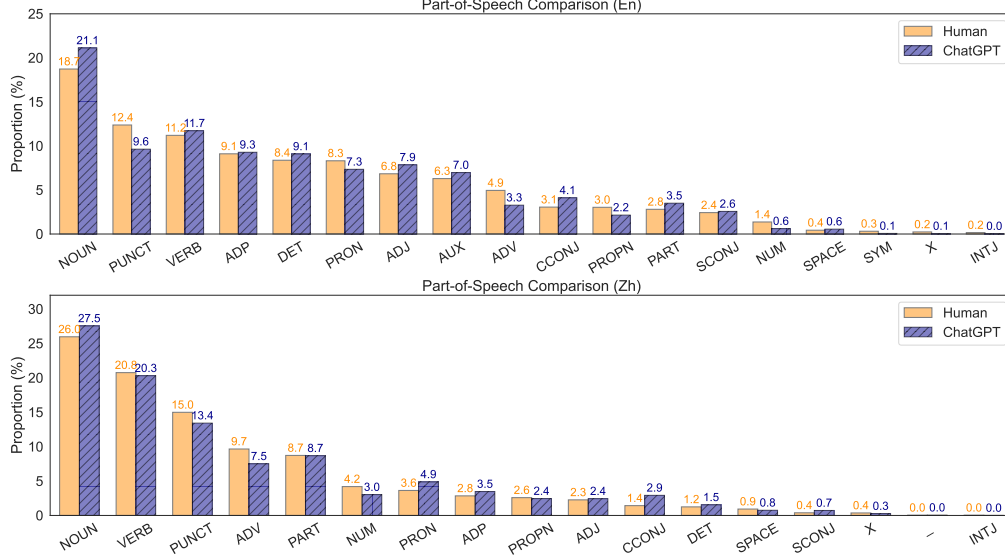


Figure 1: Part-of-Speech distribution comparison between ChatGPT and human answers. Results are sorted by POS proportion of human answers. The upper figure is for the HC3-English dataset and the lower is for the HC3-Chinese dataset.

A high proportion of nouns (NOUN) often indicates that the text is more argumentative, exhibiting informativeness and objectivity [24]. Accordingly, adposition (ADP) and adjective (ADJ) words also tend to appear more frequently [11]. The frequent co-occurrence of conjunctions (CCONJ) along with nouns, verbs, and adposition words indicates that the structure of the article and the relationships of cause-and-effect, progression, or contrast are clear. The above are also typical characteristics in academic papers or official documents [29]. We believe the RLHF training process has a great influence on ChatGPT’s writing style, which partly explains the difference in the POS tags distribution.

#### 4.2.2 Dependency Parsing

Dependency parsing is a technique that analyzes the grammatical structure of a sentence by identifying the dependencies between its words. We parse the answers in the corpus and compare the proportion of different dependency relations and their corresponding dependency distances. Figure 2 shows the comparison between humans and ChatGPT in HC3-English. Due to the limited space, the Chinese version is placed in the Appendix A.2.

The comparison of dependency relations exhibits similar characteristics to that of POS tags, where ChatGPT uses more determination, conjunction, and auxiliary relations. In terms of the dependency distance, ChatGPT has much longer distances for the punct and dep relations, which is perhaps due to the fact that ChatGPT tends to use longer sentences. However, ChatGPT has obviously shorter conj relations. According to the analysis of POS tags, ChatGPT usually uses more conjunctions than humans to make the content more logical, this may explain why the conj relations of ChatGPT are relatively shorter than humans.

#### 4.3 Sentiment Analysis

Humans are emotional beings, it is natural that our emotions are reflected in our words, to some extent. ChatGPT is learned on large-scale human-generated text, but it is further fine-tuned with human instructions. Therefore we are curious how "emotional" ChatGPT is compared with humans.

We use a multilingual sentiment classification model<sup>7</sup> fine-tuned on Twitter corpus [2] to conduct sentiment analysis for both English and Chinese comparison data. Note that deep learning-based models can be greatly influenced by some indicating words (such as "but" and "sorry" can easily

<sup>7</sup><https://huggingface.co/cardiffnlp/twitter-xlm-roberta-base-sentiment>

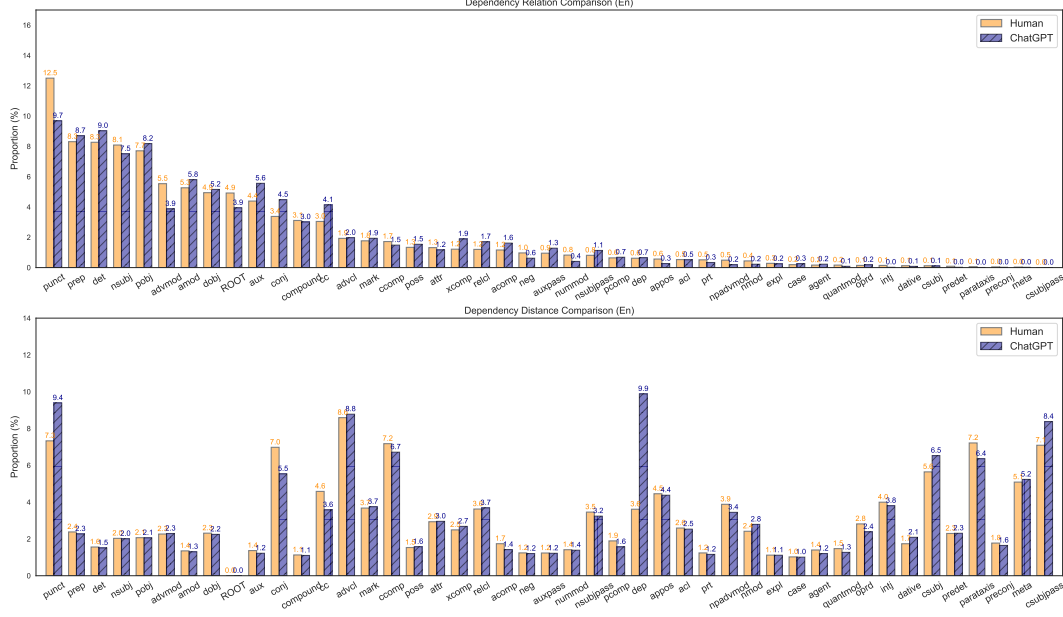


Figure 2: Top-30 dependency relations (upper) and corresponding dependency distances (lower) comparison between human and ChatGPT answers in HC3-English. Results are sorted by relations proportion of human answers.

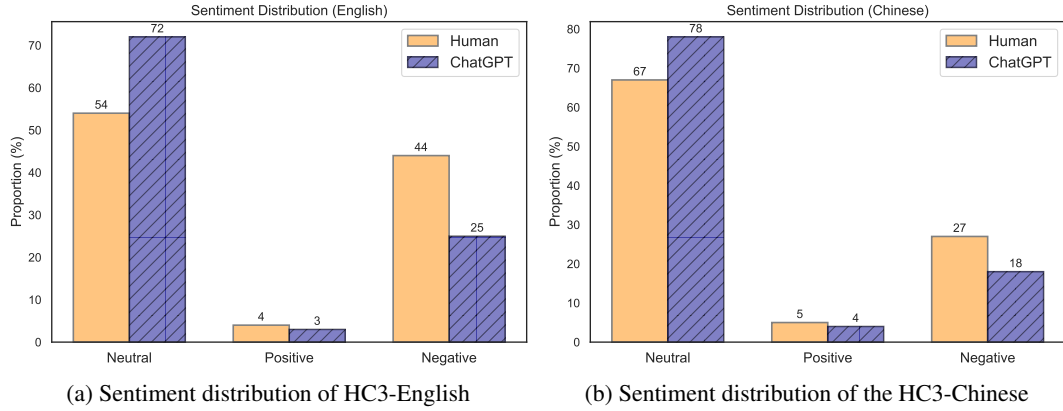


Figure 3: Proportions of three kinds of sentiments (neutral, positive, and negative) in our corpus.

fool the classifier to predict the "negative" label), making the predictions biased [16]. Therefore, the sentiment given by the classifier is only a reference to the true sentiment behind the text.

Figure 3 shows the comparison of the sentiment distribution of humans and ChatGPT. Several findings can be drawn from the results: First, we find that the proportion of neutral emotions is the largest for both humans and ChatGPT, which is in line with our expectations. However, **ChatGPT generally expresses more neutral sentiments than humans**. Then, the proportion of negative emotions is significantly higher than that of positive emotions. Notably, **humans express significantly more negative emotions than ChatGPT**. The proportion of humans' positive emotions is also slightly higher than that of ChatGPT. Overall, ChatGPT is less emotional than humans, though it is not completely emotionless.

#### 4.4 Language Model Perplexity

The perplexity (PPL) is commonly used as a metric for evaluating the performance of language models (LM). It is defined as the exponential of the negative average log-likelihood of the text under