# ChatGPT Evaluation for NLP
## A Meta Survey

Xiachong Feng
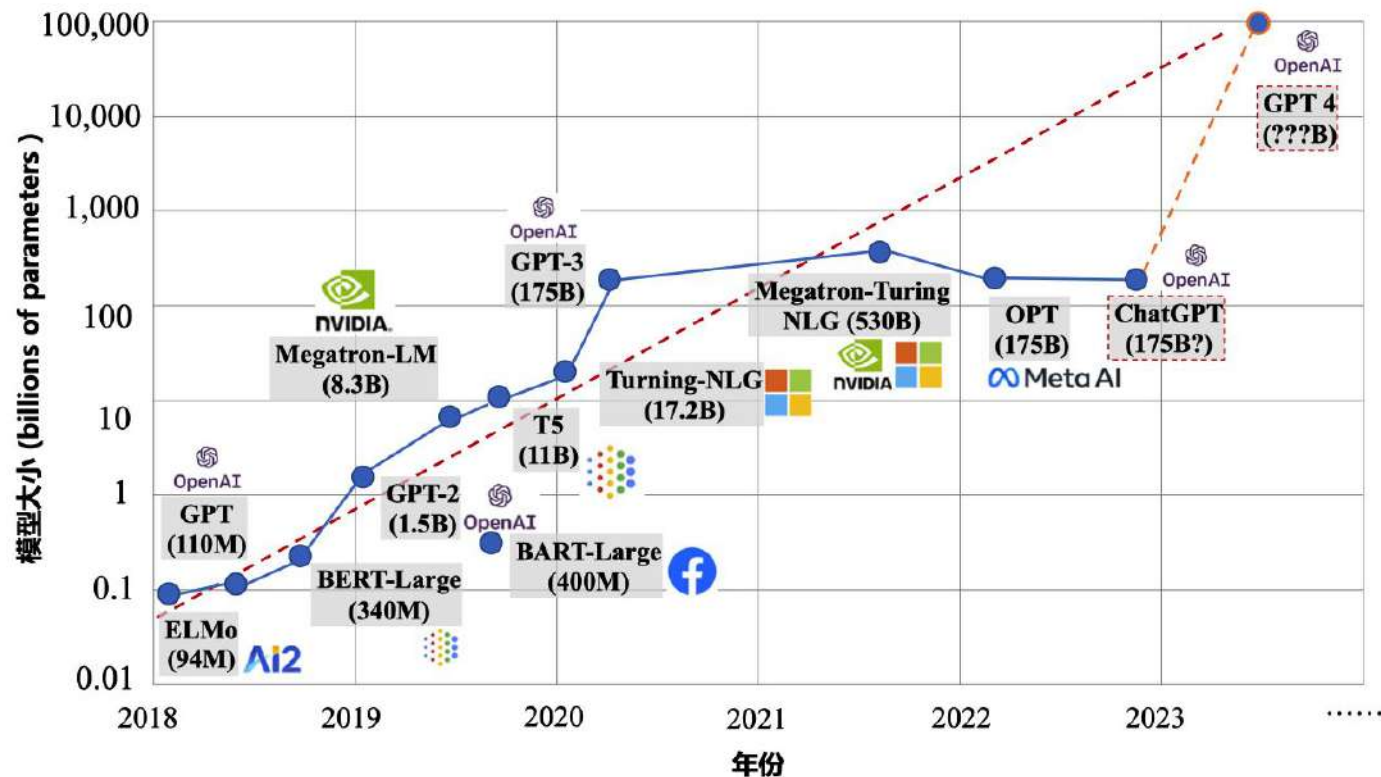Update: 2023.3.13

# Agenda

- **Introduction**

- **Social Media**

- **NLP Tasks**
    - **Summarization**
    - **Natural Language Generation Evaluation**
    - **Information Extraction**
    - **Machine Translation**
    - **Data Augmentation**
- **ChatGPT Failures**
- **Conclusion**

# Introduction

- Big-LLMs
- Instruction-tuning
- ChatGPT

# Introduction: Big-LLMs

# Introduction: Instruction-Tuning

2021.09.03   **[FLAN]** Finetuned Language Models Are Zero-Shot Learners

2021.10.15   **[T0]** Multitask Prompted Training Enables Zero-Shot Task Generalization

2022.01.28   **[CoT]** Chain of Thought Prompting Elicits Reasoning in Large Language Models

2022.03.04   **[InstructGPT]** Training language models to follow instructions with human feedback

2022.04.16   **[SUPER-NATURALINSTRUCTIONS]** Super-NaturalInstructions: Generalization via Declarative Instructions on 1600+ NLP Tasks

2022.05.24   **[Zero-shot CoT]** Large Language Models are Zero-Shot Reasoners

2022.10.20   **[FLAN-T5]** Scaling Instruction-Finetuned Language Models

  **[mT0]** Crosslingual Generalization through Multitask Finetuning

2022.11.30   **[🔥 🔥 🔥 ChatGPT]** Introducing ChatGPT

2023.02.27   **[LLaMA]** LLaMA: Open and Efficient Foundation Language Models

# Introduction: ChatGPT

# Social Media: Twitter

- ChatGPT: A Meta-Analysis after 2.5 Months 2023-2-20

# Social Media Analysis

- **Aim:** acquire insights into public opinion and sentiment on ChatGPT and understand public attitudes toward different topics related to ChatGPT.
- The dataset contains tweets across 61 languages. Over 68% of them are in English, other major languages are Japanese (6.4%), Spanish (5.3%), French (5.0%), and German (3.3%).

**#ChatGPT**

| Attribute | Detail |
|---|---|
| date range | 2022-11-30 to 2023-02-09 |
| number of tweets | 334,808 |
| language counts | 61 |
| English tweets | 228127 |
| number of users | 168,111 |

Table 1: Information of the collected Dataset

English Tweets → Sentiment Analysis

| Sentiment | Number of tweets |
|---|---|
| Positive | 100,163 |
| Neutral | 174,684 |
| Negative | 59,961 |

Table 2: Sentiment Distribution of all tweets.

*there is a relatively large proportion of positive sentiment, with 100k instances, and a smaller but still notable number of tweets of negative sentiments, with 60k instances.*

# Social Media Analysis



Observe an overall downward trend of sentiment (black solid line) during the course of ChatGPT's first 2.5 months: an initial rise in average sentiment was followed by a decrease from January 2023 onwards.

# Social Media Analysis



Observe an overall downward trend of sentiment (black solid line) during the course of ChatGPT's first 2.5 months: an initial rise in average sentiment was followed by a decrease from January 2023 onwards.

Overall tweets in English have a more positive perception of ChatGPT. This suggests that ChatGPT may be better in English, which constituted the majority of its training data; but see also our topic-based analysis below.

# Social Media Analysis



While the percentage of negative tweets is stable over time, the percentage of positive tweets decreases and there is a clear increase in tweets with the neutral sentiment.

This may indicate that the public view of ChatGPT is becoming more rational after an initial hype of this new "seemingly omnipotent" bot

# Social Media Analysis



During the course of 2.5 months after ChatGPT's debut, OpenAI announced 5 new releases claiming various updates.

Our data covers the period of the first three releases on the 15th of December 2022, the 9th of January, and the 3rd of January in 2023.

# Sentiment across Language



Figure 2: Weekly sentiment distribution averaged per language

Tweets in English have the most positive view of ChatGPT

The sentiment of English, German, and French tweets are trending downward while Spanish and Japanese tweets start from a low point and trend upwards.

# Sentiment across Topic



Figure 3: Upper: topic distribution per language. Lower: topic distribution over time for Japanese tweets.

5 major classes, which cover 86.3% of tweets in our dataset: science & technology (38.6%), learning & educational (15.2%), news & social concern (13.0%), diaries & daily life (10.2%), and business & entrepreneurs (9.3%).

The share of science & technology topic ranks the highest in all of the 5 languages.

# Sentiment across Topic



Figure 4: Sentiment distribution per topic.

business & entrepreneurs has the lowest proportion of negative tweets while the topic news & social concern contains the highest proportion of negative tweets.

# Sentiment: Human Evaluation

We manually annotated and analyzed the sentiment expressed within 40 randomly selected tweets.

- 20 random positive tweets from the period including the last two weeks of 2022 and the first week of 2023, where the general sentiment reaches the peak
- 20 random negative tweets from the second week to the fourth week of 2023, where the general sentiment declines

**For the first period [Positive]:**

ChatGPT's ability to generate <u>human-like and concise</u> text.

**For the second period [Negative]**

Potential <u>factual inaccuracies</u>, the <u>detectability of the model-generated text</u>, <u>ethical concerns</u>, <u>biased output</u> or the potential increase in <u>misinformation</u>, <u>job loss</u>

# Summarization

- Cross-Lingual Summarization via ChatGPT 2023-2-28
- Exploring the Limits of ChatGPT for Query or Aspect-based Text Summarization 2023-2-16

# Cross-lingual Summarization: Prompts

## Experimental Prompts for CLS

| | |
|---|---|
| E2E | Please summarize the following text in Chinese: *[English Doc]* |
| E2E+Interact | + Please make the Chinese summary shorter |
| Trans-Sum | Please first translate the following text to Chinese and then summarize the translated text in Chinese: *[English Doc]* |
| Trans-Sum+Interact | + Please make the Chinese summary shorter |
| Sum-Trans | Please first summarize the following text and then translate the summary to Chinese: *[English Doc]* |
| Sum-Trans+Interact | + Please make the Chinese summary shorter |

# Cross-lingual Summarization: Examples



Figure 1: An illustration of all prompts used to guide ChatGPT to perform cross-lingual summarization.

# Cross-lingual Summarization: Datasets

| Dataset | Src Lang. | Trg Lang. | Domain | Example | Doc. Length | Sum. Length |
|---------|-----------|-----------|--------|---------|-------------|-------------|
| CrossSum (Hasan et al., 2021) | English | Chinese | News | 3981 / 497 / 50 out of 497 | 814.2 | 35.6 |
| WikiLingua (Ladhak et al., 2020) | English | Chinese | How-to guide | 13211 / 1886 / 50 out of 3775 | 538.6 | 53.2 |
| | | German | | 40839 / 5833 / 50 out of 11669 | 526.1 | 63.4 |
| XSAMSum (Wang et al., 2022a) | English | Chinese | Dialogue | 14732 / 818 / 50 out of 819 | 140.1 | 27.6 |
| | | German | | 14732 / 818 / 50 out of 819 | 140.1 | 31.7 |

Table 1: Statistics of CLS datasets used in experiments. "*Src Lang.*" and "*Trg Lang*" denote the source and the target languages. "*Doc. Length*" and "*Sum. Length*" show the average length of source documents and target summaries (token level). "*Example*" lists the number of samples in each dataset w.r.t training, validation and test sets.

We randomly sample **50 documents** from the test set of each CLS dataset for evaluation.

# Cross-lingual Summarization: Results

| Method | CrossSum (En⇒Zh) | | | | WikiLingua (En⇒Zh) | | | | WikiLingua (En⇒De) | | | | XSAMSum (En⇒Zh) | | | | XSAMSum (En⇒De) | | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | R-1 | R-2 | R-L | B-S | R-1 | R-2 | R-L | B-S | R-1 | R-2 | R-L | B-S | R-1 | R-2 | R-L | B-S | R-1 | R-2 | R-L | B-S |
| mBART-50 | 26.1 | 7.4 | 22.1 | 65.4 | 32.1 | 10.4 | 26.8 | 68.5 | 26.8 | 7.7 | 20.5 | 62.5 | 40.6 | 14.4 | 33.9 | 74.5 | 42.4 | 18.9 | 35.4 | 73.7 |
| Text-davinci-003 | 18.7 | 3.6 | 14.7 | 60.2 | 23.6 | 3.8 | 17.8 | 60.9 | 18.8 | 2.6 | 12.2 | 60.7 | 24.4 | 8.0 | 20.7 | 63.4 | 35.5 | 12.4 | 27.3 | 62.4 |
| ChatGPT (e2e) | 14.2 | 3.3 | 10.3 | 60.3 | 20.9 | 5.6 | 15.5 | 62.7 | 16.9 | 2.1 | 10.7 | 60.1 | 21.3 | 5.5 | 17.1 | 63.5 | 32.0 | 10.3 | 24.5 | 61.4 |
| ChatGPT (e2e+interact) | 22.1 | 3.8 | 15.6 | 61.8 | 28.4 | 6.5 | 22.1 | 64.5 | **22.4** | 2.8 | 14.7 | **61.3** | 27.2 | 6.9 | 22.9 | 67.5 | **39.6** | **16.0** | **31.4** | **64.3** |
| ChatGPT (Trans-Sum) | 15.8 | 3.3 | 11.9 | 60.9 | 24.8 | 5.4 | 19.1 | 62.9 | 19.4 | 2.4 | 12.6 | 60.0 | 26.0 | 7.3 | 21.2 | 66.4 | 33.2 | 9.6 | 25.3 | 61.1 |
| ChatGPT (Trans-Sum+interact) | **22.6** | **4.1** | **16.9** | **62.7** | 26.1 | 5.3 | 19.7 | 63.7 | 21.6 | 2.4 | 15.1 | 60.8 | 27.4 | 6.7 | 22.4 | 67.1 | 39.4 | 13.5 | 29.4 | 63.3 |
| ChatGPT (Sum-Trans) | 16.5 | 3.8 | 12.0 | 60.8 | 27.2 | 7.3 | 20.3 | 64.3 | 21.3 | **3.5** | 14.4 | 60.9 | 26.8 | 7.7 | 21.3 | 66.7 | 31.7 | 8.8 | 23.5 | 60.8 |
| ChatGPT (Sum-Trans+interact) | 21.6 | 3.5 | 15.5 | 61.7 | **30.1** | **8.1** | **22.4** | **64.9** | 21.4 | 3.1 | **15.4** | 60.6 | **31.4** | **11.5** | **28.1** | **70.1** | 35.9 | 13.2 | 29.0 | 62.8 |

(Zero-Shot rows: Text-davinci-003 through ChatGPT (Sum-Trans+interact))

Table 2: Experimental results on CrossSum, WikiLingua and XSAMSum.

| Method | CrossSum | WikiLingua | | XSAMSum | |
|---|---|---|---|---|---|
| | En⇒Zh | En⇒Zh | En⇒De | En⇒Zh | En⇒De |
| Text-davinci-003 | 83.3 | 78.5 | 149.1 | 61.8 | 62.5 |
| e2e | 183.7 | 176.6 | 273.5 | 68.6 | 75.3 |
| e2e+interact | 66.4 | 50.0 | 80.7 | 28.7 | 42.5 |
| TransSum | 155.1 | 82.1 | 149.3 | 48.2 | 60.9 |
| TransSum+interact | 63.4 | 46.2 | 70.0 | 30.3 | 41.1 |
| SumTrans | 132.7 | 94.3 | 124.2 | 54.9 | 68.1 |
| SumTrans+interact | 57.8 | 50.1 | 71.6 | 29.3 | 37.5 |

Table 3: The average length (token level) of the generated summaries on the test set of each dataset.

**Interactive prompt is very important**

**mBART-50 Is still very strong**
(needs more carefully-designed human evaluation)

# Query-based Summarization: Dataset

| Type | Dataset | Domain | #Input Tk. | #Output Tk. | #Asp. Type |
|------|---------|--------|------------|-------------|------------|
| Query | QMSum (Zhong et al., 2021) | Meeting | 9,070(2,505*) | 70 | 1,566 |
| | SQuALITY (Wang et al., 2022) | Story | 6,052 | 252 | 437 |
| Aspect | CovidET (Zhan et al., 2022) | Reddit | 192 | 27 | 7 |
| | NEWTS (Bahrainian et al., 2022) | News | 602 | 74 | 50 |

Table 1: Statistics of the query/aspect-based summarization datasets that we used. **#Input Tk.** and **#Output Tk.** represent the number of input and output token lengths, respectively. **#Asp. Type** is the number of all aspect types. 2,505* represents the average token number in golden inputs.

# Query-based Summarization: Prompts

## SQuALITY

Q: Query. Answer the question in around 200 words. Article: story. **specific question**

Q: Query. Answer the question in around 450 words. Article: story.
Your response is too short. Please answer it in around 450 words. **general question**

## QMSum

Q: Query. Article: meeting

Q: Query. Article: golden meeting

meeting is the initial meeting, while golden meeting is the provided golden spans

## CovidET

Q: Summarize this article with respect to Aspect within one short sentence. Article0. A: Answer0. Q: Summarize this article with respect to Aspect within one short sentence. Article. A:

## NEWTS

Article. Summarize this article with respect to Aspect:

# Query-based Summarization: Results

| Datasets | Models | R-1 | R-2 | R-L | R-Lsum |
|---|---|---|---|---|---|
| CovidET | Fine-tuning | **26.19** | **6.85** | **17.86** | **20.82** |
| | ChatGPT | 20.81 | 3.99 | 15.35 | 15.36 |
| NEWTS | Fine-tuning | 31.78 | 10.83 | 20.54 | — |
| | ChatGPT | **32.54** | **11.37** | **20.74** | 20.74 |
| QMSum | Fine-tuning | **32.29** | 8.67 | **28.17** | — |
| | ChatGPT | 28.34 | **8.74** | 17.81 | 18.01 |
| QMSum(Golden) | Fine-tuning | 36.06 | 11.36 | **31.27** | — |
| | ChatGPT | **36.83** | **12.78** | 24.23 | 24.19 |
| SQuaLITY | Fine-tuning | **38.20** | **9.00** | **20.20** | — |
| | ChatGPT | 37.02 | 8.19 | 18.45 | 22.56 |
| Avg. | Fine-tuning | **32.90** | **9.34** | **23.61** | — |
| | ChatGPT | 30.94 | 8.96 | 19.22 | — |

Table 2: Comparison between ChatGPT zero-shot performance with previous Fine-tuning(FT) results.

⭐ **Fine-tuning methods are still strong.**

# Query-based Summarization: Results

| Datasets | Text | U-1-gram | U-2-gram | U-3-gram | U-4-gram | Coverage | Density | Compression |
|---|---|---|---|---|---|---|---|---|
| CovidET | *Reference* | 0.59 | 0.95 | 0.99 | 0.99 | 0.60 | 0.90 | 11.84 |
| | *ChatGPT* | 0.42 | 0.86 | 0.97 | 0.99 | 0.66 | 1.19 | 6.92 |
| NEWTS | *Reference* | 0.30 | 0.69 | 0.84 | 0.90 | 0.78 | 2.84 | 9.66 |
| | *ChatGPT* | 0.15 | 0.50 | 0.69 | 0.81 | 0.89 | 4.13 | 4.03 |
| QMSum | *Reference* | 0.19 | 0.65 | 0.88 | 0.96 | 0.87 | 2.15 | 99.76 |
| | *ChatGPT* | 0.26 | 0.73 | 0.89 | 0.95 | 0.84 | 1.97 | 44.86 |
| QMSum(Golden) | *Reference* | 0.29 | 0.69 | 0.88 | 0.95 | 0.77 | 2.05 | 15.68 |
| | *ChatGPT* | 0.26 | 0.66 | 0.85 | 0.93 | 0.79 | 2.15 | 10.74 |
| SQuaLITY | *Reference* | 0.25 | 0.85 | 0.97 | 0.99 | 0.83 | 1.55 | 32.57 |
| | *ChatGPT* | 0.33 | 0.82 | 0.93 | 0.97 | 0.81 | 1.81 | 24.83 |

Table 3: Comparison between ChatGPT zero-shot performance with references on various metrics. U-1/2/3/4-gram represents unique 1/2/3/4 gram.

⭐ **ChatGPT is smart but not that smart.**

# NLG Evaluator

- Is ChatGPT a Good NLG Evaluator? A Preliminary Study 2023-3-7

# NLG Evaluator: Prompt

## Is ChatGPT a good NLG evaluator?

Score the following [task-ins] with respect to [aspect] with one to five stars, where one star means "[ant-aspect]" and five stars means "perfect [aspect]". Note that [aspect] measures [aspect-ins].

[Conditioned Text]
[Generated Text]
Stars:

Score the following news summarization given the corresponding news with respect to fluency with one to five stars, where one star means "disfluency" and five stars means "perfect fluency". Note that fluency measures the quality of individual sentences, are they well-written and grammatically correct. Consider the quality of individual sentences.

News: [a news article]
Summary: [one generated summary]
Stars:

**ChatGPT**

I would give this news summary **four stars** for fluency. The summary is well-written and captures the main points of the news article

# NLG Evaluator: Metrics

- **Spearman correlation** (Zar, 2005) assesses the monotonic relationships between two variables;

- **Pearman correlation** (Mukaka, 2012) measures the linear relationships between two sets of data;

- **Kendall's Tau** (Kendall, 1938) evaluates the ordinal association between two measured quantities.

$$\text{Corr}_{\text{sample}} = \frac{1}{n} \sum_{1}^{n} (\rho([f_{\text{auto}}(g_{i,1}), ..., f_{\text{auto}}(g_{i,M})], [f_{\text{human}}(g_{i,1}), ..., f_{\text{human}}(g_{i,M})]))$$

$$\text{Corr}_{\text{dataset}} = \rho([f_{\text{auto}}(g_{1,1}), ..., f_{\text{auto}}(g_{n,M})], [f_{\text{human}}(g_{1,1}), ..., f_{\text{human}}(g_{n,M})])$$

Conditioned text {c1, c2, ..., cn} and M NLG models
The generated text of m-th model for the i-th condition text is denoted as g_i,m.

# NLG Evaluator: Results (Summarization for Example)

| Metrics | Coherence | | | Relevance | | | Consistency | | | Fluency | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | Spear. | Pear. | Kend. | Spear. | Pear. | Kend. | Spear. | Pear. | Kend. | Spear. | Pear. | Kend. |
| ROUGE-1 | 0.167 | 0.160 | 0.126 | 0.326 | 0.359 | 0.252 | 0.160 | 0.224 | 0.130 | 0.115 | 0.158 | 0.094 |
| ROUGE-2 | 0.184 | 0.174 | 0.139 | 0.290 | 0.327 | 0.219 | 0.187 | 0.246 | 0.155 | 0.159 | 0.185 | 0.128 |
| ROUGE-L | 0.128 | 0.102 | 0.099 | 0.311 | 0.342 | 0.237 | 0.115 | 0.189 | 0.092 | 0.105 | 0.141 | 0.084 |
| BERTScore | 0.283 | 0.310 | 0.211 | 0.311 | 0.346 | 0.243 | 0.110 | 0.152 | 0.090 | 0.192 | 0.209 | 0.158 |
| MoverScore | 0.159 | 0.167 | 0.118 | 0.318 | 0.371 | 0.244 | 0.157 | 0.224 | 0.127 | 0.129 | 0.176 | 0.105 |
| PRISM | 0.249 | 0.258 | 0.196 | 0.212 | 0.232 | 0.163 | 0.345 | 0.352 | 0.285 | 0.254 | 0.264 | 0.205 |
| BARTScore | 0.322 | 0.345 | 0.250 | 0.264 | 0.290 | 0.197 | 0.311 | 0.321 | 0.256 | 0.248 | 0.260 | 0.203 |
| BARTScore+CNN | 0.448 | 0.458 | 0.342 | 0.356 | 0.369 | 0.273 | 0.382 | 0.422 | 0.315 | 0.356 | 0.407 | 0.292 |
| BARTScore+CNN+Para | 0.424 | 0.442 | 0.325 | 0.313 | 0.364 | 0.241 | 0.401 | 0.487 | 0.332 | **0.378** | **0.448** | 0.311 |
| ChatGPT | **0.470** | **0.484** | **0.403** | **0.428** | **0.454** | **0.374** | **0.419** | **0.517** | **0.389** | 0.353 | 0.415 | **0.329** |

Table 1: Sample-level Spearman correlation (Spear.) correlation, Pearman (Pear.) correlation and Kendall's Tau (Kend.) of different aspects on SummEval.

| Metrics | Coherence | | | Relevance | | | Consistency | | | Fluency | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | Spear. | Pear. | Kend. | Spear. | Pear. | Kend. | Spear. | Pear. | Kend. | Spear. | Pear. | Kend. |
| ROUGE-1 | 0.184 | 0.193 | 0.129 | 0.302 | 0.341 | 0.217 | 0.137 | 0.175 | 0.108 | 0.080 | 0.143 | 0.062 |
| ROUGE-2 | 0.145 | 0.140 | 0.102 | 0.245 | 0.254 | 0.175 | 0.129 | 0.152 | 0.102 | 0.062 | 0.092 | 0.048 |
| ROUGE-L | 0.141 | 0.148 | 0.100 | 0.284 | 0.318 | 0.204 | 0.109 | 0.152 | 0.086 | 0.079 | 0.132 | 0.061 |
| BERTScore | 0.317 | 0.326 | 0.224 | 0.362 | 0.381 | 0.262 | 0.117 | 0.146 | 0.092 | 0.150 | 0.196 | 0.117 |
| MoverScore | 0.178 | 0.177 | 0.125 | 0.294 | 0.328 | 0.211 | 0.150 | 0.171 | 0.118 | 0.119 | 0.160 | 0.092 |
| PRISM | 0.286 | 0.281 | 0.204 | 0.280 | 0.297 | 0.202 | 0.323 | 0.297 | 0.256 | 0.236 | 0.248 | 0.184 |
| BARTScore | 0.335 | 0.357 | 0.241 | 0.363 | 0.386 | 0.263 | 0.269 | 0.276 | 0.212 | 0.187 | 0.206 | 0.146 |
| BARTScore+CNN | 0.408 | 0.434 | 0.292 | 0.394 | 0.423 | 0.286 | 0.334 | 0.377 | 0.264 | 0.285 | 0.354 | 0.223 |
| BARTScore+CNN+Para | 0.424 | 0.430 | 0.304 | 0.398 | 0.431 | 0.289 | **0.379** | 0.452 | 0.301 | **0.346** | **0.410** | **0.271** |
| ChatGPT | **0.435** | **0.438** | **0.353** | **0.448** | **0.459** | **0.366** | 0.356 | **0.515** | **0.320** | 0.300 | 0.385 | 0.268 |

Table 2: Dataset-level Spearman correlation (Spear.) correlation, Pearman (Pear.) correlation and Kendall's Tau (Kend.) of different aspects on SummEval.

**SummEval** collects 16 model-generated summaries on the CNN/DM dataset and annotates human judgments upon these summaries covering aspects of coherence, relevance, consistency and fluency.

**ChatGPT achieves state-of-the-art or competitive correlation with golden human judgments.**

Is ChatGPT a Good NLG Evaluator? A Preliminary Study

# Machine Translation

- Is ChatGPT A Good Translator? A Preliminary Study 2023-1-31
- Towards Making the Most of ChatGPT for Machine Translation 2023-3-11

# Machine Translation: Prompts

**Ask ChatGPT for providing prompts.**



W Provide ten concise prompts or templates that can make you translate.

1. "Translate this sentence from English to French: __"
2. "What does this word mean in Spanish? __"
3. "Can you help me translate this paragraph from Chinese to English? __"
4. "Please provide the German translation for this phrase: __"
5. "I need the Japanese translation of this word: __"
6. "Translate this sentence from Russian to English: __"
7. "What is the Arabic translation for this phrase? __"
8. "Can you help me translate this message from Korean to English? __"
9. "Please provide the Italian translation for this word: __"
10. "I need the Dutch translation of this sentence: __"

Figure 1: Prompts advised by ChatGPT for machine translation (Date: 2022.12.16).

Table 2: Candidate translation prompts.

| Translation Prompt |
| --- |
| TP1 Translate these sentences from [SRC] to [TGT]: |
| TP2 Answer with no quotes. What do these sentences mean in [TGT]? |
| TP3 Please provide the [TGT] translation for these sentences: |

# Machine Translation: Baseline Results

- Chinese-to-English (Zh⇒En) translation task with the test set from Flores-101
- We randomly sample 50 sentences from each set for evaluation.

Table 2: Candidate translation prompts.

| | Translation Prompt |
|---|---|
| TP1 | Translate these sentences from [SRC] to [TGT]: |
| TP2 | Answer with no quotes. What do these sentences mean in [TGT]? |
| TP3 | Please provide the [TGT] translation for these sentences: |

Table 3: Comparison of different prompts for ChatGPT to perform Chinese-to-English (Zh⇒En) translation.

| System | BLEU↑ | ChrF++↑ | TER↓ |
|---|---|---|---|
| Google | 31.66 | 57.09 | 56.21 |
| DeepL | 31.22 | 56.74 | 57.84 |
| Tencent | 29.69 | 56.24 | 57.16 |
| ChatGPT w/ TP1 | 23.25 | 53.07 | 66.03 |
| ChatGPT w/ TP2 | 24.54 | 53.05 | 63.79 |
| ChatGPT w/ TP3 | **24.73** | **53.71** | **62.84** |

Follow experiments are based on this prompt.

**Still lags behind the baselines by at least 5.0 BLEU points**

# Machine Translation: Multilingual Translation Results

Better knowledge transfer within the **same family** than between different families .

Relative to Google Translate

- ChatGPT in multilingual translation, including German (De), English (En), Romanian (Ro), and Chinese (Zh)
- The first three languages come from the same family with Latin scripts while the last is from another family with Chinese scripts
- We randomly sample 50 sentences from each set for evaluation.

Table 4: Performance of ChatGPT for multilingual translation.

High-resource task   Low-resource task

| System | De-En | | Ro-En | | Zh-En | |
|--------|-------|------|-------|------|-------|------|
| | ⇒ | ⇐ | ⇒ | ⇐ | ⇒ | ⇐ |
| Google | 45.04 | 41.16 | 50.12 | 46.03 | 31.66 | 43.58 |
| DeepL | 49.23(+9.3%) | 41.46(+0.7%) | 50.61(+0.9%) | 48.39(+5.1%) | 31.22(-1.3%) | 44.31(+1.6%) |
| Tencent | n/a | n/a | n/a | n/a | 29.69(-6.2%) | 46.06(+5.6%) |
| ChatGPT | 43.71(-2.9%) | 38.87(-5.5%) | 44.95(-10.3%) | 24.85(-46.0%) | 24.73(-21.8%) | 38.27(-12.1%) |

| System | De-Zh | | Ro-Zh | | De-Ro | |
|--------|-------|------|-------|------|-------|------|
| | ⇒ | ⇐ | ⇒ | ⇐ | ⇒ | ⇐ |
| Google | 38.71 | 21.68 | 39.05 | 25.59 | 33.31 | 32.27 |
| DeepL | 40.46(+4.5%) | 22.82(+5.2%) | 38.95(-0.2%) | 25.39(-0.7%) | 35.19(+5.6%) | 34.27(+6.1%) |
| Tencent | 40.66(+5.0%) | 19.44(-10.3%) | n/a | n/a | n/a | n/a |
| ChatGPT | 34.46(-10.9%) | 19.80(-8.6%) | 30.84(-21.0%) | 19.17(-25.0%) | 33.38(+0.2%) | 29.89(-7.3%) |

The huge resource **difference of mono-lingual data between English and Romanian** limits the language modeling capability of Romanian, which partially explains the **poor performance on English⇒Romanian**.
⭐**Romanian Language Modeling is poor!**

# Machine Translation: Pivot Prompting

Ask ChatGPT to translate the source sentence into a high-resource pivot language (i.e., English by default) first and then into the target language.

**Prompt**

```
Please provide the [PIV]
translation first and then the
[TGT] translation for these
sentences one by one:
```

Table 5: Performance of ChatGPT with pivot prompting. New results are obtained from the updated ChatGPT version on 2023.01.31. LR: length ratio.

| System | De⇒Zh | | Ro⇒Zh | |
|---|---|---|---|---|
| | BLEU | LR | BLEU | LR |
| Google | 38.71 | 0.94 | 39.05 | 0.95 |
| DeepL | 40.46 | 0.98 | 38.95 | 0.99 |
| ChatGPT (Direct) | 34.46 | 0.97 | 30.84 | 0.91 |
| ChatGPT (Direct$_{new}$) | 30.76 | 0.92 | 27.51 | 0.93 |
| ChatGPT (Pivot$_{new}$) | 34.68 | 0.95 | 34.19 | 0.98 |

WE Please provide the English translation first and then the Chinese translation for these sentences one by one:
Er sagte, dass er eine WLAN-Türklingel gebaut habe.
Dies bestätigt nicht nur, dass zumindest einige Dinosaurier Federn hatten, eine Theorie, die bereits weit verbreitet ist, sondern liefert auch Details, die Fossilien im Allgemeinen nicht liefern können, wie etwa Farbe und dreidimensionale Anordnung.
Der 802.11n-Standard läuft auf den Frequenzen 2,4 Ghz sowie 5,0 Ghz.

1. He said that he had built a Wi-Fi doorbell.
   他说他建造了一个Wi-Fi门铃。
2. This not only confirms that at least some dinosaurs had feathers, a theory that is already widespread, but also provides details that fossils generally cannot provide, such as color and three-dimensional arrangement.
   这不仅证实了至少一些恐龙有羽毛的说法，一个已经广泛存在的理论，同时也提供了化石通常无法提供的详细信息，如颜色和三维排列。
3. The 802.11n standard operates on frequencies 2.4 GHz and 5.0 GHz.
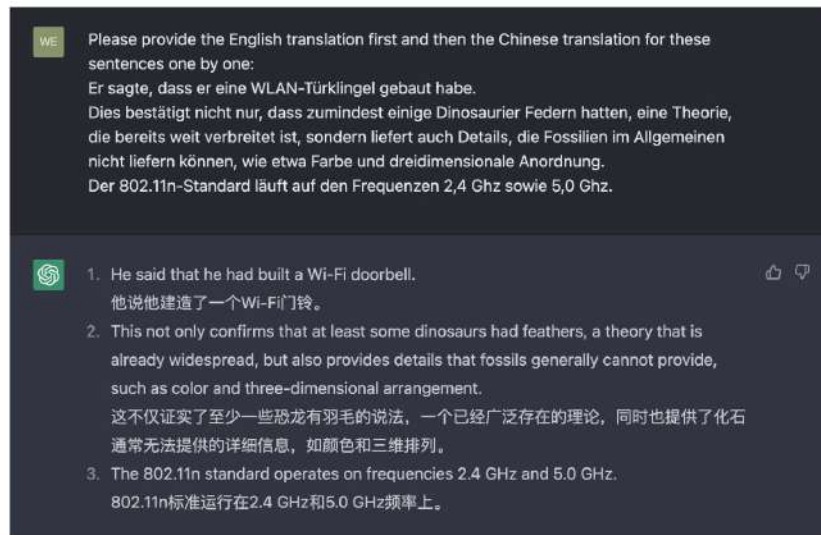   802.11n标准运行在2.4 GHz和5.0 GHz频率上。

Figure 2: Translation results by ChatGPT with pivot prompting (Date: 2023.01.31).

# Machine Translation: Translation Robustness

- To test the translation robustness, we adopt the test set of WMT19 Biomedical Translation Task (Bawden et al., 2019, i.e., Bio) and the set2 and set3 of WMT20 Robustness Task.
- We randomly sample 50 sentences from each set for evaluation.
  - WMT19 Bio test set is composed of Medline abstracts, which require domain-specific knowledge to handle the terminologies.
  - WMT20 Rob2 are comments from the social media website *reddit.com* that could contain various errors, including spelling/typographical errors, word omission/insertion/repetition, grammatical errors, spoken languages, Internet slang, and so on.

Table 6: Performance of ChatGPT for translation robustness.

| System | WMT19 Bio | WMT20 Rob2 | | WMT20 Rob3 |
|---|---|---|---|---|
| | De⇒En | En⇒Ja | Ja⇒En | De⇒En |
| Google | 37.83 | 29.72 | 19.21 | 42.91 |
| DeepL | 37.13 | 26.25 | 19.83 | 41.29 |
| ChatGPT | 33.22 | 22.36 | 18.34 | 44.59 |

WMT20 Rob3 test set that contains a **crowdsourced speech recognition corpus.** It suggests that ChatGPT, which is essentially an artificial intelligent chatting machine, **is capable of generating more natural spoken languages** than these commercial translation systems.

# Weakness and Motivation

**Previous:** Adopt simple prompts and basic settings regardless of the significant influence of the prompts' quality.

**This paper:** In this report, we aim to further elicit the capability of ChatGPT by revisiting the following three aspects and correspondingly propose two simple but effective prompts: Task-Specific Prompts (TSP) and Domain-Specific Prompts (DSP).

## Temperature

Decoding with higher temperatures displays greater linguistic variety, a diverse generation may impede its translation quality.

## Task Information

The task inconsistency (ChatGPT is a conversational system) will limit its translation ability to a certain degree. In response to this problem, we proposed **Task-Specific Prompts (TSP)** to further emphasize the task information to bridge the task gap, i.e., conversation and translation.

## Domain Information

ChatGPT can incorporate additional information, like human interactions, through the input prompts. We argue that such flexible interaction may alleviate some classical MT challenges, e.g., crossdomain generalization. Therefore, propose **Domain-Specific Prompts (DSP)** to introduce the domain navigation information to elicit ChatGPT's generalization ability across different domains.

# Machine Translation: Experimental Setting and Datasets

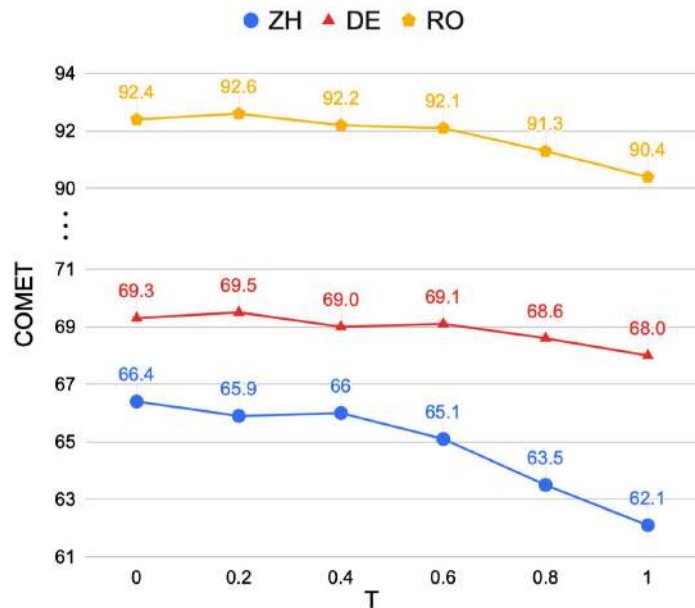| Test Set | Direction | Domain | Size |
|---|---|---|---|
| Flores-200 | Any | General | 1,012 |
| WMT19 News | En⇒Zh | News | 2,001 |
| | En⇒De | | 3,004 |
| WMT19 Bio | En⇒Zh | Biomedical | 224 |
| | Zh⇒En | | 241 |

Table 1: Data statistics and descriptions.

ChatGPT: gpt3.5-turbo-0301 models

Flores-200 for Multilingual MT

WMT19 for Cross-domain MT

We test all samples through OpenAI API.

# Machine Translation: Effect of Temperature



Figure 1: The relationship between temperature and ChatGPT's performance (in terms of COMET scores) when translating from English to other languages.

ChatGPT's performance largely depends on the temperatures, especially in difficult languages. Generally, setting **a lower temperature can result in higher performance**.

The impact of temperature is relatively small when translating to high-resource languages, while for complex languages, e.g., Chinese, it has a large degradation in performance

# Machine Translation: Effect of Task Information

Emphasizing the task information in prompts can further improve ChatGPT's performance, especially in complex tasks.

**TSP:** Prepend the sentence **"You are a machine translation system."** to "Please provide the [TGT] translation for the following sentence:"

| System | COMET | BLEU | ChrF | COMET | BLEU | ChrF |
|---|---|---|---|---|---|---|
| High-resource task | | DE⇒EN | | | EN⇒DE | |
| Google Translator | **77.7** | **47.4** | **70.5** | **70.5** | **44.4** | **68.9** |
| ChatGPT | 77.2 | 43.5 | 69.4 | 69.3 | 40.6 | 67.1 |
| ChatGPT + TSP | 77.5 | 44.1 | 69.7 | 69.4 | 40.4 | 67.0 |
| Distant language | | ZH⇒EN | | | EN⇒ZH | |
| Google Translator | **73.5** | **33.5** | **61.2** | **68.5** | **48.8** | **43.8** |
| ChatGPT | 71.3 | 26.4 | 58.3 | 66.4 | 45.1 | 39.0 |
| ChatGPT + TSP | 71.5 | 26.7 | 58.4 | 67.2 | 45.3 | 39.3 |
| Low-resource task | | RO⇒EN | | | EN⇒RO | |
| Google Translator | **82.4** | **48.0** | **71.2** | 91.6 | **43.3** | **67.0** |
| ChatGPT | 80.6 | 41.8 | 68.8 | 92.4 | 40.6 | 65.5 |
| ChatGPT + TSP | 80.8 | 41.9 | 69.0 | **92.9** | 40.8 | 65.7 |
| | | ZH⇒RO | | | RO⇒ZH | |
| Google Translator | 73.9 | **25.8** | **53.9** | 62.3 | **42.3** | **37.8** |
| ChatGPT | 73.8 | 20.9 | 51.5 | 58.9 | 37.7 | 33.3 |
| ChatGPT + TSP | **74.1** | 21.0 | 51.3 | 59.1 | 38.0 | 33.7 |

**Non-English-Centric Language Pairs**

When tackling non English-centric MT language pairs, ChatGPT tends to generate **hallucinations**.
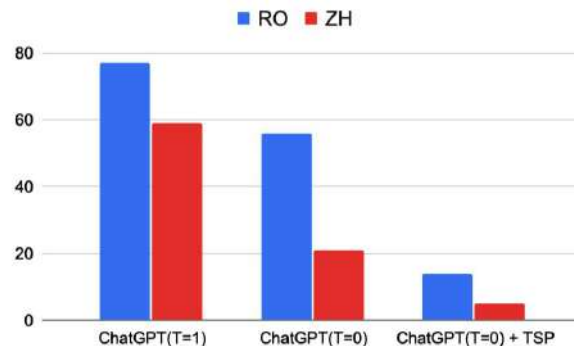


Figure 3: Number of Post-Edited sentences in non-English-centric language pairs, where a higher value means the translation contains more hallucinations. RO represents the translation for ZH⇒RO, while ZH represents the translation for ZH⇒RO.

# Machine Translation: Effect of Domain Information

| Method | Translation Prompt |
|---|---|
| **ChatGPT** | "role": "system", "content": "You are a machine translation system.", "role": "user", "content": 'Please provide the [TGT] translation for the following sentence: ' |
| **ChatGPT+DSP** | "role": "system", "content": "You are a machine translation system that translates sentences in the [DOM] domain.", "role": "user", "content": 'Please provide the [TGT] translation for the following sentence: ' |
| **ChatGPT+F-DSP** | "role": "system", "content": "You are a machine translation system that translates sentences in the [FDOM] domain.", "role": "user", "content": 'Please provide the [TGT] translation for the following sentence: ' |

Table 4: Domain-Specific translation prompts. "[DOM]" and "[FDOM]" denote the correct and incorrect domain instructions, respectively.

| System | WMT19 Bio | | | | WMT19 News | | | |
|---|---|---|---|---|---|---|---|---|
| | EN⇒ZH | | ZH⇒EN | | EN⇒ZH | | EN⇒DE | |
| | COMET | BLEU | COMET | BLEU | COMET | BLEU | COMET | BLEU |
| Google Translator | **59.4** | **38.8** | **57.1** | **36.1** | 59.3 | **43.4** | **64.1** | **33.7** |
| ChatGPT | 58.6 | 35.5 | 56.5 | 31.2 | 58.8 | 39.6 | 63.1 | 31.3 |
| ChatGPT + DSP | 58.9 | 35.8 | 57.1 | 31.0 | 59.6 | 39.8 | 63.2 | 31.5 |
| ChatGPT + F-DSP | 58.6 | 35.6 | 55.8 | 31.2 | 57.9 | 39.0 | 62.0 | 31.2 |

Table 5: Performance of ChatGPT on translation robustness, i.e., different domains. "DSP" denotes our proposed domain-specific prompting method, while "F-DSP" denotes the false domain-specific prompting, i.e., we specify wrong/unrelated domain information in the prompt. The results in green denote that "DSP" improves ChatGPT by a clear margin (0.5 (↑) score), while the red results denote the significant performance drops caused by "F-DSP".

Introducing the **correct domain information** consistently **improves** ChatGPT's performance while **wrong domain information** leads to significant **degradation** in performance.

# Machine Translation: In-context Learning

| System | EN ⇒ DE | | EN ⇒ ZH | | EN ⇒ RO | |
|---|---|---|---|---|---|---|
| | **COMET** | **BLEU** | **COMET** | **BLEU** | **COMET** | **BLEU** |
| Google Translator | **70.5** | **44.4** | **68.5** | **48.8** | 91.6 | **43.3** |
| ChatGPT | 69.4 | 40.4 | 67.2 | 45.3 | 92.9 | 40.8 |
| *Standard few-shot prompting* | | | | | | |
| -w/ 1-shot | 69.9 | 40.5 | <u>68.2</u> | 45.6 | 93.4 | <u>41.2</u> |
| -w/ 3-shot | 69.7 | <u>40.6</u> | 68.1 | <u>45.7</u> | **93.5** | 41.1 |

Table 6: Few-shot translation performance of ChatGPT on 20 samples randomly selected from Flores-200. In the standard few-shot prompting setting, we randomly sample 1/3 examples from the development set.

**Few-shot ICL can further improve ChatGPT's performance.**

# Machine Translation: CoT

| Method | Translation Prompt |
|---|---|
| **Zero-Shot CoT** | "role": "system", "content": "You are a machine translation system.", "role": "user", "content": 'Please provide the German translation for the following sentence step by step and then provide the complete sentence: ' |
| **1-Shot CoT** | "role": "system", "content": "You are a machine translation system.", "role": "user", "content": 'Please provide the German translation for the following sentence step by step and then provide the complete sentence [S] 1. [S_1] – [T_1] 2. [S_2] – [T_2] ... n. [S_n] – [T_n] The complete sentence in [TGT] is: [T] Please provide the German translation for the following sentence step by step and then provide the complete sentence:' |

Table 7: The templates of Zero-Shot CoT and 1-shot CoT. [S_n] represents the $n$-th token in source demonstration [S], [T_n] represents the $n$-th token in target demonstration [T].

| Method | EN⇒DE | | EN⇒ZH | |
|---|---|---|---|---|
| | **COMET** | **BLEU** | **COMET** | **BLEU** |
| ChatGPT | 72.4 | 36.5 | 68.3 | 41.4 |
| -w zero-shot CoT | 69.3 (↓3.1) | 35.1 (↓1.4) | 59.5 (↓8.8) | 36.2 (↓5.2) |
| -w 1-shot CoT | 69.6 (↓2.8) | 37.0 (↑0.5) | 61.1 (↓7.2) | 37.6 (↓3.8) |

Table 8: Performance of ChatGPT equipped with CoT prompting methods on English⇒German and English⇒Chinese.

The CoT prompt leads to **word-by-word translation** behavior, which is the main reason for the significant translation degradation.
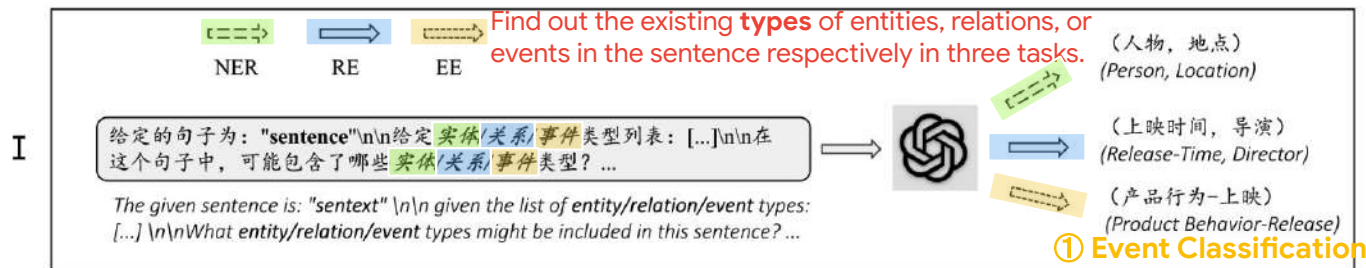
# Information Extraction

- Zero-Shot Information Extraction via Chatting with ChatGPT 2023-02-20
- Exploring the Feasibility of ChatGPT for Event Extraction 2023-03-07
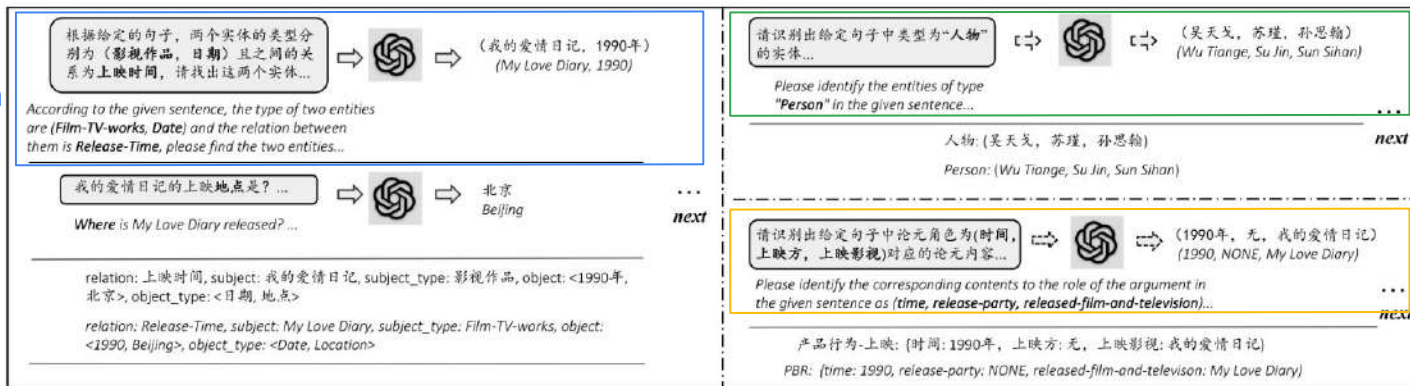
43

# Information Extraction: ChatIE



sentence: 《我的爱情日记》是1990年在北京上映的中国剧情电视剧，由吴天戈执导，苏瑾、孙思翰等主演

*"My Love Diary" is a Chinese TV series released in Beijing in 1990, directed by Wu Tiange and starred by Su Jin and Sun Sihan.*

**NER**  **RE**  **EE**

Find out the existing **types** of entities, relations, or events in the sentence respectively in three tasks.

给定的句子为："sentence"\n\n给定 *实体/关系/事件* 类型列表：[...]\n\n在这个句子中，可能包含了哪些 *实体/关系/事件* 类型？ ...

*The given sentence is: "sentext" \n\n given the list of **entity/relation/event** types: [...] \n\nWhat entity/relation/event types might be included in this sentence? ...*

（人物，地点）
*(Person, Location)*

（上映时间，导演）
*(Release-Time, Director)*

（产品行为–上映）
*(Product Behavior-Release)*

**① Event Classification**

**Entity-Relation Triple Extraction**

根据给定的句子，两个实体的类型分别为（影视作品，日期）且之间的关系为上映时间，请找出这两个实体... 

*According to the given sentence, the type of two entities are (Film-TV-works, Date) and the relation between them is Release-Time, please find the two entities...*

〈我的爱情日记，1990年〉
*(My Love Diary, 1990)*

我的爱情日记的上映地点是？ ...
*Where is My Love Diary released? ...*

北京
*Beijing*

... *next*

relation: 上映时间, subject: 我的爱情日记, subject_type: 影视作品, object: <1990年, 北京>, object_type: <日期, 地点>

*relation: Release-Time, subject: My Love Diary, subject_type: Film-TV-works, object: <1990, Beijing>, object_type: <Date, Location>*

请识别出给定句子中类型为"人物"的实体...
*Please identify the entities of type "Person" in the given sentence...*

〈吴天戈，苏瑾，孙思翰〉
*(Wu Tiange, Su Jin, Sun Sihan)*

**NER:** each turn aims to extract the entities of one type.

人物: (吴天戈, 苏瑾, 孙思翰)
*Person: (Wu Tiange, Su Jin, Sun Sihan)*

... *next*

请识别出给定句子中论元角色为(时间, 上映方, 上映影视)对应的论元内容...
*Please identify the corresponding contents to the role of the argument in the given sentence as (time, release-party, released-film-and-television)...*

(1990年, 无, 我的爱情日记)
*(1990, NONE, My Love Diary)*

**②
Argument extraction**

... *next*

产品行为–上映: (时间: 1990年, 上映方: 无, 上映影视: 我的爱情日记)
*PBR: (time: 1990, release-party: NONE, released-film-and-televison: My Love Diary)*

# Information Extraction: Results

| | RE | | | | | | NER | | | | | | EE | | | | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | DuIE2.0 | | | NYT11-HRL | | | MSRA | | | collnpp | | | DuEE1.0 | | | ACE05 | | |
| | P | R | F1 | P | R | F1 | P | R | F1 | P | R | F1 | P | R | F1 | P | R | F1 |
| fs-1 | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 | 14.7 | 7.9 | 9.7 | 2.71 | 17.2 | 4.66 | 0.4 | 0.2 | 0.3 | 0.0 | 0.0 | 0.0 |
| fs-5 | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 | 34.5 | 10.3 | 15.5 | 2.53 | 16.65 | 4.38 | 0.2 | 0.6 | 0.3 | 0.0 | 0.0 | 0.0 |
| fs-10 | 16.5 | 0.1 | 0.2 | 0.0 | 0.0 | 0.0 | 60.0 | 30.9 | 40.6 | 2.49 | 18.54 | 4.38 | 2.1 | 0.7 | 1.0 | 0.0 | 0.0 | 0.0 |
| fs-20 | 41.4 | 0.4 | 0.8 | 3.4 | 2.7 | 0.5 | 63.4 | 44.8 | 52.5 | 2.48 | 19.36 | 4.41 | 1.7 | 0.8 | 1.1 | 4.6 | 0.1 | 0.2 |
| fs-50 | 45.7 | 2.5 | 4.7 | 11.7 | 1.9 | 3.3 | 71.6 | 62.4 | 66.6 | 41.94 | 11.55 | 8.93 | 3.2 | 8.5 | 4.6 | 6.7 | 1.6 | 2.6 |
| fs-100 | 50.8 | 7.2 | 12.0 | 34.8 | 6.2 | 10.6 | 81.3 | 76.1 | 78.6 | 50.26 | 24.97 | 32.89 | 8.7 | 12.0 | 10.1 | 8.0 | 4.9 | 6.0 |
| full-shot | 68.9 | 72.2 | 70.5 | 47.9 | 55.1 | 51.3 | 96.33 | 95.63 | 95.98 | 94.18 | 94.61 | 94.39 | 50.9 | 42.8 | 46.5 | 45.3 | 54.3 | 49.4 |
| FCM | - | - | - | 43.2 | 29.4 | 35.0 | - | - | - | - | - | - | - | - | - | - | - | - |
| MultiR | - | - | - | 32.8 | 30.6 | 31.7 | - | - | - | - | - | - | - | - | - | - | - | - |
| **single** | 17.8 | 7.7 | 10.7 | 10.8 | 5.7 | 7.4 | 56.3 | **57.3** | 56.8 | 61.4 | 43.0 | 50.6 | 61.7 | 77.5 | 68.7 | 18.2 | 23.9 | 20.7 |
| **ChatIE** | **74.6** | **67.5** | **70.9** | **30.6** | **48.4** | **37.5** | **58.4** | 57.0 | **57.7** | **62.3** | **55.0** | **58.4** | **66.5** | **78.5** | **72.0** | **25.3** | **35.5** | **29.5** |

Table 1: F1 score on six datasets over two languages.

# Event Extraction

## Task Instruction

**Task Description**

This is an event extraction task where the goal is to extract structured events from the text. A structured event contains an event trigger word and an event type.

**Event types and their definitions**

**Contact.Meet**: A MEET Event occurs whenever two or more Entities come together at a single location and interact with one another face-to-face.
**Life.Die**: A DIE Event occurs whenever the life of a PERSON Entity ends.
…

**Positive Example**

**Input**: In neutral Switzerland , Foreign Minister Micheline Calmy - Rey said the United States ' decision to proceed with military action weakened the United Nations .
**Output**: [{'trigger': 'action', 'event_type': 'Conflict.Attack'}]

**Negative Example**

**Input**: Inspectors also returned to a former helicopter airfield where Iraq buried bombs it says were armed with biological weapons in 1991 .
**Output**: [{'trigger': 'returned', 'event_type': 'Life.Marry'}]

**ChatGPT**

## Evaluation Instance

**Input**: Erdogan , a leader of Turkey 's pro - Islamic movement when he was jailed , said he moderated his policies in prison .
**Expected Output**: [{'trigger': 'jailed', 'event_type': 'Justice.Arrest-Jail'}]

# Event Extraction: Results

- Dataset: ACE 2005 corpus
- We randomly select 20 samples from the raw test set to evaluate the efficacy of ChatGPT.

|  | Event Detection | | |
|---|---|---|---|
|  | **P** | **R** | **F1** |
| ChatGPT | 57.14 | 72.73 | 64.00 |
| Text2Event (T5-base) | 75.76 | 75.76 | 75.76 |
| Text2Event (T5-large) | **82.76** | 72.73 | 77.42 |
| EEQA (BERT-base) | 81.82 | **81.82** | **81.82** |

Table 1: Comparison with task-specific models.

|  | Event Detection | | |
|---|---|---|---|
|  | **P** | **R** | **F1** |
| Full Prompt | 57.14 | 72.73 | 64.00 |
| - Positive Example | 48.00 | 72.73 | 57.83 |
| - Negative Example | 62.50 | 75.76 | 68.49 |
| - Event Type Definition | 47.50 | 57.58 | 52.05 |

Table 2: Instructing with different elements.

Our results show that ChatGPT has, on average, only **51.04% of the performance of a task-specific model** such as EEQA in long-tail and complex scenarios.

The model's performance improved after **eliminating the negative sample.**

# Event Extraction: Usability of ChatGPT



Figure 2: Performance of ChatGPT with different number of attempts by four annotators.

Recruited four professional and well-educated annotators (e.g. postgraduate student on NLP research) to evaluate ChatGPT's usability.

We randomly selected ten samples from the ACE05 test set and provided each annotator with five attempts to create a task prompt that would enable ChatGPT to extract structured events from the given text.

ChatGPT is not robust enough.

ChatGPT's performance is sensitive to different prompt styles

# Data Augmentation

- ChatAug: Leveraging ChatGPT for Text Data Augmentation
  2023-2-28

# ChatAug Framework



Fig. 1. The framework of ChatAug. a (top panel): First, we apply ChatGPT for data augmentation. We input samples of all classes into ChatGPT and prompt ChatGPT to generate samples that preserves semantic consistency with existing labelled instance. b (bottom panel): In the next step, we train a BERT-based sentence classifier on the few-shot samples and the generated data samples and evaluate the model's classification performance.

ChatGPT is applied to rephrase each input sentence into six additional sentences

# ChatGPT Failures

- A Categorical Archive of ChatGPT Failures 2023-3-6

# Failure 1: Reasoning

- Critical thinking, decision making, and problem solving are all crucial activities that rely heavily on the fundamental aspect of human intelligence known as reasoning.

- **<u>Models like ChatGPT lack a "world model"</u>**, meaning they do not possess a complete understanding of the physical and social world, or the capability to reason about the connections between concepts and entities. They can only generate text based on the patterns they have learned during training.

# Failure 1.1: Spatial Reasoning

**Spatial reasoning** refers to the ability to understand and manipulate the relationships between objects, people, and places in the physical space around us.



Figure 1: The left task deals with spatial navigation and has been created by Tomer Ullman from here. The right task is a basic question regarding seat arrangement from here.

ChatGPT does possess some level of spatial understanding, as evidenced by its ability to translate the relative positions of grid boxes into language.

# Failure 1.2: Physical reasoning

**Physical reasoning** refers to the ability to understand and manipulate physical objects and their interactions in the real world.

It involves the application of physical laws and concepts to predict and explain the behavior of physical systems.

The **trophy** didn't fit in the **suitcase** because it was too small.

What was too small?

❌ Older version                ✅ ChatGPT (Jan 30, 2023)

# Failure 1.3: Temporal reasoning

**Temporal reasoning** is the ability to reason about and make predictions about events and their ordering in time.

It involves understanding the <u>temporal relationships between events</u>, <u>the duration of events,</u> and <u>the timing of events relative to each other</u>.

# Failure 1.4: Psychological reasoning

**Psychological reasoning** refers to the ability to understand and make predictions about human behavior and mental processes (a.k.a Theory of Mind).

It involves the application of psychological theories, models, and concepts to explain and predict human behavior and mental states.



Figure 4: An example failure of ChatGPT in psychological reasoning, from here.

# Failure 2: Logic

**Logic** is a branch of mathematics and philosophy that studies the principles of reasoning. It deals with the rules and methods for correct reasoning, such as syllogisms, induction, and deduction.

Mike's mum had 4 kids; 3 of them are Luis, Drake, and Matilda. What is the name of the 4th kid?

↓

ChatGPT

↓

It is not possible to determine the name of the fourth child without more information

# Failure 3: Math and Arithmetic

**Arithmetic reasoning** refers to the capability of utilizing mathematical concepts and logic to solve arithmetic problems.

It requires logical thinking and the application of mathematical principles to find the right solution to mathematical problems.

ChatGPT is limited in its capability to calculate mathematical expressions. Like most large language models, it struggles with tasks such as multiplying large numbers, finding roots, computing powers (especially with fractions), and adding or subtracting from irrational numbers (e.g. pi or e)

G   In 2015 my mother was double my age plus ten. In 2021 I was 26, how old was she in 2017?

In 2015, my mother was 26 + 10 = <<26+10=36>>36 years old.

In 2021, my mother was 26 + 10 + 6 = <<26+10+6=42>>42 years old.

In 2017, my mother was 42 - 4 = <<42-4=38>>38 years old. Answer: \boxed{38}.

*Please refer to the paper for more examples*

# Failure 4: Factual Errors

**Factual errors** refer to inaccuracies in information or statements that are not in accordance with reality or the truth.

Factual errors are often unintentional but can result in incorrect or misleading information.

- ChatGPT's output lacks accuracy in regards to scientific facts.
- It sometimes lacks knowledge of basic facts, which can be quickly obtained through a Google search.

---

**NG** Why is CPU computation faster than GPU computation for deep learning?

CPU computation is generally faster than GPU computation for deep learning because CPUs are designed to handle a wide range of tasks and can switch between different tasks quickly and efficiently. This allows them to handle complex operations, such as data preprocessing and model training, more efficiently than GPUs, which are specifically designed for parallel processing of large amounts of data. Additionally, CPUs typically have larger cache memory and more advanced memory management techniques, which can improve the speed and efficiency of deep learning algorithms.

---

**NG** Why is DNA computing faster than CPU computation for deep learning?

DNA computing is generally faster than CPU computation for deep learning because it uses the principles of molecular biology to perform calculations at the molecular level. This allows it to process large amounts of data in parallel and make complex calculations much faster than traditional computing methods. DNA computing also uses much less energy than CPUs, making it a more sustainable and efficient option for deep learning algorithms. Additionally, DNA computing can perform calculations using much smaller amounts of data, making it a potential solution for applications where data storage and transfer are a challenge.

# Failure 5: Bias and Discrimination

**Bias** in a language model refers to the systematic inaccuracies or stereotypes in the generated language output, which are influenced by the training data and reflect the societal and cultural prejudices that exist in that data.
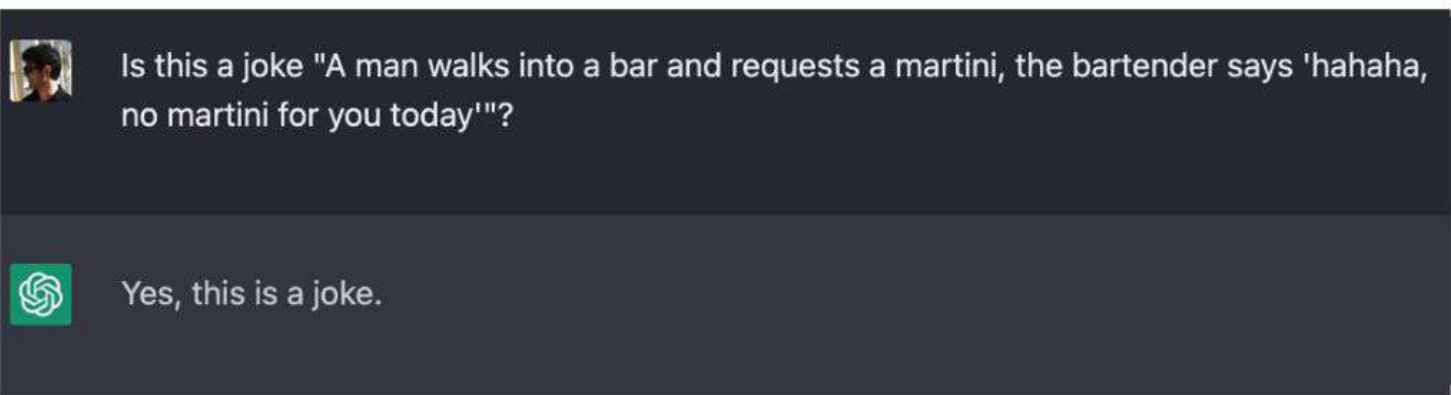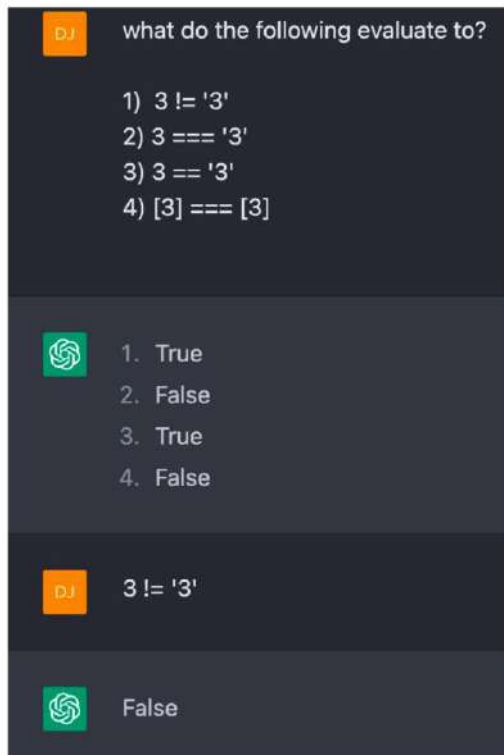
# Failure 6: Wit and Humor

**Humor** is the quality of being amusing or comical, often expressed through words or actions that entertain or make someone laugh.
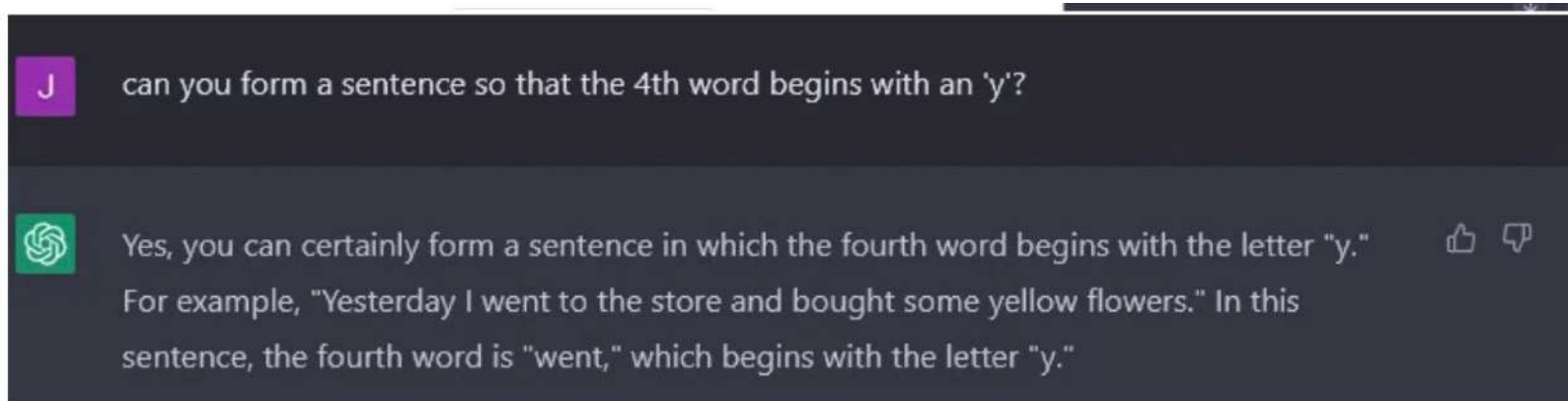


Is this a joke "A man walks into a bar and requests a martini, the bartender says 'hahaha, no martini for you today'"?

Yes, this is a joke.

# Failure 7: Coding

# Failure 8: Syntactic Structure, Spelling, and Grammar

**Syntactic structure** refers to the arrangement of words, phrases, and clauses in a sentence to form a well-defined and meaningful structure according to the rules of a particular language.

# Failure 9: Self Awareness

**Self-awareness** is the capacity to recognize oneself as an individual separate from others and to have an understanding of one's own thoughts, feelings, personality, and identity.

Self-awareness is considered an important aspect of consciousness and is closely related to self-consciousness and introspection

# Failure 10: Other Failures

1.  ChatGPT's difficulty in using idioms
2.  ChatGPT lacks real emotions and thoughts
3.  ChatGPT condenses the subject matter, but does not provide a distinctive perspective on it.
4.  ChatGPT tends to be excessively comprehensive and verbose
5.  ChatGPT lacks human-like divergences and tends to be overly literal, leading to misses in some cases
6.  ChatGPT strives to maintain a neutral stance
7.  ChatGPT's responses tend to be formal in nature due to its programming to avoid informal language.
8.  If ChatGPT is informed that its answer is incorrect, it may respond by apologizing, acknowledging its potential inaccuracies or confusion, correcting its answer, or maintaining its original response.

# Conclusion

# Takeaways

- ChatGPT updates frequently, all the current conclusions may be changed after updates.

- Due to the lack of API early on, small tests are not enough to reveal the underlying pattern.

- Prompt engineering is of vital importance.
    - Different prompts or instructions can lead to contrasting results.
    - ChatGPT is highly sensitive to different prompt styles.

- ChatGPT is super strong, but we NLPers don't need to be afraid of it.
    - Full-shot fine-tuned models can still get better results.

# References

[1]    ChatGPT: A Meta-Analysis after 2.5 Months

[2]    Exploring the Limits of ChatGPT for Query or Aspect-based Text Summarization

[3]    Cross-Lingual Summarization via ChatGPT

[4]    Is ChatGPT a Good NLG Evaluator? A Preliminary Study

[5]    Is ChatGPT A Good Translator? A Preliminary Study

[6]    Towards Making the Most of ChatGPT for Machine Translation

[7]    Zero-Shot Information Extraction via Chatting with ChatGPT

[8]    Exploring the Feasibility of ChatGPT for Event Extraction

[9]    ChatAug: Leveraging ChatGPT for Text Data Augmentation

[10]   A Categorical Archive of ChatGPT Failures

[11]   [TBD] A Multitask, Multilingual, Multimodal Evaluation of ChatGPT on Reasoning, Hallucination, and Interactivity

[12]   [TBD] Is ChatGPT a General-Purpose Natural Language Processing Task Solver?

# Thank you!