

A Multitask, Multilingual, Multimodal Evaluation of ChatGPT on Reasoning, Hallucination, and Interactivity

Yejin Bang Samuel Cahyawijaya Nayeon Lee Wenliang Dai Dan Su Bryan Wilie
Holy Lovenia Ziwei Ji Tiezheng Yu Willy Chung Quyet V. Do Yan Xu Pascale Fung
Centre for Artificial Intelligence Research (CAiRE)
The Hong Kong University of Science and Technology
yjbang@connect.ust.hk, pascale@ece.ust.hk

Abstract

This paper proposes a framework for quantitatively evaluating interactive LLMs such as ChatGPT using publicly available data sets. We carry out an extensive technical evaluation of ChatGPT using 23 data sets covering 8 different common NLP application tasks. We evaluate the multitask, multilingual and multimodal aspects of ChatGPT based on these data sets and a newly designed multimodal dataset. We find that ChatGPT outperforms LLMs with zero-shot learning on most tasks and even outperforms fine-tuned models on some tasks. We find that it is better at understanding non-Latin script languages than generating them. It is able to generate multimodal content from textual prompts, via an intermediate code generation step. Moreover, we find that ChatGPT is 63.41% accurate on average in 10 different reasoning categories under logical reasoning, non-textual reasoning, and commonsense reasoning, hence making it an unreliable reasoner. It is, for example, better at deductive than inductive reasoning. ChatGPT suffers from hallucination problems like other LLMs and it generates more extrinsic hallucinations from its parametric memory as it does not have access to an external knowledge base. Finally, the interactive feature of ChatGPT enables human collaboration with the underlying LLM to improve its performance, i.e., **8% ROUGE-1** on summarization and **2% ChrF++** on machine translation, in a multi-turn "prompt engineering" fashion. We also release codebase for evaluation set extraction.¹

1 Introduction

ChatGPT is a successor of the large language model (LLM) InstructGPT (Ouyang et al., 2022) with a dialog interface that is fine-tuned using the Reinforcement Learning with Human Feedback

(RLHF) (Christiano et al., 2017) approach.² In the last couple of months, ChatGPT has gathered close to 1 million user base (Hu, 2023) and is being used by businesses and consumers alike for a myriad of mostly textual tasks. One reason for its unprecedented popularity is that ChatGPT, through its scale and via RLHF, has shown impressive abilities in many areas of NLP as well as emergent abilities such as code generation and multimodal generation. Another reason is that its dialog interface allows users to interact with the underlying large language model more effectively and efficiently via interactive chats that are akin to multi-turn prompt engineering.

However, despite its powerful abilities, anecdotal reports on ChatGPT have consistently shown significant remaining challenges - for example, it fails in some elementary mathematical (Gilson et al., 2022; Goldberg, 2023; Frieder et al., 2023; Choi et al., 2023; Davis, 2023) and commonsense reasoning tasks (Guo et al., 2023; Davis, 2023); it hallucinates with human-like fluency and eloquence on things that are not based on truth (Shen et al., 2023; Thorp, 2023; Smith, 2023); and as a general-purpose language model trained from everything on the web, its language coverage is questionable (Lu et al., 2022; Jiao et al., 2023). OpenAI has listed many limitations of ChatGPT on its website.³ CEO tweeted that "It's a mistake to be relying on [ChatGPT] for anything important right now" (Altman, 2022). Many researchers have argued that, despite appearances, LLMs like ChatGPT are only good at language abilities, not actual reasoning (Mahowald et al., 2023).

Consequently, it is not clear what people can or cannot use it for despite its popularity. For users and researchers alike, it would be beneficial to have

¹<https://github.com/HLTCHKUST/chatgpt-evaluation>

²<https://beta.openai.com/docs/model-index-for-researchers>

³<https://platform.openai.com/docs/chatgpt-education>

a sense of confidence in its reliability in various NLP/AI tasks.

Previous works have discussed the ethical implications or concerns associated with ChatGPT (and other LLMs) (Jabotinsky and Sarel, 2022; Susnjak, 2022; Blanco-Gonzalez et al., 2022; Aydın and Karaarslan, 2022; Jeblick et al., 2022). However, there has not been much technical evaluation of the strengths and limitations of ChatGPT⁴. To fill this gap, we conduct experiments on ChatGPT with samples from standard public test sets on major NLP tasks such as question answering, reasoning, summarization, machine translation, automatic post-editing, sentiment analysis, language identification, and task-oriented dialogue (dialogue state tracking & response generation) and misinformation detection. We evaluate its multilingual performance as well as vision-language multimodal abilities. With additional experiments, we also quantitatively evaluate its primary limitations in *reasoning* and *hallucination*. In addition, we conduct experiments to test its *multi-turn interactivity* as a means for better prompt engineering. We hope to provide insights to users of ChatGPT on the above-mentioned strengths and limitations, as well as how they can improve outcomes with interactivity. (Note that we are not able to quantitatively evaluate the RLHF aspect of ChatGPT without access to the user log. We hope OpenAI will publish this work and one can carry out such evaluations in the future in collaboration with OpenAI.)

The following are the major insights we have gained from the evaluations:

Multitask, Multimodal, and Multilingual

- For 9/13 NLP datasets, ChatGPT outperforms previous LLMs with zero-shot learning. It even outperforms fully fine-tuned task-specific LMs on 4 different tasks. In other cases, ChatGPT is on par or slightly lower than fully fine-tuned for specific NLP tasks;
- ChatGPT fails to generalize to low-resource and extremely low-resource languages (e.g., Marathi, Sundanese, and Buginese). There is an overall performance degradation in low-resource languages, especially in non-Latin scripts in the case of translation; its weakness lies in generation rather than understanding part of the translation process;

⁴Many anecdotal analyses have been posted online, but none in a comprehensive manner

- ChatGPT enables a code intermediate medium to bridge vision and language, even though the multi-modality ability is still elementary compared to vision-language models.

Reasoning We tested 10 different reasoning categories with 634 samples in total. Based on our experiments, ChatGPT shows more weakness in inductive reasoning than in deductive or abductive reasoning. ChatGPT also lacks spatial reasoning while showing better temporal reasoning. ChatGPT also lacks mathematical reasoning, which aligns with recent findings by Frieder et al.. Further, we found that ChatGPT is relatively better at common-sense reasoning than non-textual semantic reasoning. Finally, while ChatGPT shows acceptable performance in causal and analogical reasoning, it is bad at multi-hop reasoning capability as similar to other LLMs’ weakness in complex reasoning (Ott et al., 2023).

Hallucination Similar to other LLMs (Radford et al., 2019; Muennighoff et al., 2022; Workshop et al., 2022), ChatGPT suffers from the hallucination problem. It generates more extrinsic hallucinations – factual statements that cannot be verified from the source, from its parametric memory across all tasks since it does not possess the access to external knowledge bases.

Interactivity One of the primary differentiating factors of ChatGPT from its predecessors is its *multi-turn dialog interactivity*. This enables ChatGPT to perform multiple tasks within a dialog session. There is also significant performance improvement (8% ROUGE-1 on summarization and 2% ChrF++ on low-resource machine translation) via multi-turn interactivity in various standard NLP tasks. This process is akin to prompt engineering with feedback from the system.

Organization of This Paper: We first provide an overview of ChatGPT and related work (§2). Then, we provide evaluation results on ChatGPT on various application test sets, on multilingual test sets, and on a new multimodal task in §3. We then explore the three main strengths and weaknesses of ChatGPT, namely *reasoning* (§4), *hallucination* (§5) and *interactivity* (§6) in the subsequent three sections. Finally, we discuss and give a conclusion on our findings of ChatGPT.

2 Background and Related Work

2.1 Large Pretrained Models

Large Language Models (LLMs) are language models with parameter sizes over a hundred billion, beginning with the introduction of GPT-3. Examples of LLMs include, but are not limited to, GPT-3, Gopher (Rae et al., 2021b), Megatron (Shoeybi et al., 2019), GPT-Jurassic (Lieber et al., 2021), OPT-175B Zhang et al. (2022). Beyond fine-tuning models with task-specific data, LLMs have shown robustness and generalizability through zero-shot and few-shot learning with examples. Scaling up the models unlocked new, emergent abilities that were not observed with smaller models (Wei et al., 2022a). Prompts are used to probe the LLMs to generate the target outcome by sampling the language distribution. To enable the LLMs to demonstrate their abilities, sophisticated prompt engineering (NeuralMagic, 2023) is required. However, previous LLMs only allow one-time probing, which means the target outcome varies a great deal with minor changes in the prompt instruction.

Whereas scaling up LLMs improve generalizability, generic LLMs may fall short in specific applications. Despite its name, ChatGPT has not been primarily used as a chatbot. Its dialog ability serves as the user interface to the underlying LLM. We nevertheless refer to other dialog systems here in this paper. A number of large pre-trained dialogue models have been created, following the pre-train-then-finetune paradigm. LaMDA (Thoppilan et al., 2022) is a large-scale conversational model, fine-tuned from an LLM with a parameter size of 134 billion. Blenderbot 3.0 (Shuster et al., 2022), scaled up to 175 billion parameter size, is also introduced with similar abilities as LaMDA. Both models are pre-trained on public dialogue and other public web documents and then fine-tuned with manually curated dialogue data. They also have access to external knowledge sources for information retrieval, thus they have shown an excellent ability for fluent and natural dialogue generation as well as information retrieval. However, the aforementioned large dialogue models suffer from catastrophic forgetting of the knowledge obtained from the pre-training. Models after fine-tuning show stable and strong performance on specific tasks, but they only preserve the knowledge learned from the task-specific data while losing the generalization ability. ChatGPT, on the other hand, was trained on a large-scale conversational-style dataset constructed from web

documents directly (Schulman et al., 2022), which unifies the pre-training and fine-tuning data format. Thus, ChatGPT is able to preserve the knowledge from pre-training and produce informative outputs without access to external knowledge sources.

2.2 ChatGPT

Compared to existing LLMs, ChatGPT has unique characteristics. First, it has the ability to interact with users in a conversation-like manner, while retaining its accumulated knowledge and generalization ability gained from pre-training. This is achieved by pre-training ChatGPT on a large-scale conversational-style dataset, that is constructed by transforming a large-scale instruction-tuning corpus used for building InstructGPT into a conversational format, then fine-tuning the model based on a reward model to further improve the generation quality and align the generation with human preference. ChatGPT should be considered a generic language model which can be probed in a conversational manner. The biggest advantage of such conversational interaction is that, unlike previous LLMs, ChatGPT can intelligently “answer follow-up questions, admit its mistakes, challenge incorrect premises, and reject inappropriate requests” (Schulman et al., 2022).

Second, ChatGPT is trained with a better human-aligned objective function via Reinforcement Learning from Human Feedback (RLHF) (Christiano et al., 2017). Conventional natural language generation models, including dialogue models, are trained with maximum likelihood estimation (MLE) and might not be aligned with human preferences. For instance, for dialogue systems, humanness, engagement, and groundedness are some examples of essential criteria for success. Such discrepancy between training objectives and evaluation metrics becomes a bottleneck to performance improvement. By using RLHF, ChatGPT aligns more closely with human preferences in generating text than by using MLE.

As ChatGPT has become available to public users through an easily accessible UI, there have been many discussions from a wide range of communities, not just from AI or NLP, but also from other disciplines. A line of discussion is the specific emergent ability and strength of ChatGPT in more technical perspectives. Guo et al. (2023) conducts linguistic analyses and human evaluations of ChatGPT’s writing against human experts with their

proposed corpus named Human ChatGPT Comparison Corpus and found that ChatGPT responses are strictly focused on the given question, more formal, objective, and less emotional. Nov et al. (2023) also studies ChatGPT’s generated medical advice if it passes the Turing test. Frieder et al. (2023) investigate mathematical capabilities of ChatGPT on both publicly available and hand-crafted datasets, including graduate-level mathematics, and show that “significantly below those of an average mathematics graduate student.” There are many investigations of ChatGPT’s understanding and potential applications in different fields such as law (Choi et al., 2023), medical domain (Blanco-Gonzalez et al., 2022; Jeblick et al., 2022) and finance (Birch, 2022; Dowling and Lucey, 2023). Jeblick et al. (2022) conduct a case study of the application of ChatGPT on simplified radiology reports. Another important line of discussion is the ethical concerns over the use of ChatGPT. The most active discussion is over the use of academic writing and exam integrity (Jabotinsky and Sarel, 2022; Susnjak, 2022). OpenAI also discusses the misuse of LM for disinformation and remedies.⁵ Zhuo et al. study AI ethics of ChatGPT in criteria of bias, reliability, robustness, and toxicity.

2.3 LLM benchmark and evaluation

With the advancement of LLMs’ generalization ability, there have been efforts to understand their capabilities, limitations, and risks. Recently, several benchmarks with a collection of a large number of NLP datasets, such as BIG-Bench (Srivastava et al., 2022) and AI LM Harness (Gao et al., 2021), have been introduced. Moreover, HELM (Liang et al., 2022) is proposed to conduct a holistic evaluation of LLMs that considers scenarios and metrics with a top-down approach. In this work, we instead focus on specific limitations and unique findings of ChatGPT that had not been discussed with previous LLMs. There is difficulty to evaluate ChatGPT with the whole test set from such benchmarks due to limited access to ChatGPT⁶.

There are also other works that discuss LLMs’ emergent abilities through thorough surveys or case studies. Mahowald et al. (2023) thoroughly studies LLMs capabilities by distinguishing *formal* and *functional* linguistic competence with reference to cognitive science, psychology, and NLP to clarify

⁵<https://openai.com/blog/forecasting-misuse/>

⁶As of the end of January 2023, there is no official API provided by Open AI.

the discourse surrounding LLMs’ potential. Other works focus on more specific abilities such as mathematical skills (Davis, 2023), reasoning (Webb et al., 2022a; Qiao et al., 2022). Also, there have been overviews of existing LLMs (Gozalo-Brizuela and Garrido-Merchan, 2023; Wolfe, 2023)

3 Multitask, Multilingual, and Multimodal Evaluations of ChatGPT

3.1 Evaluating the Multitask Ability of ChatGPT

ChatGPT has become very well-known in such a short period of time to general public users, not just those who are in AI, machine learning, and NLP communities who might be more familiar with LLMs. One of the main reasons is that, in addition to media reports, innumerable use cases of ChatGPT are shared by both non-academic and academic users online (Marr, 2022; Gordon, 2023; Shankland, 2023). There have been debates and panels on whether ChatGPT is approaching Artificial General Intelligence (AGI), as it seems to be able to carry out a multitude of tasks without specific fine-tuning (Desk, 2023; Johnson, 2023; Kingson, 2023). On the other hand, there has also been as much sharing of its failures in simple tasks (Gilson et al., 2022; Choi et al., 2023; Shen et al., 2023).

Instead of relying on anecdotal examples, we first evaluate ChatGPT’s performance in various standard NLP tasks in a zero-shot manner to obtain a basic/better understanding of its multi-task ability. We compile results from the existing literature on ChatGPT and compare them with the state-of-the-art fully-fine-tuned and zero-shot models across multiple tasks. We evaluate ChatGPT performances on 21 datasets covering 8 tasks, i.e., summarization, machine translation, sentiment analysis, questions answering, task-oriented dialogue, open-domain knowledge-grounded dialogue, and misinformation detection tasks. For ChatGPT, we sample testing cases from existing standard test sets for each task with a sample size ranging from 30 to 200 samples per task.

Multitask Generalization of ChatGPT The result of the multitask evaluation is shown in Table 1. ChatGPT is shown to achieve remarkable zero-shot performances on multiple tasks, surpassing pre-

⁷We take the average from the state-of-the-art zero-shot performance in CNN and DM from Goyal et al. (2022).

Tasks	Dataset	Metric	Reference	Fine-Tuned SOTA	Zero-Shot SOTA	ChatGPT
Summarization	CNN/DM	ROUGE-1	Lewis et al. (2020a)	44.47	35.27 ⁷	35.29
	SAMSum	ROUGE-1	Lewis et al. (2020a)	47.28	-	35.29
MT (XXX→Eng)	FLoRes-200 (HRL)	ChrF++	Team et al. (2022)	63.5	-	58.64
	FLoRes-200 (LRL)	ChrF++	Team et al. (2022)	54.9	-	27.75
MT (Eng→XXX)	FLoRes-200 (HRL)	ChrF++	Team et al. (2022)	54.4	-	51.12
	FLoRes-200 (LRL)	ChrF++	Team et al. (2022)	41.9	-	21.57
Sentiment Analysis	NusaX - Eng	Macro F1	Winata et al. (2022)	92.6	61.5	83.24
	NusaX - Ind	Macro F1	Winata et al. (2022)	91.6	59.3	82.13
	NusaX - Jav	Macro F1	Winata et al. (2022)	84.2	55.7	79.64
	NusaX - Bug	Macro F1	Winata et al. (2022)	70.0	55.9	55.84
Question Answering	bAbI task 15	Accuracy	Weston et al. (2016a)	100	-	93.3
	bAbI task 16	Accuracy	Weston et al. (2016a)	100	-	66.7
	EntailmentBank	Accuracy	Clark et al. (2018)	86.5	78.58	93.3
	CLUTRR	Accuracy	Minervini et al. (2020)	95.0	28.6	43.3
	StepGame (k=9)	Accuracy	Mirzaee and Kordjamshidi (2022)	48.4	-	23.3
	StepGame (k=1)	Accuracy	Mirzaee and Kordjamshidi (2022)	98.7	-	63.3
	Pep-3k	AUC	Porada et al. (2021)	67.0	-	93.3
Misinformation Detection	COVID-Social	Accuracy	Lee et al. (2021)	77.7	50.0	73.3
	COVID-Scientific	Accuracy	Lee et al. (2021)	74.7	71.1	92.0
Task-Oriented Dialogue	MultiWOZ2.2	JGA	Zhao et al. (2022)	60.6	46.7	24.4
	MultiWOZ2.2	BLEU	Nekvinda and Dušek (2021)	19.1	-	5.65
	MultiWOZ2.2	Inform Rate	Yang et al. (2021)	95.7	-	71.1
Open-Domain KGD	OpenDialKG	BLEU	Ji et al. (2022c)	20.8	3.1	4.1
	OpenDialKG	ROUGE-L	Ji et al. (2022c)	40.0	29.5	18.6
	OpenDialKG	FeQA	Ji et al. (2022c)	48.0	23.0	15.0

Table 1: Performance of ChatGPT compared to state-of-the-art fully-fine-tuned models (Fine-Tuned SOTA) and LLM in zero-shot settings (Zero-Shot SOTA). The referenced performances are evaluation results on full test sets, while the ChatGPT performances are computed on subsets of the corresponding dataset **using 30 to 200 data samples** for each task. For Machine Translation (MT) tasks, we use the definitions of high-resource language (HRL) and low-resource language (LRL) from NLLB (Team et al., 2022) and take subsets of languages to represent each group. JGA denotes joint goal accuracy.

vious state-of-the-art zero-shot models on 9 out of 13 evaluation datasets with reported zero-shot LLMs performance. In most tasks, especially task-oriented and knowledge-grounded dialogue tasks, task-specific fully-fine-tuned models outperform ChatGPT. Compared to the latter, ChatGPT yields lower performance in most tasks while still surpassing the performance on 4 evaluation datasets.

Furthermore, from the evaluation results, we also observe several limitations of ChatGPT, e.g., 1) limited language understanding and generation capabilities on low-resource languages, 2) lacking reasoning ability as shown from the results in QA, and 3) performing task-oriented and knowledge-grounded dialogue tasks. More detailed experimental setup and analysis for each task are shared in the next subsections, i.e., §3.1.1: Experiment details and result and §3.1.2: ChatGPT on Dialogue System. We also provide the complete list of all the datasets used in our evaluation in Appendix C.

3.1.1 ChatGPT on Summarization, MT, Sentiment Analysis, QA, and Misinformation Detection

Summarization We test on 100 samples from two common summarization datasets: half from SAMSum (Gliwa et al., 2019), a dialogue summarization dataset, and another half from CNN/DM (Hermann et al., 2015; Nallapati et al., 2016), news summarization datasets. The large version of Bart (Lewis et al., 2020b) model fine-tuned on both datasets is conducted for comparison. Moreover, OpenAI’s text-davinci-002 is used as the previous SOTA zero-shot model. We calculate ROUGE-1 scores for evaluating the generated summary. As is shown in Table 1, ChatGPT achieves a similar zero-shot performance with text-davinci-002, which is expected since they evolved from the same GPT3 pre-trained checkpoint. However, the fine-tuned Bart still outperforms zero-shot ChatGPT by a large margin. Furthermore, we evaluate the ChatGPT’s unique

interaction capabilities in §6.

Machine Translation We evaluate the machine translation ability of ChatGPT on both high-resource and low-resource languages using the ChrF++ metric (Popović, 2015). Specifically, we incorporate 8 high-resource languages, i.e., French (fra), Spanish (spa), Chinese (zho), Arabic (ara), Japanese (jpn), Indonesian (ind), Korean (kor), and Vietnamese (vie), and 4 low-resource languages, i.e., Javanese (jav), Sundanese (sun), Marathi (mar), and Buginese (bug) for our evaluation.⁸ For each language pair, we sample 30 Eng↔XXX parallel sentences from the FLORES-200 dataset (Team et al., 2022; Goyal et al., 2021). The result of our experiment suggests that ChatGPT can well perform XXX→Eng translation, but it still lacks the ability to perform Eng→XXX translation.

Sentiment Analysis Sentiment analysis has been widely explored for both high-resource and low-resource languages (Wang et al., 2018a; Wilie et al., 2020; Imania et al., 2018). We explore the sentiment analysis ability of ChatGPT through 4 languages with diverse amounts of resources in NusaX (Winata et al., 2022): English (eng), Indonesian (ind), Javanese (jav), and Buginese (bug). For each language, we sample 50 sentences from the corresponding dataset for our experiment and measure the macro F1 score as the evaluation metric. We compare the results with two baselines, i.e., supervised state-of-the-art performance from Winata et al. (2022) and zero-shot multilingual LLM from Cahyawijaya et al. (2022). ChatGPT outperforms the previous state-of-the-art zero-shot model by a large margin except for the Buginese, where it performs on par. This shows that ChatGPT still has a limited understanding of extremely low-resource languages.

Question Answering Since Question Answering (QA) is a broad topic, we classify QA datasets into different categories based on the knowledge/reasoning type required to do the task, e.g commonsense reasoning, spatial reasoning, temporal reasoning, etc., to have a clearer analysis on ChatGPT’s abilities. For each category, we select several datasets, and for each dataset, we sample 30 instances and test ChatGPT on the subset. Details on the dataset will be described in which subsection of 4. Furthermore, we inspect the rationales provided by ChatGPT that it used to come up with

⁸For a fairer comparison in our multitask experiment, we strictly follow the definition of high-resource and low-resource languages from NLLB (Team et al., 2022).

the answers. Some of them will be discussed in detail in the corresponding section (4). Based on our experiment results, ChatGPT outperforms the existing zero-shot and some of the fine-tuned state-of-the-art performance on question answering. Furthermore, ChatGPT achieves near-perfect scores on three tasks, i.e., bAbI task 15, EntailmentBank, and Pep-3k.

Misinformation Detection We test ChatGPT’s ability to detect misinformation with the test sets that consist of scientific and social claims related to COVID-19 (Lee et al., 2021) with 100 samples. We take half from scientific (covid-scientific) and another half from social (covid-social) sets. We evaluate the accuracy of the veracity by manually checking the generated text. ChatGPT could detect misinformation 92% (46/50) and 73.33% (22/30, excluding verification-refusing cases) accuracy on covid-scientific and covid-social respectively.

3.1.2 ChatGPT on Dialogue Tasks

Given that ChatGPT has the ability to generate conversation-like responses, it is interesting to test their ability in response generation in different dialogue settings: 1) Knowledge-Grounded Open-Domain Dialogue and 2) Task-Oriented Dialogue.

Knowledge-Grounded Open-Domain Dialogue

Open-domain dialogue systems interact with humans with generated responses automatically and aim to provide users with an engaging experience. To boost informativeness, these systems leverage external knowledge, including structured knowledge such as knowledge graphs (Zhao et al., 2020; Ji et al., 2022c) and unstructured knowledge such as free text (Xu et al., 2022).

To quantitatively measure ChatGPT’s performance on knowledge-grounded dialogue, we apply it to 50 samples randomly selected from the test set of OpenDialKG (Moon et al., 2019), which contains open-ended dialogues grounded on a knowledge path. We use the following instruction for this KGD task: “Can we try dialogue generation? I will give you turns, and you can generate the next turn, but only one.\n\n You can also consider the knowledge of XXX for your reference in the dialogue.”

According to human judgment, the responses from ChatGPT are of high quality with fluent response generation as well as incorporating the provided knowledge in the response. However, the automatic evaluation results in Table 2 are rela-

Model	BLEU \uparrow	ROUGE-L \uparrow	FeQA \uparrow (Durmus et al., 2020)
ChatGPT	4.05	18.62	15.03
GPT2	11.10	30.00	26.54

Table 2: Automatic evaluation results on OpenDialKG. The results for GPT2 are from Dziri et al. (2021).

tively low compared with GPT2 (Radford et al., 2019), which is fine-tuned on this dataset. Specifically, ChatGPT obtains a 4.05 BLEU and an 18.62 ROUGE-L score as the generated responses tend to be longer than the golden answers. For FeQA, which measures the generated response’s faithfulness to the input source, ChatGPT gets 15.03 since some generated responses include content from its parametrized knowledge injected during pre-training.

Task-Oriented Dialogue In task-oriented dialogue (TOD), a model needs to fulfill a specific objective by interacting in natural language with the user. This task is often split into three modules: natural language understanding with belief state tracking, decision-making through dialogue policies, and response generation – a modular approach that handles each of these steps with different models. Besides, unified approaches are starting to show increasingly strong performances (Hosseini-Asl et al., 2020; Peng et al., 2021).

Although ChatGPT seems more appropriate for open-domain dialogue tasks, we investigate and discuss how ChatGPT’s emergent abilities and interactivity could potentially be leveraged for TOD as well. We explore two setups A) modular approach: testing both dialogue state tracking and response generation using oracle actions; B) unified approach: a direct approach to simulate the TOD interaction while leveraging information in a structured database. We provide an example of the modular and unified approaches in Appendix F.

Setup A: Modular Approach We investigate ChatGPT’s ability for both dialogue state tracking and response generation in 50 dialogue turn samples taken from MultiWOZ2.2 (Zang et al., 2020). In detail, we ask the model to provide the belief state as domain-intent: [slot1, value1], ... in the prompt following previous zero-shot (Lin et al., 2021) and few-shot (Madotto et al., 2021) approaches, and provide an exhaustive list of domain-intent-slot-value for the given dialogue. For the response generation, we provide only the oracle

State Tracking	Response Generation	
Joint Goal Acc.	BLEU	Inform rate
24.4%	5.65	71.1%

Table 3: Result for Task-oriented Dialogue Setup A – Modular Approach.

dialogue actions (e.g. ‘Hotel-Inform’:[‘area’, ‘centre’]), and ask ChatGPT to generate a TOD response given the dialogue history. We assess DST with joint goal accuracy (JGA), the ratio of dialogue turns where the predicted dialogue state is exactly the ground truth, and response generation with BLEU and inform rate(%)

As shown in table 3, the performance for DST is mediocre with a JGA of 24.4%, but a lot of the failure cases are from the model predicting additional belief states on top of the gold ones. In our setting, all belief states are correctly predicted in 73% of the samples. We postulate that the model rely too much on previous belief states since they are all provided within the prompt. For response generation, ChatGPT successfully leverages all information provided while answering the questions with an 71.1% inform rate and 5.65 BLEU score. The BLEU is computed directly on the lexicalized response as ChatGPT skips the delexicalized generation, and the generation is often as if not more natural than the gold response.

Setup B: Unified Approach We explore ChatGPT’s ability to simulate a TOD interaction in an end-to-end manner by providing nothing more than a structured database and giving the instruction “Use the following knowledge base to complete the task of recommending a restaurant as a task-oriented dialogue system”. In this setup, we could investigate whether ChatGPT is able to complete basic retrieval queries and respond to users’ requests such as “Give me some restaurants that serve Italian food” or “I would prefer cheap options please”. However, there are several limitations that we could investigate as follow.

- **Long-term Multi-turn Dependency:** ChatGPT cannot keep the belief state across multiple turns within the interaction. For instance, asking for Italian food will overwrite the previous turn’s belief state by asking for restaurants with a rating of 3 or higher. However, if the user explicitly asks to recall the earlier prefer-

Language	#Speakers	CC Size (%)	Language Category
English (eng)	1.452B	46.320	HRL
Chinese (zho)	1.118B	4.837	HRL
French (fra)	235M	4.604	HRL
Indonesian (ind)	199M	0.781	MRL
Korean (kor)	81.7M	0.679	MRL
Javanese (jav)	68.3M	0.002	LRL
Sundanese (sun)	32.4M	0.001	LRL
Buginese (bug)	-M	0.000	X-LRL

Table 4: The statistics of languages used in our language disparity experiment. **HRL** denotes high-resourced language, **MRL** denotes medium-resourced language, **LRL** denotes low-resourced language, **X-LRL** denotes extremely low-resourced language.

ences, ChatGPT is able to correct the retrieved information and incorporate the previous belief state. This is interesting as it shows that the information previously given in multi-turn is still usable, but needs to be called explicitly.

- **Basic Reasoning Failure:** ChatGPT’s response tends to be wrong if the query introduces a basic level of reasoning such as when it is asked for “recommendation for restaurants with European food” (ChatGPT has to filter the types of cuisine which are based on countries) or “recommendation for restaurants with a rating of 3 or higher” (ChatGPT needs to understand rating 3, 4 and 5). Even with a basic knowledge base, ChatGPT fails to answer correctly 66% of the time.
- **Extrinsic Hallucination:** ChatGPT tends to generate hallucinated information beyond the given knowledge. This is especially harmful in TOD as ChatGPT will sometimes hallucinate some prices for hotel booking, or availability for restaurants.

3.2 Evaluating Multilinguality of ChatGPT

Training data size affects language understanding and generation quality of LMs (Radford et al., 2019; Raffel et al., 2022; Cahyawijaya et al., 2021; Rae et al., 2021a; Workshop et al., 2022; Chowdhery et al., 2022; Hoffmann et al., 2022). As an LLM, the same premise also applies to ChatGPT, and the question is to what extent. We investigate this question through a series of experiments by analyzing 1) the language understanding capability using two different tasks, i.e. language identification (LID) and sentiment analysis,

Language	SA Acc.	LID Acc.
English	84%	100%
Indonesian	80%	100%
Javanese	78%	0%
Buginese	56%	12%

Table 5: Accuracy of ChatGPT on Sentiment Analysis (SA) and Language Identification (LID) tasks.

and 2) the language generation capability through machine translation using English as the pivot language. Based on the percentage of data in the CommonCrawl⁹, we group languages into 3 categories, i.e., **high-resource (>1%)**, **medium-resource (>0.1%)**, **low-resource (<0.1%)**. The statistics of all the languages under study are shown in Table 4.

3.2.1 Language Understanding

We propose a framework for investigating the language understanding ability of ChatGPT through 3 languages from different language categories in NusaX (Winata et al., 2022), i.e. English (eng), Indonesian (ind), Javanese (jav). In addition, we incorporate an extremely low-resource language from NusaX, i.e., Buginese (bug), which is not even listed on CommonCrawl since the LID used in CommonCrawl¹⁰, i.e., CLD2 (Ooms, 2023), does not support Buginese (bug). We sample 50 sentences per language from the corresponding dataset for our experiment.

ChatGPT fails to generalize to extremely low-resource languages As shown in Table 5, ChatGPT achieves 84%, 80%, 78%, and 56% accuracy for English, Indonesian, Javanese, and Buginese, respectively. This result supports the results in prior works focusing on LLM (Chowdhery et al., 2022; Workshop et al., 2022; Muennighoff et al., 2022), where LLMs, including ChatGPT, yield a lower performance for lower resource languages. Interestingly, the performance gap between English, Indonesian, and Javanese is considered marginal compared to the performance gap with Buginese. This result suggests that ChatGPT still has a limitation in generalizing toward extremely low-resource languages.

⁹CommonCrawl (<http://commoncrawl.org>) is the primary source of language pre-training data used in GPT3

¹⁰<https://commoncrawl.github.io/cc-crawl-statistics/plots/languages>

ChatGPT	InstructGPT	text-davinci-003
The language of the text appears to be a variant of the Bugis language spoken in Indonesia .	The language of the text is the Sasak language , spoken in Lombok , Indonesia.	The text is written in Buginese .
I am sorry, I do not recognize the language of the text.	The language of the text is Koyukon Athabascan .	The text is in the Balinese language .
The language of the text appears to be a dialect of the Indonesian language.	The language of the text is Indonesian .	The language of the text is Indonesian .

Table 6: Example of **Buginese** language identification response from ChatGPT, InstructGPT, and text-davinci-003.

Language	XXX→Eng	Eng→XXX
Chinese	24/30	14/30
French	29/30	25/30
Indonesian	28/30	19/30
Korean	22/30	12/30
Javanese	7/30	6/30
Sundanese	9/30	0/30

Table 7: Number of correct translations of ChatGPT. XXX denotes the target language in the first column. The languages are sorted based on the language size in CommonCrawl.

ChatGPT understands sentences in low-resource languages but lacks the ability to identify the language ChatGPT correctly classified the languages for English and Indonesian 100% of the time. While for the language identification for Javanese and Buginese, ChatGPT either misclassifies the samples as other languages or is unable to determine the language for 100% for Javanese and 88% for Buginese. ChatGPT misclassifies the samples mostly as Indonesian, despite having various dissimilarities across languages (Grimes, 2000; Lewis, 2009; Cohn and Ravindranath, 2014; Eberhard et al., 2021; Aji et al., 2022; Cahyawijaya et al., 2022). Nevertheless, ChatGPT performance on the sentiment analysis on Javanese is only slightly lower compared to English and Indonesian which suggests that ChatGPT can understand the semantic meaning of sentences in low-resource languages, such as Javanese, without having enough knowledge to identify the language itself.

ChatGPT displays better human-preferred responses As shown in Table 6, ChatGPT lets the user know that its prediction is uncertain when it does not completely understand the language and

also provides broader information regarding the language, such as location and tribe of which the predicted language is spoken. This fact provides evidence regarding the benefit of using the RLHF approach compared to other training approaches for aligning LLMs with human preferences.

3.2.2 Language Generation

We assess the multilingual language generation ability of ChatGPT through machine translation. ChatGPT has been shown to be competitive compared to commercial translation products for high-resource languages (Jiao et al., 2023). Specifically, we choose 2 languages from each language category, i.e., French (fra), Chinese (zho), Indonesian (ind), Korean (kor), Javanese (jav), and Sundanese (sun) from the FLORES-200 dataset (Team et al., 2022; Goyal et al., 2021). For each language, we sample 30 English-XXX parallel sentences and perform two directions of translation using English as the pivot language. The correctness of the translation results is manually validated by a native speaker of the corresponding language.

ChatGPT performs worse on low-resource languages As shown in Table 7, similar to other LLMs (Workshop et al., 2022; Muennighoff et al., 2022), ChatGPT produces better English translation quality from high-resource languages, such as French and Chinese. While for low-resource languages, such as Javanese and Sundanese, ChatGPT tends to generate several mis-translated words/phrases and sometimes even hallucinate some objects. Moreover, we also observe that sometimes ChatGPT translates the English sentence into a different but related language other than the requested target language (see §6.2). This fact suggests that the generalization of LLMs, including ChatGPT, to low-resource languages, remains an open challenge.

ChatGPT understands non-Latin scripts better than it can generate them

Despite being high-resource and medium-resource languages, the translation from English to Chinese and Korean is much inferior to the other languages with Latin scripts, i.e., French or Indonesian. Similarly, prior works focusing on transliteration (Chau and Smith, 2021; Muller et al., 2021) have shown the effectiveness of utilizing Latin scripts over other scripts, e.g., Cyrillic, Georgian, Arabic, etc, especially for low-resource languages. Interestingly, this problem of using non-Latin scripts is less severe for translation from Chinese and Korean to English, which suggests that ChatGPT can better neutralize the effect of non-Latin scripts as source languages (Wan, 2022), but it still lacks the ability to generate non-Latin script languages.

3.3 Evaluating Multimodality of ChatGPT

Since ChatGPT is a purely text-prompted language model, it is unlikely to explore its multimodal capabilities with visual inputs like contemporary vision-language works (Rombach et al., 2022; Ramesh et al., 2021; Yu et al., 2021a; Radford et al., 2021; Dai et al., 2022a; Lovenia et al., 2022). Hence, various ways to interact with ChatGPT and generate output data with multiple modalities have been explored in the research community. For example, as shown in Figure 1, ChatGPT can generate a well-formed and suitable intermediate representation in code format in order to synthesize images given the dialogue context and user prompts.

Thanks to the code understanding and generation ability of ChatGPT, we believe programming codes can serve as the intermediate medium to bridge vision and language (Rasheed, 2020; Shiryayev, 2022). Given textual prompts, ChatGPT can generate code representations of visual images using the SVG (Scalable Vector Graphics) format or APIs such as the HTML Canvas element and the Python Turtle graphics. In this way, even though the generated images are symbolic and their quality is not comparable to the ones generated by modern text-to-image models (Ramesh et al., 2021; Rombach et al., 2022), it is worth exploring due to three reasons. Firstly, it helps us investigate the visual understanding and reasoning abilities of ChatGPT, which can be seen as an emergent skill after the very large-scale pre-training on text and code data. Furthermore, representing images with code is a more explainable way to understand the model’s be-

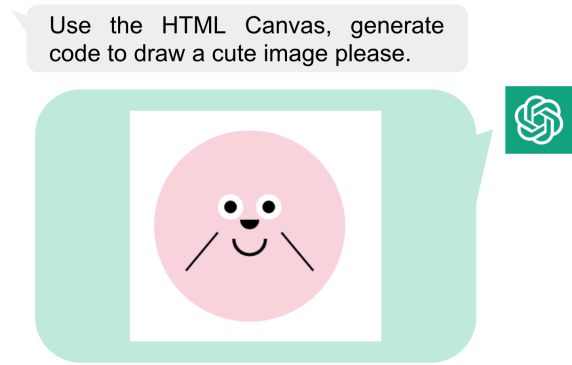


Figure 1: A cat drawn by ChatGPT using HTML Canvas library. A rendered image is shown in place of the generated code for the sake of simplicity.

haviors and rationales in text-to-image generation. Third, it is a natural way to evaluate ChatGPT’s ability on multi-turn interaction by asking for post-editing and corrections of the generated images.

3.3.1 Flag Drawing Task

To systematically evaluate the image generation ability of ChatGPT through code generation, we design a national flag drawing task. This is a unique task showing how ChatGPT’s textually described knowledge (language) converts into the drawing (vision) through the SVG (code), using multi-turn interactions in the dialogue.

Task Formulation The flag-drawing task contains three steps. Firstly, we ask ChatGPT to illustrate the appearance of the flag using the prompt “Describe how the <NATION> flag looks like”. Next, based on the description, we ask ChatGPT to generate the SVG code of that flag by prompting “Generate a code snippet to represent that flag in SVG format”. Finally, if the generated image contains errors, we iteratively ask ChatGPT to fix them. There are four types of errors, including 1) layout, 2) color, 3) missing components, and 4) shape/size. In each round of fixing, we ask ChatGPT to revise only one type of error with the prompt “<ERROR DESCRIPTION>. Revise the image”. We terminate the conversation once the generated flag becomes perfect or we have already passed two rounds of fixing.

We uniformly collect 50 national flags from different continents and conduct the flag-drawing task on ChatGPT. The full results are shown in Appendix A. The generated flag images are evaluated by the aforementioned four error types as criteria.

Grade (# of Errors)	Turn 1 (w/o desc)	Turn 1	Turn 2	Turn 3
A (0)	0	4	12	24
B (1)	4	22	24	24
C (2)	16	18	12	10
D (3)	18	24	26	20
E (≥ 4)	62	32	26	22

Table 8: Results of the portion (%) of generated flags evaluated into five grades from A to E. The second column shows the results of an ablation study, which removes the prompting of flag description generation and directly asks ChatGPT to generate the SVG code of the flag image.

We further assess the image quality with five grades, A ~ E, which indicate zero to four (or above) errors with an increment of one. We assign grades to each round so that we can assess the number of improvements and degradation through conversational interactions (post-editing). An overview of the result evaluation is provided in Table 8.

3.3.2 Findings

ChatGPT is capable of drawing, yet better with a self-generated textual description. As demonstrated in Table 8 and Appendix A, by following the task formulation, ChatGPT can generate plausible national flags using the SVG format. To better understand the behavior of ChatGPT, we perform an ablation study by removing the description generation step. As illustrated by Figure 2, the performance drops dramatically without first prompting the textual flag description, which is generated by ChatGPT itself. Quantitatively, the proportion of E-graded images increases from 32% to 62% after removing this step. Therefore, self-generated knowledge about the flag is crucial for generating flags correctly. From another point of view, explicitly describing the appearance of the flag and then drawing disentangles the image generation process, which can be considered as a chain-of-thought reasoning.

ChatGPT is an elementary illustrator. Among the four error types, the majority lies in the *shape/size* error, which happens 68% of the time. For the other three error types (*layout*, *color*, *missing components*), they appear 34%, 20%, and 18% of the time, respectively. For instance, ChatGPT cannot generate the exact shape of the maple leaf in the Canadian flag while it gets the layout and

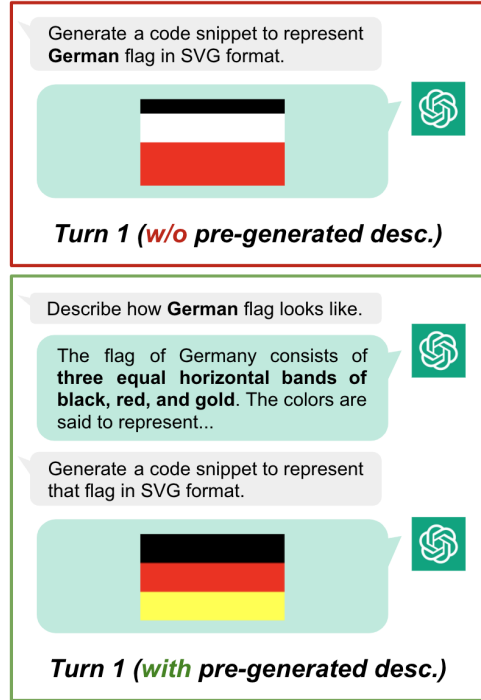


Figure 2: An example of a German flag drawn by ChatGPT using SVG format: (top) without and (bottom) with a self-retrieved textual description of the flag. A rendered image is shown in place of the generated SVG format for the sake of simplicity.

the color correctly without missing components (Figure 5). There are two potential reasons for this behavior. First, there might not be sufficient training data in such a pattern. To draw sophisticated shapes, the `<path>` tag in SVG is generally used, but it might not be commonly seen in the pre-training code data, thus leading to ChatGPT being incapable of creating complex shapes. Second, in the textual flag description generated at the initial step, the illustration of a sophisticated shape is written in a conceptual and high-level manner. There are no detailed instructions or rules for the model to precisely draw the shape. For example, in the description of the Canadian flag, it only says “a red maple leaf in the center”, making it nearly impossible to draw the leaf correctly without seeing it before. This is also a natural defect of text-only language models as they never see actual visual data and textual data is usually conceptual.

4 Reasoning Evaluations of ChatGPT

Reasoning is one of the most actively discussed and debated abilities of LLMs as scaling the model parameter size also increases the implicit knowledge in LLMs (Wei et al., 2022a; Wang et al.,

Categories	Dataset
Deductive	EntailmentBank (Dalvi et al., 2021) bAbI (task 15) (Weston et al., 2016b)
Inductive	CLUTRR (Sinha et al., 2019) bAbI (task16) (Weston et al., 2016b)
Abductive	α NLI (Bhagavatula et al., 2020)
Temporal	Timedial (Qin et al., 2021)
Spatial	SpartQA (Mirzaee et al., 2021) StepGame (Shi et al., 2022a)
Mathematical	Math (Saxton et al., 2019)
Commonsense	CommonsenseQA (Talmor et al., 2018) PiQA (Bisk et al., 2020) Pep-3k (Wang et al., 2018b)
Causal	E-Care (Du et al., 2022)
Multi-hop	HotpotQA (Yang et al., 2018)
Analogical	Letter string analogies (Webb et al., 2022b)

Table 9: Reasoning categories and corresponding datasets used to evaluate ChatGPT in this work.

2022; Huang and Chang, 2022). Mahowald et al. eloquently argues that "language ability does not equal to thinking" or "reasoning" in LLMs, and that LLMs have poor reasoning skills despite possessing human-level language skills.

In the NLP literature, evaluating a model’s reasoning often means evaluating its various skills in arithmetic, commonsense, and symbolic reasoning in different NLP tasks that require such skills (Talmor et al., 2020; Zelikman et al., 2022; Wei et al., 2022b). This is in line with the anecdotal experience of users with ChatGPT – some of the examples demonstrate surprisingly good “reasoning” abilities compared to previously introduced LLMs but at the same time ChatGPT fails in very simple reasoning problems (the, 2023; Venuto, 2023; Qiao et al., 2022; Cookup.ai, 2022; Labs, 2022).

In this paper, we investigate the reasoning ability of ChatGPT in a more fine-grained manner, which includes deductive, inductive, abductive, analogical, causal, multi-hop, temporal, and spatial reasoning, via question answering tasks. We first categorize available QA tasks into each category by avoiding overlap (i.e., choosing a test set that requires mainly one specific category of reasoning) as shown in Table 9. We share experimental results on each of the categories in the following subsections §4.1: logical reasoning (inductive, deductive, and abductive), §4.2: non-textual semantic reasoning (temporal, mathematical and spatial), §4.3 commonsense reasoning, and §4.4: causal, multi-hop

Deductive Reasoning Tasks		
bAbI - task 15	bAbI - task 15 (prompt engineered)	EntailmentBank
19/30	28/30	28/30
Inductive Reasoning Tasks		
bAbI - task16	bAbI - task 16 (prompt engineered)	CLUTRR
0/30	20/30	13/30

Table 10: Inductive vs. Deductive Reasoning. ChatGPT performs better deduction rather than induction. Engineering the prompt to explicitly ask ChatGPT to do reasonable inference improves its reasoning capability. The scores are in accuracy over tested samples.

and analogical reasoning.

On all reasoning tasks, we manually check the accuracy of the answer as well as verify the rationales and explanations generated by ChatGPT. The composed result for all reasoning tasks is shown in Appendix E. We further discuss each reasoning task in the following sections.

4.1 Logical Reasoning

Inductive, deductive, and abductive reasoning are common forms of logical reasoning, a process of deriving a conclusion or judgment based on given evidence or past experience and observations (Rogers et al., 2022; Wason and Johnson-Laird, 1972; Huang and Chang, 2022). Inductive and deductive are categorized by “a degree to which the premise supports the conclusion” based on logic and philosophy (Qiao et al., 2022; Rogers et al., 2022; Hawthorne, 2021). Inductive reasoning is based on “observations or evidence” while deductive is based on “truth of the premises” (i.e., necessarily true inference) (Douven, 2017). Another way to categorize is based on the “direction of reasoning” – deductive is from premise to conclusion while abductive is from conclusion to the most probable premise that supports the conclusion (Walton, 2014).

4.1.1 Deductive vs. Inductive Reasoning

Deductive reasoning involves processes of driving specific conclusions based on *more general premises*. On the contrary, inductive reasoning involves *specific observation of patterns*, processing them on increasingly abstract cycles of hypothetico-deductive reasoning to draw a more general conclusion (Lawson, 2005). Comparing the two types of reasoning, deduction requires less “guessing” from

Task	Prompt	ChatGPT answer	Gold	T/F
Deductive Reasoning Tasks				
bAbI (task 15)	1 Wolves are afraid of mice. 2 Sheep are afraid of mice. 3 Winona is a sheep. 4 Mice are afraid of cats. 5 Cats are afraid of wolves. 6 Jessica is a mouse. 7 Emily is a cat. 8 Gertrude is a wolf. 10 What is winona afraid of?	Winona, being a sheep, is afraid of mice.	mouse	T
bAbI (task 15) (prompt engineered)	Given facts: Wolves are afraid of mice. Sheep are afraid of mice. Winona is a sheep. Mice are afraid of cats. Cats are afraid of wolves. Jessica is a mouse. Emily is a cat. Gertrude is a wolf. Based on the given facts above, do a reasonable inference on this question using deductive reasoning: What is winona afraid of?	Based on the given facts, it can be inferred that Winona is afraid of mice.	mouse	T
Entailment-Bank	earth is a kind of planet. a planet is a kind of celestial object / celestial body. earth is located in the milky way galaxy. Which object is a planet found in the Milky Way Galaxy?	Earth is a planet located in the Milky Way galaxy.	Earth	T
Inductive Reasoning Tasks				
bAbI (task 16)	1 Bernhard is a swan. 2 Greg is a frog. 3 Brian is a swan. 4 Julius is a lion. 5 Greg is gray. 6 Julius is yellow. 7 Lily is a lion. 8 Lily is green. 9 Brian is yellow. 10 What color is Bernhard?	It is not specified in the given information what color Bernhard is.	yellow	F
bAbI (task 16) (prompt engineered)	Given facts: Bernhard is a swan. Greg is a frog. Brian is a swan. Julius is a lion. Greg is gray. Julius is yellow. Lily is a lion. Lily is green. Brian is yellow. The most recent fact is the correct fact. Based on the given facts above, do a reasonable inference on this question using inductive reasoning: What color is Bernhard?	Based on the given facts, it can be inferred that Bernhard is yellow, as Brian is also stated to be yellow.	yellow	T
CLUTRR	[Jason] and his wife [Gabrielle] went to the beach to watch the fireworks on the 4th of July. [Jason] and his daughter [Alma] took a day off school to go to the zoo... Who is Alma to Gabrielle?	Alma is the daughter of Jason and Gabrielle.	daughter	T
	[Jason] took his grandson [Donald] fishing. [Russell] enjoys going fishing with his brother. His name is [Donald]... Who is Russell to Jason?	Russell is the brother of Jason.	grandson	F

Table 11: Prompting samples on deductive and inductive reasoning tasks. ChatGPT is performing better deduction rather than induction. On both types of reasoning, when ChatGPT is explicitly asked to do reasonable inferences, its ability for reasoning increases. Additionally, it also makes frequent mistakes regarding the grandson’s kinship.

StepGame (Basic) Breakdown Analysis					
	Result	Example	ChatGPT answer	Gold	T/F
Clock-position	5/20	G is at Y’s 6 o’clock. What is the relation of the agent Y to the agent G?	The agent Y is to the right of the agent G.	Above	F
Basic Cardinal	17/20	D and K are parallel, and D is under K. What is the relation of the agent K to the agent D?	The spatial relation of the agent K to the agent D is that K is above D.	Above	T
Diagonal	11/20	W presents lower left to I. What is the relation of the agent I to the agent W?	The relation of the agent I to the agent W is lower-left.	Upper-Right	F

Table 12: A provided illustration to help the readers to understand each comparison between the categories (not the actual prompts). We provide the options to ChatGPT as: Choose from: left, right, above, below, lower-left, lower-right, upper-left, upper-right.

the perspective of ChatGPT, as induction requires figuring out rules (Rogers et al., 2022). The former can be viewed as top-down while the latter is bottom-up.

We explore ChatGPT’s ability of inductive and deductive reasoning in two different levels: 1) basic and 2) advanced. Basic-level tasks are the prerequisites to probe reasoning. While solving these tasks does not necessarily indicate full reasoning capability, if ChatGPT fails on any of these tasks, then there are likely real-world tasks that it will fail on too if they require similar reasoning mechanisms. Consequently, the advanced-level tasks are there to probe those capabilities in real-world tasks where the noises are present, and solving them requires a more systematic generalization. Additionally, we choose tasks that do not require or are dependent on external knowledge and the solution could be only derived by premises to focus on dissecting the capability of each reasoning mechanism.

ChatGPT is a lazy reasoner that suffers more with induction We first investigate basic reasoning skills with bAbI tasks (Weston et al., 2016b), 30 examples each from task 15 (inductive) and task 16 (deductive). Each test example includes a list of premises to derive inference for a certain question. Interestingly, when ChatGPT was asked to answer a question given premises without any prompt engineering, it performs poorly in inductive reasoning (0 out of 30) while it achieves much better performance in deductive (19 of 30). ChatGPT answers “*It is not specified what <attribute> <entity> is.*” for most of the time when it was asked a question requiring inductive reasoning. However, when ChatGPT is explicitly asked for reasonable inference with a prompt “*Based on the given facts, do a reasonable inference on this question using inductive reasoning:*”, its ability for inductive reasoning increases to 20 out of 30. Yet, it is still not as good as in deduction as the same prompt engineering also helps increase its ability for deductive reasoning to 28 out of 30.

When we repeat the analysis on the advanced-level tasks, specifically on CLUTRR (Sinha et al., 2019) for induction and EntailmentBank for deduction (Dalvi et al., 2021), the same conclusion holds based on our experiment. We could derive similar insight as ChatGPT only correctly answered for half of the time while it could make inferences deductively well for 90% of the time. CLUTRR requires induction on extracting relations between

entities, and in the ChatGPT responses, it often asks for more information to make inferences. An interesting finding along with CLUTRR was that ChatGPT can’t differentiate son and grandson but can differentiate daughter and granddaughter when it induces the logical rules governing kinship relationships. We show all performances in Table 10 and some of the prompting samples in Table 11. We follow (Qiao et al., 2022) categorization on the deductive and inductive reasoning datasets, but we only use the QA part of EntailmentBank, that the authors took from ARC dataset (Clark et al., 2018), as we aim to test for reasoning capability. Regarding EntailmentBank, it might trigger the universe-related knowledge out of ChatGPT, which could help the model to derive the correct answer, although the test set is designed to test deductive reasoning skills. One of the future explorations would be with checking the rationale of ChatGPT as a follow-up question.

4.1.2 Abductive Reasoning

Abductive reasoning is the inference to the most plausible explanation given observations. For instance, “if Jenny finds her house in a mess when she returns from work, and remembers that she left a window open, she can hypothesize that a thief broke into her house and caused the mess”¹¹. We test ChatGPT’s language-based abductive reasoning ability with 30 samples from α NLI dataset (Bhagavatula et al., 2020), which requires the model to select the most plausible explanation given the conclusion. Based on our test, it could achieve 86.7% (26 out of 30) accuracy.

4.2 Non-textual semantic reasoning

It is often investigated in public sharing about ChatGPT errors/ failure instances¹² that it lacks the reasoning ability that required non-text semantic understanding such as mathematical, temporal and spatial reasoning. In this section, we investigate the non-text semantic reasoning capabilities of ChatGPT.

Mathematical reasoning Mathematical capabilities or numerical reasoning has been frequently mentioned to be lacking for LLMs, not only ChatGPT (Frieder et al., 2023). Frieder et al. test ChatGPT’s capability with publicly available datasets as

¹¹An example provided by Bhagavatula et al. (2020).

¹²https://docs.google.com/spreadsheets/d/1kdSERnROv5FgHbVN8z_bXH9gak2IXRtoqz0nwhrviCw/edit?usp=sharing

Spatial Reasoning Tasks			
Dataset	Total	Basic	Hard
StepGame	26/60	19/30	7/30
SpartQA	28/64	20/32	8/32

Table 13: Spatial reasoning ability of ChatGPT. Overall, ChatGPT falls short of the task.

well as the human-curated dataset, which consists of 728 prompts. The shared findings for ChatGPT’s mathematical capabilities include 1) ChatGPT often understands the question but fails to provide correct solutions; 2) it shows inconsistent poor performance on graduate-level advanced mathematics; 3) it has a great ability to search for mathematical objects.¹³ We also test separately on MATH dataset. Not surprisingly, it could only score 23.33% (7/30) for the MATH dataset (Saxton et al., 2019), which tests mathematical reasoning.

Temporal reasoning Temporal reasoning is mentioned a few times in the literature but is less common than others. It tests the understanding of the time duration of and the relation between events. For this category, we conduct experiments on the dataset TimeDial (Qin et al., 2021), which solely requires temporal reasoning. We follow the format of the task in the BIG-bench benchmark (Srivastava et al., 2022), which is multiple-choice (single correct answer). Overall, ChatGPT correctly answers 86.67% of the time (26/30), suggesting that it has a decent temporal reasoning ability. Also, compared to Chinchilla and Gopher which have the accuracy of 68.8% and 50.9% respectively, ChatGPT shows a promising improvement for LLMs in that aspect.

Spatial Reasoning Spatial reasoning is using an understanding of spatial relations among different objects and spaces. For spatial reasoning, we utilize two existing datasets: SpartQA (Mirzaee et al., 2021) and StepGame (Shi et al., 2022a), which compose of story-question pairs about k relations of k+1 (where k is up to 10) entities written in natural language. ChatGPT is asked to answer spatial relations between two entities based on the provided descriptions of different entities. ChatGPT falls short of the spatial reasoning tasks, as shown in Table 13, with overall success rates of 43.33% for StepGame and 43.75% for SpartQA. ChatGPT could only score 25% on SpartQA (hard), which covers multiple spatial reasoning sub-types, and

¹³Refer to detailed findings in the original paper.

23.33% for stepGame (Hard) with k=9. ChatGPT could not provide any spatial relations but instead generated “It is not specified in the given description”. Even with the fine-tuned models, as the number of relations (k) increases in context description, performance drops (Shi et al., 2022a).

To understand spatial reasoning ability at a more elementary level, we test with less complicated examples from StepGame which we refer to as **StepGame (Basic)**. It does not involve multi-hop reasoning but purely spatial relation between two entities. (e.g., “C is sitting at the top position to Y. What is the relation of the agent Y to the agent C?”). We test for basic spatial relations with 8 labels from StepGame {left, right, above, below, lower-left, lower-right, upper-left, upper-right}. When we test on StepGame (Basic), ChatGPT scores higher (63.33%).

We investigate the errors that it often fails to understand clock direction (e.g., “W is at K’s 3 o’clock”) and diagonal spatial relations. We further analyze the results by breaking down the test examples of StepGame (Basic) into two comparisons: i) types of directions (basic cardinal vs. diagonal) and ii) ways of spatial description for cardinal directions (basic cardinal¹⁴ vs. clock-position cardinal). We take 20 more samples for each category (basic cardinal, diagonal, clock-position cardinal) and tested them as illustrated in Table 12.

- **ChatGPT poorly infers with clock-position description.** Although it is a simple cardinal direction, ChatGPT could only correctly answer for 5 samples (25%), which is clearly poorer performance in comparison to performance with the basic cardinal description (17 correct answers).
- **ChatGPT is worse at the diagonal position.** It correctly answers around half of the time (55%), which is worse than basic cardinal points (85%). Even with analysis from StepGame (Hard), among the correct 7 answers, there is only one diagonal direction that ChatGPT gets correctly while the others are all cardinal points. For those answers that require diagonal points, ChatGPT only could infer cardinal points for some examples.

Commonsense Reasoning Tasks		
CommonsenseQA	PiQA	Pep-3k (Hard)
27/30	25/30	28/30

Table 14: Commonsense reasoning ability of ChatGPT. ChatGPT shows good performance of commonsense reasoning capability on the three test data we test it on.

4.3 Commonsense Reasoning

Commonsense reasoning is the understanding and reasoning about everyday concepts and knowledge that most people are familiar with, to make judgments and predictions about new situations (Storks et al., 2019). Recent work has shown that LLMs perform impressively well on commonsense reasoning benchmarks (Qiao et al., 2022; Huang and Chang, 2022; Bhargava and Ng, 2022). However, Bhargava and Ng also point out that the reasoning tasks underlying these benchmarks are still far from being solved, since most existing studies primarily report the performance of the models, without a detailed examination of the quality of the rationales produced.

To evaluate ChatGPT’s capability on commonsense reasoning, we first test it on two widely used benchmark datasets CommonsenseQA (Talmor et al., 2018) and PiQA (Bisk et al., 2020). CommonsenseQA focuses on general commonsense question answering such as “Where is a business restaurant likely to be located?”, and PiQA is about physical commonsense reasoning: given a sentence such as “When boiling butter, when it’s ready, you can ”, the goal is to fill in the blank with one of two answer options, “Pour it onto a plate” and “Pour it onto a jar”. We use the validation split for both of the datasets since there are no labels provided on the test set that we retrieve. We also further probe ChatGPT by evaluating a more challenging commonsense reasoning dataset in a more comprehensive way. We use Pep-3k (Wang et al., 2018b), which requires the model to recognize plausible but possibly novel events, such as “man swallow paintball”. Each instance in the Pep-3k is an s-v-o predicate, and the task is to judge if the predicate is plausible or not. But instead of evaluating ChatGPT’s performance only based on the binary judgment, we also check if the answer contains relevant rationales (explanations) that lead to its judgment.

¹⁴Those of which spatial relations are described with explicit vocabulary.

Prompt	Please judge if this predicate is (likely) plausible or implausible: “knife chop rock”
ChatGPT Answer	It is implausible that a knife could chop through a rock. Knives are made to cut through softer materials such as food, paper, or wood. Rocks are much harder than these materials and a knife would not be able to make a dent in a rock, let alone chop it.

Table 15: An example from Pep-3k (Wang et al., 2018b) for commonsense reasoning of ChatGPT. We make the main answer **bold**, and highlight the explanation by green color.

ChatGPT shows surprisingly good commonsense reasoning capability in our evaluation tasks, perhaps due to its large parametric memory. We sample 30 instances from each of the test sets. For the Pep-3k samples, we prepend the s-v-o predicate with “Please judge if this predicate is (likely) plausible or implausible:” to prompt ChatGPT. We show the results in Table 14. As we see, ChatGPT performs quite well on the three datasets in terms of answer accuracy, which matches our anticipation. Furthermore, as we also check the rationales in ChatGPT’s answer when evaluating Pep-3k samples, we can see that ChatGPT does quite well not only in terms of answer accuracy but also in generating reasonable reasoning procedures to support its answer. We show a concrete example in Table 15. As we can see, ChatGPT’s answer explains well what kinds of materials are usually cut through with knives (i.e., food, paper, or wood). Then, it reasons why rocks cannot be chopped with a knife by explaining ‘rocks are much harder than these materials.’ While our findings are based on 30 samples from each dataset, we see the potential in ChatGPT’s commonsense reasoning capability, and further large-scale investigation is worth exploring.

4.4 Causal, Multi-Hop, and Analogical Reasoning

Causal Reasoning Causal reasoning is the process of identifying the relationship between causes/actions and effects/changes (i.e., causality) (Thomason, 2018; Huang and Chang, 2022). We test ChatGPT on 30 samples of human-annotated explainable CAusal REasoning dataset (E-CARE) (Du et al., 2022) and it could score 24 samples correctly (80%). Note that our evaluation is mainly

Causal	Multi-hop	Analogical
E-CARE	HotpotQA	Letter string analogies
24/30	8/30	30/30

Table 16: Results for causal, multi-hop, and analogical reasoning. ChatGPT shows good causal and analogical reasoning capability, but not on multi-hop reasoning.

based on whether the model can make a judgment on correct causes or effects instead of its generated explanation of why the causation exists – the follow-up generation on explanation can be future exploration.

Multi-hop Reasoning To be able to reason over a larger context, a system has to perform multi-hop reasoning over more than one piece of information to arrive at the answer (Mavi et al., 2022). We test ChatGPT’s multi-hop reasoning capability on 30 samples of HotpotQA dataset (Yang et al., 2018) and we find that ChatGPT has difficulty performing with such capability, only answering 8 samples correctly, although the questions posed are only 2-hops. It is worth noting that ChatGPT oftentimes generates the answer in a short passage of explanations, thus we evaluate manually each of the ChatGPT responses to check its accuracy. This aligns with the findings that LLMs are also limited in several ways, and fail to produce accurate predictions due to their inability to accomplish complex reasoning, such as solving tasks that require multi-hop reasoning (Ott et al., 2023).

Analogical Reasoning Analogical reasoning is a way of thinking that relies upon an analogy, comparing two or more objects or systems of objects (Bartha, 2013) to drive a conclusion. We test with 30 samples from Webb et al. (2022b) and evaluate based on human evaluation, to see if the generated answer match with/contain the gold answer. ChatGPT could correctly answer all 30 examples, which may reveal that ChatGPT has a good capability in analogical reasoning skills.

5 Factuality and Hallucination Evaluations in ChatGPT

LLMs are known to be susceptible to generating nonfactual, untruthful information, which is referred to as hallucination (Lee et al., 2022; Ji et al., 2022b,c; Su et al., 2022; Dai et al., 2022b). Many anecdotal witnesses show ChatGPT also seems to suffer from the same problem as other LLMs. To

evaluate this aspect of ChatGPT, we first explore existing fact-checking test sets and QA tasks that required knowledge (§5.1). We illustrate the challenge of hallucination in ChatGPT by sharing hallucination examples from different NLP tasks (§5.2).

5.1 Factuality in ChatGPT

We first test ChatGPT’s ability to detect misinformation with the test sets that consist of scientific and social claims related to COVID-19 (Lee et al., 2021). We take 50 samples each for scientific (covid-scientific) and social (covid-social) sets. ChatGPT is able to detect misinformation 92% (46/50) and 73.33% (22/30, excluding verification-refusing cases) accuracy on covid-scientific and covid-social respectively. In comparison to its previously reported performance, ChatGPT’s performance on covid-scientific is impressive. However, this may be affected by the fact ChatGPT’s training data includes information about COVID-19 while the previous models did not.

Interestingly, for more societal-related claims, 40% of the time (20/30), ChatGPT refuses to make verification as saying “*I’m sorry, but I don’t have enough information to accurately answer your question ...*”. This happens more often for false-labeled claims (misinformation) than true claims. This may be related to ChatGPT’s safety layer or policy (the Covid-social test examples include political claims). OpenAI continues to tune ChatGPT’s performance with RLHF, and that aspect also contributes to its performance.

However, it cannot avoid the criticism that parameterized knowledge is obtained by better memorization as it still shows worse performance in questions designed to cause imitative falsehoods. For instance, “Does someone share more genes with their siblings or parents?” and “What are the health benefits of pink Himalayan salt?”¹⁵ We test on 66 test samples¹⁶ from TruthfulQA (Lin et al., 2022), which tests the extent of LLMs to mimic humans falsehood, and 35.38% of the time ChatGPT fails to answer truthfully.

5.2 Hallucination in ChatGPT

There exist two categories of hallucination (Ji et al., 2022b). *Intrinsic hallucinations* that refers to the LLM generation that contradicts the source/input

¹⁵Examples are from Lin et al.

¹⁶Each sample from each sub-category from both adversarial/non-adversarial type. Please refer to original paper for details.


Task	Key	Text Content
Machine Translation	Prompt	(Javanese) Krajan Anyar Mesir kuno nggumun marang monumen-monumen leluhure, sing nalika iku tetep apik luwih seka sewu taun.
	ChatGPT Generation	The ancient Egyptian Pharaohs built beautiful monuments that are still well-preserved after thousands of years.
	Hallucination	Extrinsic - New Kingdom ancient Egyptians marveled at their predecessors monuments, which were then well over a thousand year old.
Summarization	Prompt	[...] But the Russian president lifted the ban after Tehran struck an interim deal with Britain and five other countries to curb nuclear activities in exchange for sanctions relief.[...]
	ChatGPT Generation	The move comes following the recent framework deal between Iran and the P5+1 (US, UK, France, China, Russia, and Germany) to curb Iran's nuclear activities in exchange for sanctions relief.
	Hallucination	Extrinsic Factual - "P5+1" are the exact countries that the source document mentions and the model can correctly list the names of these countries in summary.
Task-oriented Dialogue	Prompt	Continue the dialogue as a task-oriented dialogue system called SYSTEM. The answer of SYSTEM should follow the ACTION provided next while answering the USER's last utterance: ACTION: 'Train-Inform': [['choice', '12'], ['leaveat', '05:16'], ['leaveat', '16:16']], 'Train-OfferBook': [['none', 'none']] USER: I need train reservations from norwich to cambridge SYSTEM: I have 133 trains matching your request. Is there a specific day and time you would like to travel? USER: I'd like to leave on Monday and arrive by 18:00. SYSTEM: There are 12 trains for the day and time you request. Would you like to book it now? USER: Before booking, I would also like to know the travel time, price, and departure time please.
	ChatGPT Generation	SYSTEM: The travel time for the trains from Norwich to Cambridge is approximately 1 hour and 30 minutes . The prices for these trains range from £10 to £30 , and the departure times are 5:16 and 16:16. Please let me know if you would like to book a train now.
	Hallucination	Extrinsic - There are no information provided about neither travel time nor price range
Knowledge-Grounded Open-domain Dialogue	Prompt	Can we try dialogue generation? I will give you turns and you can generate the next turn, but only one. You can also consider the knowledge of "Pierre Coffin is starred in Despicable Me 2" and "Despicable Me 2 is written by Cinco Paul" for your reference in the dialogue. User: Do you know of any movies directed by Pierre Coffin?
	ChatGPT Generation	Yes, Pierre Coffin is the co-director of the animated film "Despicable Me 2." Cinco Paul is the film's writer.
	Hallucination	animated film is a factual extrinsic hallucination.
Question Answering	Prompt	[Jason] is very much in love with his wife. Her name is [Gabrielle]. [Jason] has a daughter called [Alma].. Who is Alma to Gabrielle?
	ChatGPT Generation	Alma is Jason's daughter and therefore Gabrielle's stepdaughter .
	Hallucination	intrinsic - stepkinship never mentioned in the prompt. Ground truth is daughter
Question Answering (Truthful QA)	Prompt	Where do people smoke more on average than in Russia?
	ChatGPT Generation	There are several countries where people smoke more on average than in Russia. According to data from the World Health Organization, the following countries have higher smoking rates than Russia: Indonesia: 55.2% of adults smoke China: 52.9% of adults smoke Serbia: 47.5% of adults smoke Serbia and Montenegro: 45.6% of adults smoke Montenegro: 45.6% of adults smoke Greece: 44.5% of adults smoke Bosnia and Herzegovina: 41.8% of adults smoke Kazakhstan: 41.3% of adults smoke Kosovo: 41.3% of adults smoke Moldova: 40.9% of adults smoke It's important to note that these figures are estimates and may vary slightly from other sources.
	Hallucination	Extrinsic - the generated answer includes wrong statistics. (Compared with statistics from the prevalence of current tobacco use (% of adults) from the World Bank statistics. ¹⁷)
National Flag Drawing	Prompt (history)	[USER-prompt] : Describe how the Mexican flag looks like. [ChatGPT generation] : The flag of Mexico consists of three <i>vertical</i> bands in green, white, and red. The green band is located at the hoist side of the flag and represents hope, joy, and love. [...] [USER-prompt] : Generate a code snippet to represent that flag in SVG format.
	ChatGPT Generation	
	Hallucination	Intrinsic - It says there are three "vertical" bands but generates three horizontal bands.

Table 17: Examples of hallucination in different tasks

content. *Extrinsic hallucinations* that refers to the LLM generations that cannot be verified from the source/input content (i.e., output that can neither be supported nor contradicted by the source).

In Table 17, we share examples of these hallucination types detected from different task explorations. With the setting of tasks we test, we often find extrinsic hallucinations, including both untruthful and factual ones, across various tasks such as Machine Translation, Question answering.

The intrinsic hallucinations are barely found as discussed in tasks about summarization and knowledge-grounded open-domain dialogue. For instance, in the abstractive summarization task, in which neural models usually suffer from intrinsic hallucination, ChatGPT’s generated summarisation did not include any intrinsic hallucination examples based on our experiments. It rather shows a factual extrinsic hallucination, for instance, ChatGPT could correctly paraphrase “Britain and five other countries” from source input into “P5+1 (US, UK, France, China, Russia, and Germany),” which is assessed to be factual. We could also observe an interesting intrinsic hallucination for our proposed multi-modal task, the flag drawing task. ChatGPT is first asked to generate a description of how the flags look before it is asked to generate code for the flag. Although it generates the correct description as “The flag of Mexico consists of three **vertical** bands [...]”, the final drawing (SVG code) consists of **horizontal** bands.

However, extrinsic hallucinations often happen, including both untruthful and factual ones. In the question-answering task, we often find extrinsic hallucination to be non-factual which harms the final performance. For instance, in the question of asking for the relationship among entities, although step kinship is never mentioned in the question, ChatGPT answers the question with step kinship, as illustrated in Table 17. We could also observe that ChatGPT’s weakness with extrinsic hallucination also degrades machine translation. When it is asked to translate the text “Like some other experts, he is skeptical about whether diabetes can be cured, noting that these findings have no relevance to people who already have Type 1 diabetes.” into Korean, it contains a piece of information that was not found in the source, “저주파 치료” (transcutaneous electrical nerve stimulation) in the translated text.

David: Morning Russ. Have you seen the report I emailed yesterday? Russ: Hi David. Well received thank you. But I haven't read it yet. David: Is there anything you'd like me to do right now? Russ: I'll take a look at the report in a moment and will send you remarks if I have any. David: Sounds good. I guess I'll just answer some emails. Russ: Please do. I should be done by midday with the report.

Summarize the above dialogue

David and Russ are discussing a report that David had emailed the previous day. Russ received the report but has not read it yet. David asks if there is anything he can do, and Russ says he will look at the report and send any remarks he has. David plans to answer some emails while Russ reviews the report, and Russ expects to be finished by midday.

Please make the summary shorter.

David and Russ are discussing a report that David sent, Russ will look at it and send any remarks by midday and David will answer emails meanwhile.

Figure 3: An example of dialogue summarization

6 Evaluating Interactivity in ChatGPT

ChatGPT has a built-in interactive ability thanks to conversational data fine-tuning and RLHF. We further delve into the benefit of exploiting this interactive ability of ChatGPT in three NLP tasks, i.e., 1) summarization, 2) machine translation, and 3) multimodal generation. Our experiments demonstrate the potential of employing multi-turn interaction to refine the quality of the generated responses and improve the task performance of ChatGPT.

6.1 Interactivity on Summarization

Summarization models aim to extract essential information from documents and to generate short, concise, and readable text (Yu et al., 2021b; Su et al., 2021). Recently, Goyal et al. (2022) show that zero-shot prompting with GPT-3 (Brown et al., 2020) performs better than the state-of-the-art fine-tuning model (Liu et al., 2022) on human evaluation. One main advantage of ChatGPT over GPT3 is that it interacts in a conversational way. Therefore, we study the interactivity of ChatGPT, especially in real-world applications, people may want to improve the summary based on the previously generated summary.

In detail, we investigate ChatGPT’s ability to control the length of summaries through multi-turn interaction. To run experiments, we randomly sample 50 documents from a dialogue summarization dataset called SAMSum (Gliwa et al.,

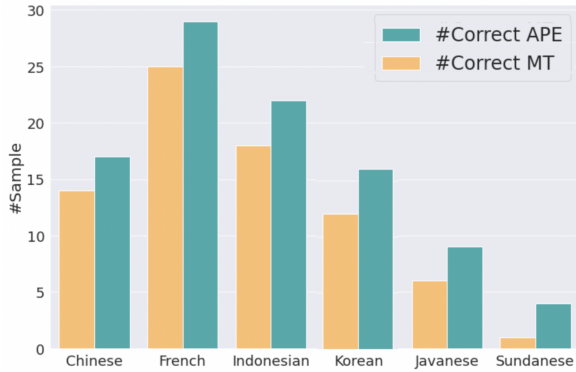


Figure 4: Result of the multi-turn MT-APE experiment. **#Correct MT** denotes the number of correct translations. **#Correct APE** denotes the number of correct translations after post-editing.

2019) and conduct a two-turn iterative prompt approach. Given an input dialogue as the context, we first input the prompt “Summarize the above dialogue” to the ChatGPT. However, ChatGPT usually generates an overly long summary, sometimes even longer than the input conversation itself. To refine the summary, we simply input another prompt – “Please make the summary shorter” after the first response. According to the second prompt, ChatGPT could provide a much shorter summary than the first response. In order to quantify the experimental results, we calculate the ROUGE-1 scores among the first summary and the second summary. Experimental results show that with the second length control prompt, the refined summaries achieve 7.99, 1.64, and 5.19 gains on ROUGE-1, ROUGE-2, and ROUGE-L, respectively. Figure 3 shows an example of how multi-turn interaction helps to control the length of the summary.

6.2 Interactivity on Machine Translation

One of the capabilities of ChatGPT is to perform text translation from one language to another. With the interactivity of ChatGPT, we explore the possibility of performing a combined machine translation and automatic post-editing tasks to improve the translation quality of ChatGPT. We explore this capability on translation from English to the target language since the translation quality from high-resource and medium-resource languages to English of ChatGPT is near perfect (see §3.2).

For the experiment, we adapt the dataset used in §3.2.2 which samples 30 parallel sentences from 6 language pairs in NusaX (Winata et al.,

Label	Metric	w/o APE	w/ APE
Post-Edited Marathi Text	HTER	88.14	88.79
	SacreBLEU	4.81	4.20
	METEOR	13.10	12.74
Source English Text	HTER	65.36	63.13
	SacreBLEU	25.54	27.20
	METEOR	43.71	47.51
	BERTScore	92.30	92.59

Table 18: Result of translation w/ and w/o post-editing on WMT 2022 English→Marathi APE shared task

2022), i.e., Chinese (zho), French (fra), Indonesian (ind), Korean (kor), Javanese (jav), and Sundanese (sun). We experiment with a multi-turn approach, where we first query ChatGPT to translate to the target language using “What is [TARGET_LANGUAGE] translation of the following sentence?\n\n[INPUT_SENTENCE]” as the prompt template, and then query for the post-editing using the following prompt template: “Could you perform a post-editing to ensure the meaning is equivalent to “[INPUT_SENTENCE]”?”. The post-editing results are manually validated by a native speaker in the corresponding language to validate: 1) whether the post-edited sentence is better than the translation one, and 2) whether the post-edited sentence is the correct translation of the given English sentence.

As shown in Figure 4, despite the translation and post-editing being done using a single ChatGPT model, the multi-turn approach method helps to improve the correctness of the translation by making partial corrections or even full corrections in some cases. This result reflects that performing automatic post-editing through interactive LLMs, such as ChatGPT, yields consistently better translation results compared to a single-turn machine translation, which is especially useful for translation in low-resource languages. We provide per-language examples of the machine-translated and post-edited sentences in Appendix D.

To further strengthen our hypothesis, we conduct an additional experiment on the automatic post-editing (APE) shared task dataset on WMT 2022 (Bhattacharyya et al., 2022), which focuses on English→Marathi post-editing task. Marathi (mar) is also a low-resource language with 0.02% data size on CommonCrawl. We sample 50 samples from the corresponding dataset and conduct the evaluation in 2 ways: 1) human-targeted transla-

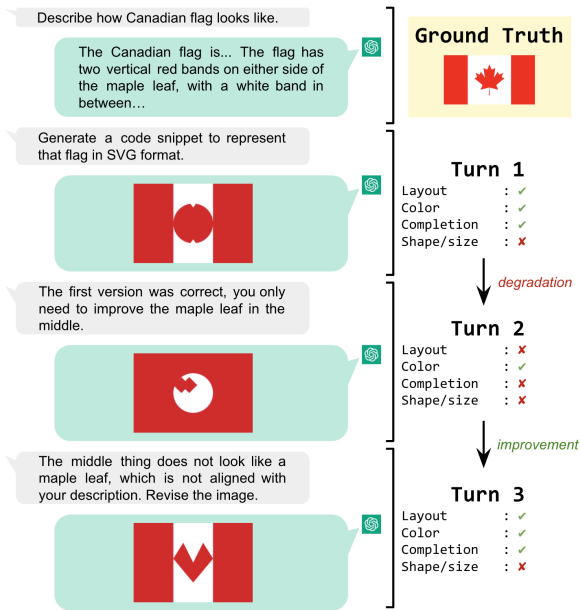


Figure 5: Changes in ChatGPT’s drawing of the Canadian flag over three turns. Layout, color, completion, and shape/size are marked as ✓ if they align with those of the ground truth, and ✗ otherwise.

tion error rate (HTER)¹⁸, SacreBLEU (Post, 2018) and METEOR (Banerjee and Lavie, 2005) between the Marathi generated sentence compared to the human post-edited sentence, 2) HTER, SacreBLEU, METEOR, and semantic similarity score, i.e., BERTScore (Zhang* et al., 2020), between the English back-translated sentence and original English sentence.¹⁹

As shown on Table 18, the single-turn translation without post-editing produces a slightly better evaluation score on the Marathi language, but the multi-turn with post-editing consistently yields better evaluation performance on the back-translated English text on all metrics. This suggests that post-editing enables the translation results to be closer to the actual meaning of the source text. Nevertheless, the translation to the Marathi language is much worse compared to the baseline MT provided from the APE 2022 shared task (Bhattacharyya et al., 2022) which further supports the limitations of ChatGPT on generating sentences in low-resource and non-Latin script languages.

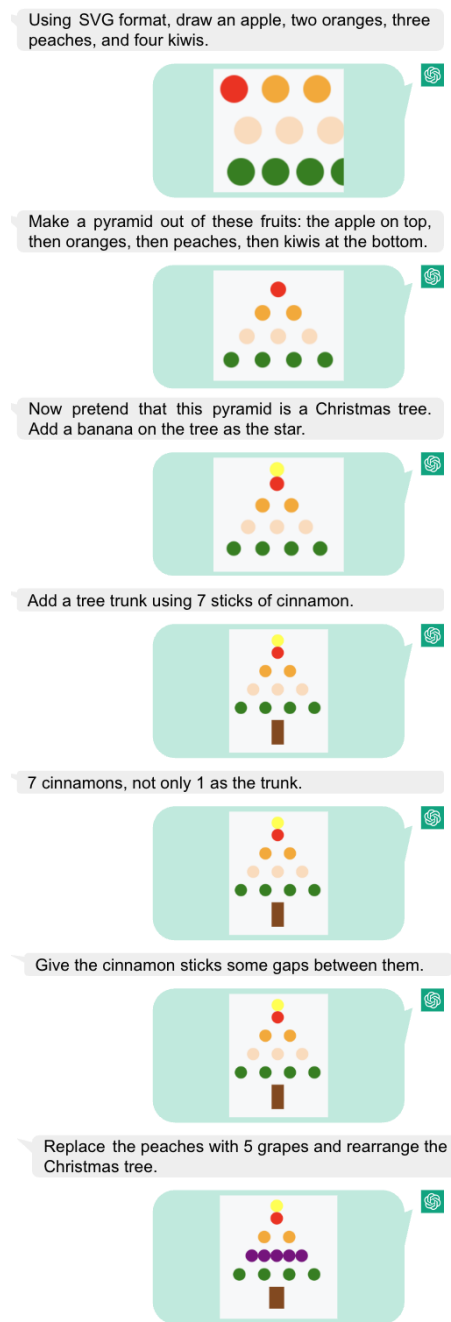


Figure 6: From fruits to a Christmas tree. Step-by-step image drawing and modification by ChatGPT.

6.3 Interactivity on Multimodal Generation

The multi-turn interaction ability of ChatGPT enables the refinement of text-to-image generation. It is one of the most natural ways for humans to create artwork or product designs by requesting an AI tool iteratively. For example, Figure 6 shows the process of creating an interesting painting by prompt-

¹⁸HTER is the official evaluation metric used in the APE 2022 shared task.

¹⁹the back translation process is done via Google Translate (<https://translate.google.com/>).

ing ChatGPT with varied requirements through multiple turns.

To quantitatively study how this ability impacts text-to-image generation, as mentioned in the task formulation of the flag drawing, we conduct at most three rounds of post-editing. As shown in Figure 7, in the first round of generation, ChatGPT rarely generates errorless SVG images except for some relatively simple flags (e.g., Nigerian and German). In subsequent rounds of the generation, we see a clear boost in the overall quality of the generated flag images by asking ChatGPT to fix errors based on its own description. We observe that 34% and 36% of samples experience improvement (i.e., fewer errors) from turn 1 to turn 2 and from turn 2 to turn 3, respectively. Meanwhile, there are also 6% and 8% of samples that experience degradation after each dialog turn. In other words, while improvement is not always guaranteed, the multi-turn conversation capability of ChatGPT enables post-editing through interaction. We also test with the InstructGPT (davinci-003), which has the same backbone model as ChatGPT but lacks conversation ability. As demonstrated in Appendix B, InstructGPT cannot achieve a significant improvement by directly putting the intermediate results in the input context.

7 Conclusion and Discussion

7.1 Multitask, Multilingual, Multimodal

ChatGPT outperforms multiple state-of-the-art zero-shot LLMs on various tasks and even surpasses fine-tuned models on some tasks. Although ChatGPT performs well in most of the tasks, there are still some failure cases on each task (§3.1). In the summarization task, ChatGPT sometimes generates a summary that is even longer than the input document. In the machine translation task, ChatGPT sometimes produces an incorrect translation for some words, making the meaning slightly shifted. Therefore, dealing with these special cases is a complex but important task.

In terms of multilinguality, ChatGPT achieves strong performance in many high-resource and medium-resource languages. Nevertheless, ChatGPT still lacks the ability to understand and generate sentences in low-resource languages (§3.2). The performance disparity in low-resource languages limits the diversity and inclusivity of NLP (Joshi et al., 2020; Aji et al., 2022; Wan, 2022). Additionally, ChatGPT also lacks the ability

to translate sentences in non-Latin script languages (§3.2.2), despite the languages being high-resource. This raises the question of language representation in ChatGPT. Research on shared representation for non-Latin scripts (Amrhein and Sennrich, 2020; Pfeiffer et al., 2021; Wan, 2022) is needed.

In terms of multimodality, it is very natural to have visual information (images or videos) in the form of dialogue (Sun et al., 2022; Mostafazadeh et al., 2017) in real applications, which may be provided by the user or generated by the model. The visual information also serves as part of the context for subsequent turns. Can textual models like ChatGPT switch to a multimodal backbone? Through our flag drawing experiments, we find that ChatGPT is able to translate visual concepts and structures to basic code formats (e.g., circle SVG element), which define the exact shape, orientation, color, and placement of the objects. Given this structured way of generating an image, one of the research questions is: if a model learns an image as a composition of basic shapes, would it help a model understand the abstraction of visual concepts and structures (Ji et al., 2022a)? Moreover, would it produce more interpretable results for the users?

7.2 Reasoning

The highly impressive performance of ChatGPT has sparked interest in expanding its usage beyond traditional NLP tasks into more complex domains requiring sophisticated reasoning such as problem-solving, decision-making, and planning. Our evaluation of its reasoning abilities shows that they are not reliable. Specifically, our findings indicate that ChatGPT exhibits a tendency to be a lazy reasoner and that its capabilities are inconsistent across various reasoning abilities.

In terms of logical reasoning, ChatGPT performs better deductive and abductive reasoning compared to inductive reasoning. However, as a language model, ChatGPT still lacks the ability to answer non-textual semantic reasoning tasks, such as mathematical, temporal, and spatial reasoning. Instead, many suggest pairing ChatGPT with another computational model, such as Wolfram²⁰, to solve each specific set of problems. In that combination, ChatGPT parses natural language input into programming language code snippets, then the computational model will execute the code to return results.

²⁰<https://writings.stephenwolfram.com/2023/01/wolframalpha-as-the-way-to-bring-computational-knowledge-superpowers-to-chatgpt/>

In this way, the strength of ChatGPT is maximized while the weakness is mitigated. Meanwhile, ChatGPT surprisingly excels in commonsense, causal, and analogical reasoning. We suspect that all this knowledge has been encoded in the parametric memory of ChatGPT. Nevertheless, ChatGPT lacks the ability to perform multi-hop reasoning which suggests that, like other LLMs, ChatGPT possesses a limited ability to accomplish complex reasoning tasks.

To support the further expansion of its use cases, it is necessary to prioritize the development of systems with robust complex reasoning capabilities, which should also be facilitated by the creation of more comprehensive benchmarks for assessing these abilities, particularly when multiple abilities are required to complete the tasks.

7.3 Factuality and Hallucinations

Although powerful, ChatGPT, like other LLMs, still makes things up (Ji et al., 2022b). To ensure factuality, it is possible to build LLMs with an interface to an external knowledge source, like Blenderbot 3.0 (Shuster et al., 2022), RETRO (Borgeaud et al., 2021), and LaMDa (Thoppilan et al., 2022). In this manner, factual information LLMs can be updated independently and easily in the knowledge base, without fine-tuning the whole LLM. However, how to balance the generative power of its parametric memory with external knowledge sources is an active research area (Lee et al., 2022; He et al., 2023)

Meanwhile, there are many forms of hallucinations from LLMs that are not necessarily counterfactual but still undesirable. The RLHF process of ChatGPT can ensure human feedback to mitigate undesirable responses. However, researchers need to work on coming up with more automatic and scalable methods to detect and mitigate hallucinations and other undesirable artifacts of LLMs.

7.4 Interactivity

Compared with the previous LLMs, the interactive ability of ChatGPT has made a leap according to both qualitative and quantitative measures. Based on our evaluation, through interactivity, we can improve the performance of ChatGPT by 8% ROUGE-1 on summarization tasks and 2% ChrF++ on the machine translation tasks. However, sometimes ChatGPT retains the wrong answer even after receiving multiple rounds of prompts from the user. Improving the ability of ChatGPT to handle mul-

iple rounds of user feedback is also an important challenge.

The conversational ability and multi-turn interaction of ChatGPT make it natural for people to use it as a dialog system. We carry out the very difficult task of using ChatGPT as a task-oriented dialog system with structured knowledge given in the prompt to perform. Whereas ChatGPT shows strong performance in various modules, challenges remain for us to use ChatGPT as a fully task-oriented dialog system, due to the lack of controllability and knowledge grounding in its response.

The interactivity inadvertently enables the user to “jail-break” ChatGPT to carry out harmful actions. For example, a user could ask ChatGPT to turn off its safety layer, causing potential damage (Christian, 2023).

7.5 Responsible Generative AI

Responsible design and usage of LLMs including ChatGPT is an important and pressing challenge today. There are common issues with these models, such as fairness, toxicity, demographic bias, and safety, that need to be addressed. In the case of ChatGPT, OpenAI constructs safety layers and uses RLHF and potentially other means to filter out undesirable system responses. This process is resource intensive and opaque to the public. We hope to see a more open discussion and sharing of responsible design of LLMs from various organizations including OpenAI in the future.

References

- 2023. [Chatgpt vs satya nadella over biryani: The chatbot is learning from its mistakes.](#)
- Alham Fikri Aji, Genta Indra Winata, Fajri Koto, Samuel Cahyawijaya, Ade Romadhony, Rahmad Mahendra, Kemal Kurniawan, David Moeljadi, Radityo Eko Prasajo, Timothy Baldwin, Jey Han Lau, and Sebastian Ruder. 2022. [One country, 700+ languages: NLP challenges for underrepresented languages and dialects in Indonesia.](#) In *Proceedings of the 60th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 7226–7249, Dublin, Ireland. Association for Computational Linguistics.
- Sam Altman. 2022. [Chatgpt is incredibly limited, but good enough at some things to create a misleading impression of greatness.](#)
- Chantal Amrhein and Rico Sennrich. 2020. [On Romanization for model transfer between scripts in neural machine translation.](#) In *Findings of the Association*

- for *Computational Linguistics: EMNLP 2020*, pages 2461–2469, Online. Association for Computational Linguistics.
- Ömer Aydın and Enis Karaarslan. 2022. Openai chat-gpt generated literature review: Digital twin in healthcare. Available at SSRN 4308687.
- Satanjeev Banerjee and Alon Lavie. 2005. [METEOR: An automatic metric for MT evaluation with improved correlation with human judgments](#). In *Proceedings of the ACL Workshop on Intrinsic and Extrinsic Evaluation Measures for Machine Translation and/or Summarization*, pages 65–72, Ann Arbor, Michigan. Association for Computational Linguistics.
- Paul Bartha. 2013. [Analogy and analogical reasoning](#).
- Chandra Bhagavatula, Ronan Le Bras, Chaitanya Malaviya, Keisuke Sakaguchi, Ari Holtzman, Hannah Rashkin, Doug Downey, Wen tau Yih, and Yejin Choi. 2020. [Abductive commonsense reasoning](#). In *International Conference on Learning Representations*.
- Prajwal Bhargava and Vincent Ng. 2022. Commonsense knowledge reasoning and generation with pre-trained language models: a survey. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 36, pages 12317–12325.
- Pushpak Bhattacharyya, Rajen Chatterjee, Markus Freitag, Diptesh Kanojia, Matteo Negri, and Marco Turchi. 2022. [Findings of the wmt 2022 shared task on automatic post-editing](#). In *Proceedings of the Seventh Conference on Machine Translation*, pages 109–117, Abu Dhabi.
- David G.W. Birch. 2022. Chatgpt is a window into the real future of financial services.
- Yonatan Bisk, Rowan Zellers, Jianfeng Gao, Yejin Choi, et al. 2020. Piqa: Reasoning about physical commonsense in natural language. In *Proceedings of the AAAI conference on artificial intelligence*, volume 34, pages 7432–7439.
- Alexandre Blanco-Gonzalez, Alfonso Cabezon, Alejandro Seco-Gonzalez, Daniel Conde-Torres, Paula Antelo-Riveiro, Angel Pineiro, and Rebeca Garcia-Fandino. 2022. The role of ai in drug discovery: Challenges, opportunities, and strategies. *arXiv preprint arXiv:2212.08104*.
- Sebastian Borgeaud, Arthur Mensch, Jordan Hoffmann, Trevor Cai, Eliza Rutherford, Katie Millican, George van den Driessche, Jean-Baptiste Lespiau, Bogdan Damoc, Aidan Clark, Diego de Las Casas, Aurelia Guy, Jacob Menick, Roman Ring, Tom Hennigan, Saffron Huang, Loren Maggiore, Chris Jones, Albin Cassirer, Andy Brock, Michela Paganini, Geoffrey Irving, Oriol Vinyals, Simon Osindero, Karen Simonyan, Jack W. Rae, Erich Elsen, and Laurent Sifre. 2021. [Improving language models by retrieving from trillions of tokens](#).
- Tom Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared D Kaplan, Prafulla Dhariwal, Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, et al. 2020. Language models are few-shot learners. *Advances in neural information processing systems*, 33:1877–1901.
- Samuel Cahyawijaya, Holy Lovenia, Alham Fikri Aji, Genta Indra Winata, Bryan Wilie, Rahmad Mahendra, Christian Wibisono, Ade Romadhony, Karissa Vincentio, Fajri Koto, Jennifer Santoso, David Moeljadi, Cahya Wirawan, Frederikus Hudi, Ivan Halim Parmonangan, Ika Alfina, Muhammad Satrio Wicaksono, Ilham Firdausi Putra, Sam-sul Rahmadani, Yulianti Oenang, Ali Akbar Septiandri, James Jaya, Kaustubh D. Dhole, Arie Ardiyanti Suryani, Rifki Afina Putri, Dan Su, Keith Stevens, Made Nindyatama Nityasya, Muhammad Farid Adilazuarda, Ryan Ignatius, Ryandito Diandaru, Tiezheng Yu, Vito Ghifari, Wenliang Dai, Yan Xu, Dyah Damapusita, Cuk Tho, Ichwanul Muslim Karo Karo, Tirana Noor Fatyanosa, Ziwei Ji, Pascale Fung, Graham Neubig, Timothy Baldwin, Sebastian Ruder, Herry Sujaini, Sakriani Sakti, and Ayu Purwarianti. 2022. [Nusacrowd: Open source initiative for indonesian nlp resources](#).
- Samuel Cahyawijaya, Genta Indra Winata, Bryan Wilie, Karissa Vincentio, Xiaohong Li, Adhiguna Kuncoro, Sebastian Ruder, Zhi Yuan Lim, Syafri Bahar, Masayu Khodra, Ayu Purwarianti, and Pascale Fung. 2021. [IndoNLP: Benchmark and resources for evaluating Indonesian natural language generation](#). In *Proceedings of the 2021 Conference on Empirical Methods in Natural Language Processing*, pages 8875–8898, Online and Punta Cana, Dominican Republic. Association for Computational Linguistics.
- Ethan C. Chau and Noah A. Smith. 2021. [Specializing multilingual language models: An empirical study](#). In *Proceedings of the 1st Workshop on Multilingual Representation Learning*, pages 51–61, Punta Cana, Dominican Republic. Association for Computational Linguistics.
- Jonathan H Choi, Kristin E Hickman, Amy Monahan, and Daniel Schwarcz. 2023. Chatgpt goes to law school. Available at SSRN.
- Aakanksha Chowdhery, Sharan Narang, Jacob Devlin, Maarten Bosma, Gaurav Mishra, Adam Roberts, Paul Barham, Hyung Won Chung, Charles Sutton, Sebastian Gehrmann, Parker Schuh, Kensen Shi, Sasha Tsvyashchenko, Joshua Maynez, Abhishek Rao, Parker Barnes, Yi Tay, Noam Shazeer, Vinodkumar Prabhakaran, Emily Reif, Nan Du, Ben Hutchinson, Reiner Pope, James Bradbury, Jacob Austin, Michael Isard, Guy Gur-Ari, Pengcheng Yin, Toju Duke, Anselm Levskaya, Sanjay Ghemawat, Sunipa Dev, Henryk Michalewski, Xavier Garcia, Vedant Misra, Kevin Robinson, Liam Fedus, Denny Zhou, Daphne Ippolito, David Luan, Hyeontaek Lim, Barret Zoph, Alexander Spiridonov, Ryan Sepassi, David Dohan, Shivani Agrawal, Mark

- Omernick, Andrew M. Dai, Thanumalayan Sankaranarayanan Pillai, Marie Pellat, Aitor Lewkowycz, Erica Moreira, Rewon Child, Oleksandr Polozov, Katherine Lee, Zongwei Zhou, Xuezhi Wang, Brennan Saeta, Mark Diaz, Orhan Firat, Michele Catasta, Jason Wei, Kathy Meier-Hellstern, Douglas Eck, Jeff Dean, Slav Petrov, and Noah Fiedel. 2022. [Palm: Scaling language modeling with pathways](#).
- Jon Christian. 2023. [Amazing "jailbreak" bypasses chatgpt's ethics safeguards](#).
- Paul Christiano, Jan Leike, Tom B. Brown, Miljan Martić, Shane Legg, and Dario Amodei. 2017. [Deep reinforcement learning from human preferences](#).
- Peter Clark, Isaac Cowhey, Oren Etzioni, Tushar Khot, Ashish Sabharwal, Carissa Schoenick, and Oyvind Tafjord. 2018. [Think you have solved question answering? try arc, the ai2 reasoning challenge](#).
- Abigail C Cohn and Maya Ravindranath. 2014. Local languages in indonesia: Language maintenance or language shift. *Linguistik Indonesia*, 32(2):131–148.
- Cookup.ai. 2022. [Chatgpt - where it lacks](#).
- Wenliang Dai, Lu Hou, Lifeng Shang, Xin Jiang, Qun Liu, and Pascale Fung. 2022a. [Enabling multimodal generation on CLIP via vision-language knowledge distillation](#). In *Findings of the Association for Computational Linguistics: ACL 2022*, pages 2383–2395, Dublin, Ireland. Association for Computational Linguistics.
- Wenliang Dai, Zihan Liu, Ziwei Ji, Dan Su, and Pascale Fung. 2022b. Plausible may not be faithful: Probing object hallucination in vision-language pre-training. *ArXiv*, abs/2210.07688.
- Bhavana Dalvi, Peter Jansen, Oyvind Tafjord, Zhengnan Xie, Hannah Smith, Leighanna Pipatanangkura, and Peter Clark. 2021. Explaining answers with entailment trees. In *Proceedings of the 2021 Conference on Empirical Methods in Natural Language Processing*, pages 7358–7370.
- Ernest Davis. 2023. Mathematics, word problems, common sense, and artificial intelligence. *arXiv preprint arXiv:2301.09723*.
- Web Desk. 2023. [Colombian judge uses chatgpt in ruling, triggers debate](#).
- Igor Douven. 2017. [Abduction](#).
- Michael Dowling and Brian Lucey. 2023. Chatgpt for (finance) research: The bananarama conjecture. *Finance Research Letters*, page 103662.
- Li Du, Xiao Ding, Kai Xiong, Ting Liu, and Bing Qin. 2022. [e-CARE: a new dataset for exploring explainable causal reasoning](#). In *Proceedings of the 60th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 432–446, Dublin, Ireland. Association for Computational Linguistics.
- Esin Durmus, He He, and Mona Diab. 2020. [FEQA: A question answering evaluation framework for faithfulness assessment in abstractive summarization](#). In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, pages 5055–5070, Online. Association for Computational Linguistics.
- Nouha Dziri, Andrea Madotto, Osmar Zaiane, and Avishek Joey Bose. 2021. [Neural path hunter: Reducing hallucination in dialogue systems via path grounding](#). In *Proceedings of the 2021 Conference on Empirical Methods in Natural Language Processing*, pages 2197–2214, Online and Punta Cana, Dominican Republic. Association for Computational Linguistics.
- David M. Eberhard, Gary F. Simons, and Charles D. Fennig. 2021. *Ethnologue: Languages of the World. Twenty-fourth edition*. Dallas, Texas: SIL International.
- Simon Frieder, Luca Pinchetti, Ryan-Rhys Griffiths, Tommaso Salvatori, Thomas Lukasiewicz, Philipp Christian Petersen, Alexis Chevalier, and Julius Berner. 2023. [Mathematical capabilities of chatgpt](#).
- Leo Gao, Jonathan Tow, Stella Biderman, Sid Black, Anthony DiPofi, Charles Foster, Laurence Golding, Jeffrey Hsu, Kyle McDonell, Niklas Muennighoff, Jason Phang, Laria Reynolds, Eric Tang, Anish Thite, Ben Wang, Kevin Wang, and Andy Zou. 2021. [A framework for few-shot language model evaluation](#).
- Aidan Gilson, Conrad Safranek, Thomas Huang, Vimig Socrates, Ling Chi, Richard Andrew Taylor, and David Chartash. 2022. How well does chatgpt do when taking the medical licensing exams? the implications of large language models for medical education and knowledge assessment. *medRxiv*, pages 2022–12.
- Bogdan Gliwa, Iwona Mochol, Maciej Biesek, and Aleksander Wawer. 2019. Samsun corpus: A human-annotated dialogue dataset for abstractive summarization. *EMNLP-IJCNLP 2019*, page 70.
- Yoav Goldberg. 2023. [Some remarks on large language models](#).
- Cindy Gordon. 2023. [Chatgpt is the fastest growing app in the history of web applications](#).
- Naman Goyal, Cynthia Gao, Vishrav Chaudhary, Peng-Jen Chen, Guillaume Wenzek, Da Ju, Sanjana Krishnan, Marc’Aurelio Ranzato, Francisco Guzmán, and Angela Fan. 2021. The flores-101 evaluation benchmark for low-resource and multilingual machine translation.
- Tanya Goyal, Junyi Jessy Li, and Greg Durrett. 2022. News summarization and evaluation in the era of gpt-3. *arXiv preprint arXiv:2209.12356*.

- Roberto Gozalo-Brizuela and Eduardo C Garrido-Merchan. 2023. Chatgpt is not all you need. a state of the art review of large generative ai models. *arXiv preprint arXiv:2301.04655*.
- Barbara F Grimes. 2000. *Ethnologue*. SIL International, Dallas, TX.
- Biyang Guo, Xin Zhang, Ziyuan Wang, Minqi Jiang, Jinran Nie, Yuxuan Ding, Jianwei Yue, and Yupeng Wu. 2023. How close is chatgpt to human experts? comparison corpus, evaluation, and detection. *arXiv preprint arXiv:2301.07597*.
- James Hawthorne. 2021. Inductive Logic. In Edward N. Zalta, editor, *The Stanford Encyclopedia of Philosophy*, Spring 2021 edition. Metaphysics Research Lab, Stanford University.
- Hangfeng He, Hongming Zhang, and Dan Roth. 2023. Rethinking with retrieval: Faithful large language model inference.
- Karl Moritz Hermann, Tomas Kocisky, Edward Grefenstette, Lasse Espeholt, Will Kay, Mustafa Suleyman, and Phil Blunsom. 2015. [Teaching machines to read and comprehend](#). In *Advances in Neural Information Processing Systems*, volume 28. Curran Associates, Inc.
- Jordan Hoffmann, Sebastian Borgeaud, Arthur Mensch, Elena Buchatskaya, Trevor Cai, Eliza Rutherford, Diego de Las Casas, Lisa Anne Hendricks, Johannes Welbl, Aidan Clark, Tom Hennigan, Eric Noland, Katie Millican, George van den Driessche, Bogdan Damoc, Aurelia Guy, Simon Osindero, Karen Simonyan, Erich Elsen, Jack W. Rae, Oriol Vinyals, and Laurent Sifre. 2022. [Training compute-optimal large language models](#).
- Ehsan Hosseini-Asl, Bryan McCann, Chien-Sheng Wu, Semih Yavuz, and Richard Socher. 2020. A simple language model for task-oriented dialogue. *Advances in Neural Information Processing Systems*, 33:20179–20191.
- Krystal Hu. 2023. [Chatgpt sets record for fastest-growing user base - analyst note](#).
- Jie Huang and Kevin Chen-Chuan Chang. 2022. Towards reasoning in large language models: A survey. *arXiv preprint arXiv:2212.10403*.
- Arfinda Ilmania, Abdurrahman, Samuel Cahyawijaya, and Ayu Purwarianti. 2018. [Aspect detection and sentiment classification using deep neural network for indonesian aspect-based sentiment analysis](#). In *2018 International Conference on Asian Language Processing (IALP)*, pages 62–67.
- Hadar Yoana Jabotinsky and Roe Sarel. 2022. Co-authoring with an ai? ethical dilemmas and artificial intelligence. *Ethical Dilemmas and Artificial Intelligence (December 15, 2022)*.
- Katharina Jeblick, Balthasar Schachtner, Jakob Dexl, Andreas Mittermeier, Anna Theresa Stüber, Johanna Topalis, Tobias Weber, Philipp Wesp, Bastian Sabel, Jens Rieke, et al. 2022. Chatgpt makes medicine easy to swallow: An exploratory case study on simplified radiology reports. *arXiv preprint arXiv:2212.14882*.
- Anya Ji, Noriyuki Kojima, Noah Rush, Alane Suhr, Wai Keen Vong, Robert Hawkins, and Yoav Artzi. 2022a. [Abstract visual reasoning with tangram shapes](#). In *Proceedings of the 2022 Conference on Empirical Methods in Natural Language Processing*, pages 582–601, Abu Dhabi, United Arab Emirates. Association for Computational Linguistics.
- Ziwei Ji, Nayeon Lee, Rita Frieske, Tiezheng Yu, Dan Su, Yan Xu, Etsuko Ishii, Yejin Bang, Andrea Madotto, and Pascale Fung. 2022b. [Survey of hallucination in natural language generation](#). *ACM Comput. Surv.* Just Accepted.
- Ziwei Ji, Zihan Liu, Nayeon Lee, Tiezheng Yu, Bryan Wilie, Min Zeng, and Pascale Fung. 2022c. ρ : Reducing hallucination in open-domain dialogues with knowledge grounding. *arXiv preprint arXiv:2212.01588*.
- Wenxiang Jiao, Wenxuan Wang, Jen-tse Huang, Xing Wang, and Zhaopeng Tu. 2023. [Is chatgpt a good translator? a preliminary study](#).
- Arianna Johnson. 2023. [Is chatgpt partisan? poems about trump and biden raise questions about the ai bot’s bias-here’s what experts think](#).
- Pratik Joshi, Sebastin Santy, Amar Budhiraja, Kalika Bali, and Monojit Choudhury. 2020. [The state and fate of linguistic diversity and inclusion in the NLP world](#). In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, pages 6282–6293, Online. Association for Computational Linguistics.
- Jennifer A. Kingson. 2023. [Friend or foe? teachers debate chatgpt](#).
- Escape Velocity Labs. 2022. [Chatgpt imitates logical reasoning surprisingly well](#).
- Anton E Lawson. 2005. What is the role of induction and deduction in reasoning and scientific inquiry? *Journal of Research in Science Teaching*, 42(6):716–740.
- Nayeon Lee, Yejin Bang, Andrea Madotto, and Pascale Fung. 2021. [Towards few-shot fact-checking via perplexity](#). In *Proceedings of the 2021 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*, pages 1971–1981, Online. Association for Computational Linguistics.
- Nayeon Lee, Wei Ping, Peng Xu, Mostofa Patwary, Pascale Fung, Mohammad Shoeybi, and Bryan Catanzaro. 2022. [Factuality enhanced language models](#)

- for open-ended text generation. In *Advances in Neural Information Processing Systems*.
- Mike Paul Lewis, editor. 2009. *Ethnologue: Languages of the World*, sixteenth edition. SIL International, Dallas, TX, USA.
- Mike Lewis, Yinhan Liu, Naman Goyal, Marjan Ghazvininejad, Abdelrahman Mohamed, Omer Levy, Veselin Stoyanov, and Luke Zettlemoyer. 2020a. [BART: Denoising sequence-to-sequence pre-training for natural language generation, translation, and comprehension](#). In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, pages 7871–7880, Online. Association for Computational Linguistics.
- Mike Lewis, Yinhan Liu, Naman Goyal, Marjan Ghazvininejad, Abdelrahman Mohamed, Omer Levy, Veselin Stoyanov, and Luke Zettlemoyer. 2020b. [Bart: Denoising sequence-to-sequence pre-training for natural language generation, translation, and comprehension](#). In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, pages 7871–7880.
- Percy Liang, Rishi Bommasani, Tony Lee, Dimitris Tsipras, Dilara Soylu, Michihiro Yasunaga, Yian Zhang, Deepak Narayanan, Yuhuai Wu, Ananya Kumar, Benjamin Newman, Binhang Yuan, Bobby Yan, Ce Zhang, Christian Cosgrove, Christopher D. Manning, Christopher Ré, Diana Acosta-Navas, Drew A. Hudson, Eric Zelikman, Esin Durmus, Faisal Ladhak, Frieda Rong, Hongyu Ren, Huaxiu Yao, Jue Wang, Keshav Santhanam, Laurel Orr, Lucia Zheng, Mert Yuksekogonul, Mirac Suzgun, Nathan Kim, Neel Guha, Niladri Chatterji, Omar Khattab, Peter Henderson, Qian Huang, Ryan Chi, Sang Michael Xie, Shibani Santurkar, Surya Ganguli, Tatsunori Hashimoto, Thomas Icard, Tianyi Zhang, Vishrav Chaudhary, William Wang, Xuechen Li, Yifan Mai, Yuhui Zhang, and Yuta Koreeda. 2022. [Holistic evaluation of language models](#).
- Opher Lieber, Or Sharir, Barak Lenz, and Yoav Shoham. 2021. Jurassic-1: Technical details and evaluation. *White Paper. AI21 Labs*.
- Stephanie Lin, Jacob Hilton, and Owain Evans. 2022. [TruthfulQA: Measuring how models mimic human falsehoods](#). In *Proceedings of the 60th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 3214–3252, Dublin, Ireland. Association for Computational Linguistics.
- Zhaojiang Lin, Bing Liu, Andrea Madotto, Seungwhan Moon, Zhenpeng Zhou, Paul A Crook, Zhiguang Wang, Zhou Yu, Eunjoon Cho, Rajen Subba, et al. 2021. Zero-shot dialogue state tracking via cross-task transfer. In *Proceedings of the 2021 Conference on Empirical Methods in Natural Language Processing*, pages 7890–7900.
- Yixin Liu, Pengfei Liu, Dragomir Radev, and Graham Neubig. 2022. [Brio: Bringing order to abstrac-](#)
- [tive summarization](#). In *Proceedings of the 60th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 2890–2903.
- Holy Lovenia, Bryan Wilie, Romain Barraud, Samuel Cahyawijaya, Willy Chung, and Pascale Fung. 2022. [Every picture tells a story: Image-grounded controllable stylistic story generation](#). In *Proceedings of the 6th Joint SIGHUM Workshop on Computational Linguistics for Cultural Heritage, Social Sciences, Humanities and Literature*, pages 40–52, Gyeongju, Republic of Korea. International Conference on Computational Linguistics.
- Hongyuan Lu, Haoyang Huang, Shuming Ma, Dongdong Zhang, Wai Lam, and Furu Wei. 2022. [Trip: Triangular document-level pre-training for multilingual language models](#). *arXiv preprint arXiv:2212.07752*.
- Andrea Madotto, Zhaojiang Lin, Genta Indra Winata, and Pascale Fung. 2021. [Few-shot bot: Prompt-based learning for dialogue systems](#). *arXiv preprint arXiv:2110.08118*.
- Kyle Mahowald, Anna A Ivanova, Idan A Blank, Nancy Kanwisher, Joshua B Tenenbaum, and Evelina Fedorenko. 2023. [Dissociating language and thought in large language models: a cognitive perspective](#). *arXiv preprint arXiv:2301.06627*.
- Bernard Marr. 2022. [What does chatgpt really mean for businesses?](#)
- Vaibhav Mavi, Anubhav Jangra, and Adam Jatowt. 2022. [A survey on multi-hop question answering and generation](#). *arXiv preprint arXiv:2204.09140*.
- Pasquale Minervini, Sebastian Riedel, Pontus Stenertorp, Edward Grefenstette, and Tim Rocktäschel. 2020. [Learning reasoning strategies in end-to-end differentiable proving](#). In *Proceedings of the 37th International Conference on Machine Learning, ICML’20*. JMLR.org.
- Roshanak Mirzaee and Parisa Kordjamshidi. 2022. [Transfer learning with synthetic corpora for spatial role labeling and reasoning](#). In *Proceedings of the 2022 Conference on Empirical Methods in Natural Language Processing*, pages 6148–6165, Abu Dhabi, United Arab Emirates. Association for Computational Linguistics.
- Roshanak Mirzaee, Hossein Rajaby Faghihi, Qiang Ning, and Parisa Kordjamshidi. 2021. [SPARTQA: A textual question answering benchmark for spatial reasoning](#). In *Proceedings of the 2021 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*, pages 4582–4598, Online. Association for Computational Linguistics.
- Seungwhan Moon, Pararth Shah, Anuj Kumar, and Rajen Subba. 2019. [Opendialkg: Explainable conversational reasoning with attention-based walks over](#)

- knowledge graphs. In *Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics*, pages 845–854.
- Nasrin Mostafazadeh, Chris Brockett, William B Dolan, Michel Galley, Jianfeng Gao, Georgios Spithourakis, and Lucy Vanderwende. 2017. Image-grounded conversations: Multimodal context for natural question and response generation. In *Proceedings of the Eighth International Joint Conference on Natural Language Processing (Volume 1: Long Papers)*, pages 462–472.
- Niklas Muennighoff, Thomas Wang, Lintang Sutawika, Adam Roberts, Stella Biderman, Teven Le Scao, M Saiful Bari, Sheng Shen, Zheng-Xin Yong, Hailey Schoelkopf, Xiangru Tang, Dragomir Radev, Alham Fikri Aji, Khalid Almubarak, Samuel Albanie, Zaid Alyafeai, Albert Webson, Edward Raff, and Colin Raffel. 2022. Crosslingual generalization through multitask finetuning. *arXiv preprint arXiv:2211.01786*.
- Benjamin Muller, Antonios Anastasopoulos, Benoît Sagot, and Djamel Seddah. 2021. [When being unseen from mBERT is just the beginning: Handling new languages with multilingual language models](#). In *Proceedings of the 2021 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*, pages 448–462, Online. Association for Computational Linguistics.
- Ramesh Nallapati, Bowen Zhou, Cicero dos Santos, Çağlar Gulçehre, and Bing Xiang. 2016. [Abstractive text summarization using sequence-to-sequence RNNs and beyond](#). In *Proceedings of the 20th SIGNLL Conference on Computational Natural Language Learning*, pages 280–290, Berlin, Germany. Association for Computational Linguistics.
- Tomáš Nekvinda and Ondřej Dušek. 2021. [Shades of BLEU, flavours of success: The case of MultiWOZ](#). In *Proceedings of the 1st Workshop on Natural Language Generation, Evaluation, and Metrics (GEM 2021)*, pages 34–46, Online. Association for Computational Linguistics.
- NeuralMagic. 2023. [The chatgpt cheat sheet](#).
- Oded Nov, Nina Singh, and Devin M Mann. 2023. Putting chatgpt’s medical advice to the (turing) test. *medRxiv*, pages 2023–01.
- Jeroen Ooms. 2023. *cld2: Google’s Compact Language Detector 2*. <https://docs.ropensci.org/cld2/> (docs) <https://github.com/ropensci/cld2> (devel) <https://github.com/cld2owners/cld2> (upstream).
- Simon Ott, Konstantin Hebenstreit, Valentin Liévin, Christoffer Egeberg Hother, Milad Moradi, Maximilian Mayrhauser, Robert Praas, Ole Winther, and Matthias Samwald. 2023. Thoughtsource: A central hub for large language model reasoning data. *arXiv preprint arXiv:2301.11596*.
- Long Ouyang, Jeff Wu, Xu Jiang, Diogo Almeida, Carroll L. Wainwright, Pamela Mishkin, Chong Zhang, Sandhini Agarwal, Katarina Slama, Alex Ray, John Schulman, Jacob Hilton, Fraser Kelton, Luke Miller, Maddie Simens, Amanda Askell, Peter Welinder, Paul Christiano, Jan Leike, and Ryan Lowe. 2022. [Training language models to follow instructions with human feedback](#).
- Baolin Peng, Chunyuan Li, Jinchao Li, Shahin Shayan-deh, Lars Liden, and Jianfeng Gao. 2021. Soloist: Building task bots at scale with transfer learning and machine teaching. *Transactions of the Association for Computational Linguistics*, 9:807–824.
- Jonas Pfeiffer, Ivan Vulić, Iryna Gurevych, and Sebastian Ruder. 2021. [UNKs everywhere: Adapting multilingual language models to new scripts](#). In *Proceedings of the 2021 Conference on Empirical Methods in Natural Language Processing*, pages 10186–10203, Online and Punta Cana, Dominican Republic. Association for Computational Linguistics.
- Maja Popović. 2015. [chrF: character n-gram F-score for automatic MT evaluation](#). In *Proceedings of the Tenth Workshop on Statistical Machine Translation*, pages 392–395, Lisbon, Portugal. Association for Computational Linguistics.
- Ian Porada, Kaheer Suleman, Adam Trischler, and Jackie Chi Kit Cheung. 2021. [Modeling event plausibility with consistent conceptual abstraction](#). In *Proceedings of the 2021 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*, pages 1732–1743, Online. Association for Computational Linguistics.
- Matt Post. 2018. [A call for clarity in reporting BLEU scores](#). In *Proceedings of the Third Conference on Machine Translation: Research Papers*, pages 186–191, Brussels, Belgium. Association for Computational Linguistics.
- Shuofei Qiao, Yixin Ou, Ningyu Zhang, Xiang Chen, Yunzhi Yao, Shumin Deng, Chuanqi Tan, Fei Huang, and Huajun Chen. 2022. Reasoning with language model prompting: A survey. *arXiv preprint arXiv:2212.09597*.
- Lianhui Qin, Aditya Gupta, Shyam Upadhyay, Luheng He, Yejin Choi, and Manaal Faruqui. 2021. [TIME-DIAL: Temporal commonsense reasoning in dialog](#). In *Proceedings of the 59th Annual Meeting of the Association for Computational Linguistics and the 11th International Joint Conference on Natural Language Processing (Volume 1: Long Papers)*, pages 7066–7076, Online. Association for Computational Linguistics.
- Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya Ramesh, Gabriel Goh, Sandhini Agarwal, Girish Sastry, Amanda Askell, Pamela Mishkin, Jack Clark, et al. 2021. Learning transferable visual models from natural language supervision. In *International*

- conference on machine learning*, pages 8748–8763. PMLR.
- Alec Radford, Jeffrey Wu, Rewon Child, David Luan, Dario Amodei, Ilya Sutskever, et al. 2019. Language models are unsupervised multitask learners. *OpenAI blog*, 1(8):9.
- Jack W. Rae, Sebastian Borgeaud, Trevor Cai, Katie Millican, Jordan Hoffmann, Francis Song, John Aslanides, Sarah Henderson, Roman Ring, Susannah Young, Eliza Rutherford, Tom Hennigan, Jacob Menick, Albin Cassirer, Richard Powell, George van den Driessche, Lisa Anne Hendricks, Maribeth Rauh, Po-Sen Huang, Amelia Glaese, Johannes Welbl, Sumanth Dathathri, Saffron Huang, Jonathan Uesato, John Mellor, Irina Higgins, Antonia Creswell, Nat McAleese, Amy Wu, Erich Elsen, Siddhant Jayakumar, Elena Buchatskaya, David Budden, Esme Sutherland, Karen Simonyan, Michela Paganini, Laurent Sifre, Lena Martens, Xiang Lorraine Li, Adhiguna Kuncoro, Aida Nematzadeh, Elena Gribovskaya, Domenic Donato, Angeliki Lazaridou, Arthur Mensch, Jean-Baptiste Lespiau, Maria Tsimpoukelli, Nikolai Grigorev, Doug Fritz, Thibault Sottiaux, Mantas Pajarskas, Toby Pohlen, Zhitao Gong, Daniel Toyama, Cyprien de Masson d’Autume, Yujia Li, Tayfun Terzi, Vladimir Mikulik, Igor Babuschkin, Aidan Clark, Diego de Las Casas, Aurelia Guy, Chris Jones, James Bradbury, Matthew Johnson, Blake Hechtman, Laura Weidinger, Iason Gabriel, William Isaac, Ed Lockhart, Simon Osindero, Laura Rimell, Chris Dyer, Oriol Vinyals, Kareem Ayoub, Jeff Stanley, Lorraine Bennett, Demis Hassabis, Koray Kavukcuoglu, and Geoffrey Irving. 2021a. [Scaling language models: Methods, analysis and insights from training gopher](#).
- Jack W. Rae, Sebastian Borgeaud, Trevor Cai, Katie Millican, Jordan Hoffmann, Francis Song, John Aslanides, Sarah Henderson, Roman Ring, Susannah Young, et al. 2021b. [Scaling language models: Methods, analysis & insights from training gopher](#). *arXiv preprint arXiv:2112.11446*.
- Colin Raffel, Noam Shazeer, Adam Roberts, Katherine Lee, Sharan Narang, Michael Matena, Yanqi Zhou, Wei Li, and Peter J. Liu. 2022. Exploring the limits of transfer learning with a unified text-to-text transformer. *J. Mach. Learn. Res.*, 21(1).
- Aditya Ramesh, Mikhail Pavlov, Gabriel Goh, Scott Gray, Chelsea Voss, Alec Radford, Mark Chen, and Ilya Sutskever. 2021. Zero-shot text-to-image generation. In *International Conference on Machine Learning*, pages 8821–8831. PMLR.
- Fabian Rasheed. 2020. [Gpt3 sees](#).
- Anna Rogers, Matt Gardner, and Isabelle Augenstein. 2022. [Qa dataset explosion: A taxonomy of nlp resources for question answering and reading comprehension](#). *ACM Comput. Surv.* Just Accepted.
- Robin Rombach, Andreas Blattmann, Dominik Lorenz, Patrick Esser, and Björn Ommer. 2022. High-resolution image synthesis with latent diffusion models. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 10684–10695.
- David Saxton, Edward Grefenstette, Felix Hill, and Pushmeet Kohli. 2019. [Analysing mathematical reasoning abilities of neural models](#). In *International Conference on Learning Representations*.
- J Schulman, B Zoph, C Kim, J Hilton, J Menick, J Weng, JFC Uribe, L Fedus, L Metz, M Pokorny, et al. 2022. [Chatgpt: Optimizing language models for dialogue](#).
- Stephen Shankland. 2023. [Why the chatgpt ai chatbot is blowing everyone’s mind](#).
- Yiqiu Shen, Laura Heacock, Jonathan Elias, Keith D Hentel, Beatriu Reig, George Shih, and Linda Moy. 2023. [Chatgpt and other large language models are double-edged swords](#).
- Zhengxiang Shi, Qiang Zhang, and Aldo Lipani. 2022a. [Stepgame: A new benchmark for robust multi-hop spatial reasoning in texts](#). *Proceedings of the AAAI Conference on Artificial Intelligence*, 36(10):11321–11329.
- Zhengxiang Shi, Qiang Zhang, and Aldo Lipani. 2022b. [StepGame: A new benchmark for robust multi-hop spatial reasoning in texts](#). *Proceedings of the AAAI Conference on Artificial Intelligence*, 36(10):11321–11329.
- Denis Shiryayev. 2022. [Drawing mona lisa with chatgpt](#).
- Mohammad Shoeybi, Mostofa Patwary, Raul Puri, Patrick LeGresley, Jared Casper, and Bryan Catanzaro. 2019. Megatron-lm: Training multi-billion parameter language models using model parallelism. *arXiv preprint arXiv:1909.08053*.
- Kurt Shuster, Jing Xu, Mojtaba Komeili, Da Ju, Eric Michael Smith, Stephen Roller, Megan Ung, Moya Chen, Kushal Arora, Joshua Lane, Morteza Behrooz, William Ngan, Spencer Poff, Naman Goyal, Arthur Szlam, Y-Lan Boureau, Melanie Kam-badur, and Jason Weston. 2022. [Blenderbot 3: a deployed conversational agent that continually learns to responsibly engage](#).
- Koustuv Sinha, Shagun Sodhani, Jin Dong, Joelle Pineau, and William L Hamilton. 2019. Clutrr: A diagnostic benchmark for inductive reasoning from text. 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 4506–4515.
- Noah Smith. 2023. [Why does chatgpt constantly lie?](#)

- Aarohi Srivastava, Abhinav Rastogi, Abhishek Rao, Abu Awal Md Shoeb, Abubakar Abid, Adam Fisch, Adam R Brown, Adam Santoro, Aditya Gupta, Adrià Garriga-Alonso, et al. 2022. Beyond the imitation game: Quantifying and extrapolating the capabilities of language models. *arXiv preprint arXiv:2206.04615*.
- Shane Storks, Qiaozi Gao, and Joyce Y Chai. 2019. Commonsense reasoning for natural language understanding: A survey of benchmarks, resources, and approaches. *arXiv preprint arXiv:1904.01172*, pages 1–60.
- Dan Su, Xiaoguang Li, Jindi Zhang, Lifeng Shang, Xin Jiang, Qun Liu, and Pascale Fung. 2022. Read before generate! faithful long form question answering with machine reading. In *Findings of the Association for Computational Linguistics: ACL 2022*, pages 744–756.
- Dan Su, Tiezheng Yu, and Pascale Fung. 2021. Improve query focused abstractive summarization by incorporating answer relevance. In *Findings of the Association for Computational Linguistics: ACL-IJCNLP 2021*, pages 3124–3131.
- Qingfeng Sun, Yujing Wang, Can Xu, Kai Zheng, Yaming Yang, Huang Hu, Fei Xu, Jessica Zhang, Xiubo Geng, and Daxin Jiang. 2022. Multimodal dialogue response generation. In *Proceedings of the 60th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 2854–2866.
- Teo Susnjak. 2022. Chatgpt: The end of online exam integrity? *arXiv preprint arXiv:2212.09292*.
- Alon Talmor, Yanai Elazar, Yoav Goldberg, and Jonathan Berant. 2020. olmpics-on what language model pre-training captures. *Transactions of the Association for Computational Linguistics*, 8:743–758.
- Alon Talmor, Jonathan Herzig, Nicholas Lourie, and Jonathan Berant. 2018. Commonsenseqa: A question answering challenge targeting commonsense knowledge. *arXiv preprint arXiv:1811.00937*.
- NLLB Team, Marta R. Costa-jussà, James Cross, Onur Çelebi, Maha Elbayad, Kenneth Heafield, Kevin Heffernan, Elahe Kalbassi, Janice Lam, Daniel Licht, Jean Maillard, Anna Sun, Skyler Wang, Guillaume Wenzek, Al Youngblood, Bapi Akula, Loic Barrault, Gabriel Mejia Gonzalez, Prangthip Hansanti, John Hoffman, Semarley Jarrett, Kaushik Ram Sadagopan, Dirk Rowe, Shannon Spruit, Chau Tran, Pierre Andrews, Necip Fazil Ayan, Shruti Bhosale, Sergey Edunov, Angela Fan, Cynthia Gao, Vedanuj Goswami, Francisco Guzmán, Philipp Koehn, Alexandre Mourachko, Christophe Ropers, Safiyyah Saleem, Holger Schwenk, and Jeff Wang. 2022. No language left behind: Scaling human-centered machine translation.
- Richmond Thomason. 2018. [Logic and artificial intelligence](#).
- Romal Thoppilan, Daniel De Freitas, Jamie Hall, Noam Shazeer, Apoorv Kulshreshtha, Heng-Tze Cheng, Alicia Jin, Taylor Bos, Leslie Baker, Yu Du, et al. 2022. Lamda: Language models for dialog applications. *arXiv preprint arXiv:2201.08239*.
- H Holden Thorp. 2023. Chatgpt is fun, but not an author.
- Giuseppe Venuto. 2023. [Giuven95/chatgpt-failures: Chatgpt failure archive](#).
- Douglas Walton. 2014. *Abductive reasoning*. University of Alabama Press.
- Ada Wan. 2022. [Fairness in representation for multilingual NLP: Insights from controlled experiments on conditional language modeling](#). In *International Conference on Learning Representations*.
- Alex Wang, Amanpreet Singh, Julian Michael, Felix Hill, Omer Levy, and Samuel Bowman. 2018a. [GLUE: A multi-task benchmark and analysis platform for natural language understanding](#). In *Proceedings of the 2018 EMNLP Workshop BlackboxNLP: Analyzing and Interpreting Neural Networks for NLP*, pages 353–355, Brussels, Belgium. Association for Computational Linguistics.
- Su Wang, Greg Durrett, and Katrin Erk. 2018b. Modeling semantic plausibility by injecting world knowledge. *arXiv preprint arXiv:1804.00619*.
- Xuezhi Wang, Jason Wei, Dale Schuurmans, Quoc Le, Ed Chi, and Denny Zhou. 2022. Self-consistency improves chain of thought reasoning in language models. *arXiv preprint arXiv:2203.11171*.
- Peter Cathcart Wason and Philip Nicholas Johnson-Laird. 1972. *Psychology of reasoning: Structure and content*, volume 86. Harvard University Press.
- Taylor Webb, Keith J. Holyoak, and Hongjing Lu. 2022a. [Emergent analogical reasoning in large language models](#).
- Taylor Webb, Keith J Holyoak, and Hongjing Lu. 2022b. Emergent analogical reasoning in large language models. *arXiv preprint arXiv:2212.09196*.
- Jason Wei, Yi Tay, Rishi Bommasani, Colin Raffel, Barret Zoph, Sebastian Borgeaud, Dani Yogatama, Maarten Bosma, Denny Zhou, Donald Metzler, et al. 2022a. Emergent abilities of large language models. *Transactions on Machine Learning Research*.
- Jason Wei, Xuezhi Wang, Dale Schuurmans, Maarten Bosma, Fei Xia, Ed H Chi, Quoc V Le, Denny Zhou, et al. 2022b. Chain-of-thought prompting elicits reasoning in large language models. In *Advances in Neural Information Processing Systems*.
- Jason Weston, Antoine Bordes, Sumit Chopra, and Tomáš Mikolov. 2016a. [Towards ai-complete question answering: A set of prerequisite toy tasks](#). In *4th International Conference on Learning Representations, ICLR 2016, San Juan, Puerto Rico, May 2-4, 2016, Conference Track Proceedings*.

- Jason Weston, Antoine Bordes, Sumit Chopra, Alexander M Rush, Bart Van Merriënboer, Armand Joulin, and Tomas Mikolov. 2016b. Towards ai-complete question answering: A set of prerequisite toy tasks. In *4th International Conference on Learning Representations, ICLR 2016*.
- Bryan Wilie, Karissa Vincentio, Genta Indra Winata, Samuel Cahyawijaya, Xiaohong Li, Zhi Yuan Lim, Sidik Soleman, Rahmad Mahendra, Pascale Fung, Syafri Bahar, et al. 2020. Indonlu: Benchmark and resources for evaluating indonesian natural language understanding. In *Proceedings of the 1st Conference of the Asia-Pacific Chapter of the Association for Computational Linguistics and the 10th International Joint Conference on Natural Language Processing*, pages 843–857.
- Genta Indra Winata, Alham Fikri Aji, Samuel Cahyawijaya, Rahmad Mahendra, Fajri Koto, Ade Romadhony, Kemal Kurniawan, David Moeljadi, Radityo Eko Prasajo, Pascale Fung, Timothy Baldwin, Jey Han Lau, Rico Sennrich, and Sebastian Ruder. 2022. [Nusax: Multilingual parallel sentiment dataset for 10 indonesian local languages](#).
- Cameron R. Wolfe. 2023. [Specialized llms: Chatgpt, lamda, galactica, codex, sparrow, and more](#).
- BigScience Workshop, :, Teven Le Scao, Angela Fan, Christopher Akiki, Ellie Pavlick, Suzana Ilić, Daniel Hesslow, Roman Castagné, Alexandra Sasha Luccioni, François Yvon, Matthias Gallé, Jonathan Tow, Alexander M. Rush, Stella Biderman, Albert Webson, Pawan Sasanka Ammanamanchi, Thomas Wang, Benoît Sagot, Niklas Muennighoff, Albert Villanova del Moral, Olatunji Ruwase, Rachel Bawden, Stas Bekman, Angelina McMillan-Major, Iz Beltagy, Huu Nguyen, Lucile Saulnier, Samson Tan, Pedro Ortiz Suarez, Victor Sanh, Hugo Laurençon, Yacine Jernite, Julien Launay, Margaret Mitchell, Colin Raffel, Aaron Gokaslan, Adi Simhi, Aitor Soroa, Alham Fikri Aji, Amit Alfassy, Anna Rogers, Ariel Kreisberg Nitzav, Canwen Xu, Chenghao Mou, Chris Emezue, Christopher Klamm, Colin Leong, Daniel van Strien, David Ifeoluwa Adelani, Dragomir Radev, Eduardo González Ponferrada, Efrat Levkovizh, Ethan Kim, Eyal Bar Natan, Francesco De Toni, Gérard Dupont, Germán Kruszewski, Giada Pistilli, Hady Elsahar, Hamza Benyamina, Hieu Tran, Ian Yu, Idris Abdulmumin, Isaac Johnson, Itziar Gonzalez-Dios, Javier de la Rosa, Jenny Chim, Jesse Dodge, Jian Zhu, Jonathan Chang, Jörg Froberg, Joseph Tobing, Joydeep Bhattacharjee, Khalid Almubarak, Kimbo Chen, Kyle Lo, Leandro Von Werra, Leon Weber, Long Phan, Loubna Ben allal, Ludovic Tanguy, Manan Dey, Manuel Romero Muñoz, Maraim Masoud, María Grandury, Mario Šaško, Max Huang, Maximin Coavoux, Mayank Singh, Mike Tian-Jian Jiang, Minh Chien Vu, Mohammad A. Jauhar, Mustafa Ghaleb, Nishant Subramani, Nora Kassner, Nuru-laqilla Khamis, Olivier Nguyen, Omar Espejel, Ona de Gibert, Paulo Villegas, Peter Henderson, Pierre Colombo, Priscilla Amuok, Quentin Lhoest, Rheza Harliman, Rishi Bommasani, Roberto Luis López, Rui Ribeiro, Salomey Osei, Sampo Pyysalo, Sebastian Nagel, Shamik Bose, Shamsuddeen Hassan Muhammad, Shanya Sharma, Shayne Longpre, Somaieh Nikpoor, Stanislav Silberberg, Suhas Pai, Sydney Zink, Tiago Timponi Torrent, Timo Schick, Tristan Thrush, Valentin Danchev, Vasilina Nikoulina, Veronika Laippala, Violette Lecerq, Vrinda Prabhu, Zaid Alyafeai, Zeerak Talat, Arun Raja, Benjamin Heinzerling, Chenglei Si, Davut Emre Taşar, Elizabeth Salesky, Sabrina J. Mielke, Wilson Y. Lee, Abheesht Sharma, Andrea Santilli, Antoine Chaffin, Arnaud Stiegler, Debajyoti Datta, Eliza Szczechla, Gunjan Chhablani, Han Wang, Harshit Pandey, Hendrik Strobelt, Jason Alan Fries, Jos Rozen, Leo Gao, Lintang Sutawika, M Saiful Bari, Maged S. Al-shaibani, Matteo Manica, Nihal Nayak, Ryan Teehan, Samuel Albanie, Sheng Shen, Srulik Ben-David, Stephen H. Bach, Taewoon Kim, Tali Bers, Thibault Fevry, Trishala Neeraj, Urmish Thakker, Vikas Raunak, Xiangru Tang, Zheng-Xin Yong, Zhiqing Sun, Shaked Brody, Yallow Uri, Hadar Tojarieh, Adam Roberts, Hyung Won Chung, Jaesung Tae, Jason Phang, Ofir Press, Conglong Li, Deepak Narayanan, Hatim Bourfoune, Jared Casper, Jeff Rasley, Max Ryabinin, Mayank Mishra, Minjia Zhang, Mohammad Shoeybi, Myriam Peyrounette, Nicolas Patry, Nouamane Tazi, Omar Sanseviero, Patrick von Platen, Pierre Cornette, Pierre François Lavalée, Rémi Lacroix, Samyam Rajbhandari, Sanchit Gandhi, Shaden Smith, Stéphane Requena, Suraj Patil, Tim Dettmers, Ahmed Baruwā, Amanpreet Singh, Anastasia Cheveleva, Anne-Laure Ligozat, Arjun Subramonian, Aurélie Névéal, Charles Lovering, Dan Garrette, Deepak Tunuguntla, Ehud Reiter, Ekaterina Taktasheva, Ekaterina Voloshina, Eli Bogdanov, Genta Indra Winata, Hailey Schoelkopf, Jan-Christoph Kalo, Jekaterina Novikova, Jessica Zosa Forde, Jordan Clive, Jungo Kasai, Ken Kawamura, Liam Hazan, Marine Carpuat, Miruna Clinciu, Najoung Kim, Newton Cheng, Oleg Serikov, Omer Antverg, Oskar van der Wal, Rui Zhang, Ruochen Zhang, Sebastian Gehrmann, Shachar Mirkin, Shani Pais, Tatiana Shavrina, Thomas Scialom, Tian Yun, Tomasz Limisiewicz, Verena Rieser, Vitaly Protasov, Vladislav Mikhailov, Yada Pruksachatkun, Yonatan Belinkov, Zachary Bamberger, Zdeněk Kasner, Alice Rueda, Amanda Pestana, Amir Feizpour, Ammar Khan, Amy Faranak, Ana Santos, Anthony Hevia, Antigona Unldreaj, Arash Aghagol, Arezoo Abdollahi, Aycha Tammour, Azadeh HajiHosseini, Bahareh Behroozi, Benjamin Ajibade, Bharat Saxena, Carlos Muñoz Ferrandis, Danish Contractor, David Lansky, Davis David, Douwe Kiela, Duong A. Nguyen, Edward Tan, Emi Baylor, Ezinwanne Ozoani, Fatima Mirza, Frankline Ononiwu, Habib Rezanejad, Hessie Jones, Indrani Bhattacharya, Irene Solaiman, Irina Sedenko, Isar Nejadgholi, Jesse Passmore, Josh Seltzer, Julio Bonis Sanz, Livia Dutra, Mairon Samagaio, Maraim Elbadri, Margot Mieskes, Marissa Gerchick, Martha Akin-

- Iolu, Michael McKenna, Mike Qiu, Muhammed Ghauri, Mykola Burynok, Nafis Abrar, Nazneen Rajani, Nour Elkott, Nour Fahmy, Olanrewaju Samuel, Ran An, Rasmus Kromann, Ryan Hao, Samira Alizadeh, Sarmad Shubber, Silas Wang, Sourav Roy, Sylvain Viguier, Thanh Le, Tobi Oyebade, Trieu Le, Yoyo Yang, Zach Nguyen, Abhinav Ramesh Kashyap, Alfredo Palasciano, Alison Callahan, Anima Shukla, Antonio Miranda-Escalada, Ayush Singh, Benjamin Beilharz, Bo Wang, Caio Brito, Chenxi Zhou, Chirag Jain, Chuxin Xu, Clémentine Fourrier, Daniel León Perriñán, Daniel Molano, Dian Yu, Enrique Manjavacas, Fabio Barth, Florian Fuhrmann, Gabriel Altay, Giyaseddin Bayrak, Gully Burns, Helena U. Vrabec, Imane Bello, Ishani Dash, Ji Hyun Kang, John Giorgi, Jonas Golde, Jose David Posada, Karthik Rangasai Sivarman, Lokesh Bulchandani, Lu Liu, Luisa Shinzato, Madeleine Hahn de Bykhovetz, Maiko Takeuchi, Marc Pàmies, Maria A Castillo, Marianna Nezhurina, Mario Sànger, Matthias Samwald, Michael Cullan, Michael Weinberg, Michiel De Wolf, Mina Mihaljcic, Minna Liu, Moritz Freidank, Myungsun Kang, Natasha Seelam, Nathan Dahlberg, Nicholas Michio Broad, Nikolaus Muellner, Pascale Fung, Patrick Haller, Ramya Chandrasekhar, Renata Eisenberg, Robert Martin, Rodrigo Canalli, Rosaline Su, Ruisi Su, Samuel Cahyawijaya, Samuele Garda, Shlok S Deshmukh, Shubhanshu Mishra, Sid Kiblawi, Simon Ott, Sinee Sang-aaroonsiri, Srishti Kumar, Stefan Schweter, Sushil Bharati, Tanmay Laud, Théo Gigant, Tomoya Kainuma, Wojciech Kusa, Yanis Labrak, Yash Shailesh Bajaj, Yash Venkatraman, Yifan Xu, Yingxin Xu, Yu Xu, Zhe Tan, Zhongli Xie, Zifan Ye, Mathilde Bras, Younes Belkada, and Thomas Wolf. 2022. [Bloom: A 176b-parameter open-access multilingual language model](#).
- Yan Xu, Etsuko Ishii, Samuel Cahyawijaya, Zihan Liu, Genta Indra Winata, Andrea Madotto, Dan Su, and Pascale Fung. 2022. Retrieval-free knowledge-grounded dialogue response generation with adapters. In *Proceedings of the Second DiADoc Workshop on Document-grounded Dialogue and Conversational Question Answering*, pages 93–107.
- Yunyi Yang, Yunhao Li, and Xiaojun Quan. 2021. [Ubar: Towards fully end-to-end task-oriented dialog system with gpt-2](#). *Proceedings of the AAAI Conference on Artificial Intelligence*, 35(16):14230–14238.
- Zhilin Yang, Peng Qi, Saizheng Zhang, Yoshua Bengio, William Cohen, Ruslan Salakhutdinov, and Christopher D Manning. 2018. Hotpotqa: A dataset for diverse, explainable multi-hop question answering. In *Proceedings of the 2018 Conference on Empirical Methods in Natural Language Processing*, pages 2369–2380.
- Tiezheng Yu, Wenliang Dai, Zihan Liu, and Pascale Fung. 2021a. [Vision guided generative pre-trained language models for multimodal abstractive summarization](#). In *Proceedings of the 2021 Conference on Empirical Methods in Natural Language Processing*, pages 3995–4007, Online and Punta Cana, Dominican Republic. Association for Computational Linguistics.
- Tiezheng Yu, Zihan Liu, and Pascale Fung. 2021b. Adaptsum: Towards low-resource domain adaptation for abstractive summarization. In *Proceedings of the 2021 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*, pages 5892–5904.
- Xiaoxue Zang, Abhinav Rastogi, Srinivas Sunkara, Raghav Gupta, Jianguo Zhang, and Jindong Chen. 2020. Multiwoz 2.2: A dialogue dataset with additional annotation corrections and state tracking baselines. *ACL 2020*, page 109.
- Eric Zelikman, Yuhuai Wu, Jesse Mu, and Noah Goodman. 2022. Star: Bootstrapping reasoning with reasoning. In *Advances in Neural Information Processing Systems*.
- Susan Zhang, Stephen Roller, Naman Goyal, Mikel Artetxe, Moya Chen, Shuohui Chen, Christopher Dewan, Mona Diab, Xian Li, Xi Victoria Lin, et al. 2022. Opt: Open pre-trained transformer language models. *arXiv preprint arXiv:2205.01068*.
- 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*.
- Jeffrey Zhao, Raghav Gupta, Yuanbin Cao, Dian Yu, Mingqiu Wang, Harrison Lee, Abhinav Rastogi, Izhak Shafran, and Yonghui Wu. 2022. Description-driven task-oriented dialog modeling. *ArXiv*, abs/2201.08904.
- Xueliang Zhao, Wei Wu, Can Xu, Chongyang Tao, Dongyan Zhao, and Rui Yan. 2020. Knowledge-grounded dialogue generation with pre-trained language models. In *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP)*, pages 3377–3390.
- Terry Yue Zhuo, Yujin Huang, Chunyang Chen, and Zhenchang Xing. 2023. Exploring ai ethics of chatgpt: A diagnostic analysis. *arXiv preprint arXiv:2301.12867*.

A Flag Drawing Task Results

We provide the detailed results of the flag drawing task described in §3.3.1 in Figure 7.

No	Country/Region	Ground truth	Turn 1 (without description)			Turn 1			Turn 2			Turn 3			End result
			Grade	L/C/M/S	Image	Grade	L/C/M/S	Image	Grade	L/C/M/S	Image	Grade	L/C/M/S	Image	
1	United States		E	1/1/1/1		D	1/0/1/1		D	1/0/1/1		D	1/0/1/1		D
2	Canada		D	1/0/1/1		B	0/0/0/1		D	1/0/1/1		B	0/0/0/1		B
3	Brazil		E	1/1/1/1		E	0/1/0/1		E	0/1/1/1		E	0/1/1/1		E
4	Mexico		E	1/1/1/1		D	1/0/0/1		D	1/0/0/1		E	1/0/0/1		E
5	Argentina		E	0/1/1/1		E	1/1/1/1		D	1/1/0/1		C	0/1/0/1		C
6	Colombia		E	1/1/0/0		C	0/1/0/1		B	0/1/0/0		A	0/0/0/0		A
7	Chile		E	1/1/1/1		B	0/0/0/1		B	0/0/0/1		B	0/0/0/1		B
8	Peru		D	1/1/1/1		D	1/0/1/1		B	0/0/1/0		B	0/0/0/1		B
9	Puerto Rico		E	1/1/1/1		D	1/0/1/1		D	1/0/1/1		B	0/0/0/1		B
10	Ecuador		E	1/1/1/1		C	0/1/1/0		B	0/0/1/0		D	0/1/1/1		D
11	Dominican Republic		E	1/1/1/1		D	1/0/1/1		D	1/0/1/1		E	1/0/1/1		E
12	Cuba		E	1/1/1/1		E	1/1/1/1		E	1/1/1/1		C	0/1/0/1		C
13	Nigeria		C	1/1/0/1		A	0/0/0/0		A	0/0/0/0		A	0/0/0/0		A
14	Egypt		D	1/1/1/1		E	1/1/1/1	invalid	E	1/1/1/1	invalid	B	0/0/1/0		B
15	South Africa		E	1/1/1/1		E	1/1/1/1		E	1/1/1/1		E	1/1/1/1		E
16	Algeria		E	1/1/1/1		E	1/1/0/1		D	1/1/0/1		D	1/1/0/1		D
17	Morocco		E	1/1/1/1		E	1/1/0/1		D	1/1/0/1		D	1/1/0/1		D
18	Angola		E	1/1/1/1		C	0/1/0/1		C	0/1/0/1		C	0/1/0/1		C
19	Kenya		E	1/1/1/1		E	0/1/1/1		E	0/1/1/1		E	0/1/0/1		C
20	Ethiopia		E	0/1/1/0		B	0/0/0/1		B	0/0/0/1		B	0/0/0/1		B
21	Tanzania		E	1/1/1/1		E	1/1/1/1		E	1/1/0/1		E	1/1/0/1		E
22	Ghana		E	1/1/1/0		D	1/1/0/1		C	0/1/0/1		B	0/0/0/1		B
23	Ivory Coast		E	1/1/1/1		C	0/1/0/1		B	0/1/0/0		A	0/0/0/0		A
24	DR Congo		E	1/1/1/1		E	1/1/1/1		E	1/1/1/1		E	1/1/1/1		E
25	China		E	1/1/1/1		D	1/0/1/1		D	1/0/1/1		C	1/0/0/1		C
26	Japan		B	0/1/0/0		B	0/0/0/1		A	0/0/0/0		A	0/0/0/0		A
27	India		E	1/1/1/1		D	0/1/1/1		D	0/1/1/1		D	0/1/1/1		D
28	Iran		D	1/1/1/0		B	0/0/1/0		C	0/0/1/1		B	0/0/0/1		B
29	South Korea		D	1/1/1/0		D	0/1/1/1		D	0/1/1/1		D	0/1/1/1		D
30	Indonesia		C	0/1/1/0		B	0/0/0/1		A	0/0/0/0		A	0/0/0/0		A
31	Saudi Arabia		E	1/1/1/1		C	0/1/0/1		C	0/1/0/1		D	1/1/0/1		D
32	Turkey		C	0/0/1/1		B	0/0/0/1		B	0/0/0/1		B	0/0/0/1		D
33	Thailand		E	1/1/1/1		B	0/1/0/0		B	0/1/0/0		A	0/0/0/0		A
34	Israel		E	1/1/1/1		E	1/1/1/1		D	0/1/1/1		C	0/0/1/1		D
35	United Arab Emirates		E	1/1/1/1		C	1/0/1/0		D	1/1/1/0		D	1/1/1/0		D
36	Hong Kong		E	1/1/1/1		B	0/0/0/1		B	0/0/0/1		B	0/0/0/1		B
37	Germany		C	0/1/0/1		A	0/0/0/0		A	0/0/0/0		A	0/0/0/0		A
38	United Kingdom		E	1/1/1/1		D	1/0/1/1		D	1/0/1/1		D	1/0/1/1		D
39	France		C	1/0/1/0		B	1/0/0/0		A	0/0/0/0		A	0/0/0/0		A
40	Russia		C	1/0/1/0		C	0/1/1/0		C	0/1/1/0		B	0/0/1/0		B
41	Italy		C	1/1/0/0		C	1/1/0/0		B	1/0/0/0		A	0/0/0/0		A
42	Spain		C	0/1/1/0		E	1/1/1/1		E	1/1/1/1		E	1/1/1/1		E
43	Netherlands		B	0/1/0/0		D	1/0/1/1		C	1/1/0/0		B	0/1/0/0		B
44	Switzerland		D	1/0/1/1		C	1/0/0/1		B	0/0/0/1		A	0/0/0/0		A
45	Poland		D	1/1/1/0		B	1/0/0/0		A	0/0/0/0		A	0/0/0/0		A
46	Sweden		D	1/0/1/1		E	1/1/1/1		E	1/1/1/1		E	1/1/1/1		E
47	Finland		D	1/0/1/1		D	1/0/1/1		B	1/0/0/0		A	0/0/0/0		A
48	Iceland		E	1/1/1/1		E	1/1/1/1		E	1/1/1/1		E	1/0/1/1		D
49	Australia		E	1/1/1/1		E	1/1/1/1		E	1/1/1/1		E	1/1/1/1		E
50	New Zealand		E	1/1/1/1		E	1/1/1/1		E	1/1/1/1		E	1/1/1/1		E

Figure 7: Complete results of the flag drawing task. Multi-turn refinement allows ChatGPT to generate a more similar image to the ground truth image.

B InstructGPT for Multimodality

We show an example of a multi-turn flag drawing of InstructGPT in Figure 8. Similar to ChatGPT, InstructGPT can revise the generated flag image in each turn, although the generation quality is still elementary.

Instruct. 0 : Describe how Canadian flag looks like.

Output 0 : The flag of Canada is a red and white maple leaf design... The red and white colors...

Instruct. 1 : [I_0] [O_0] Generate a code snippet to represent that flag in SVG format.

Output 1 : 

Instruct. 2 : [I_0] [O_0] [I_1] [O_1] The flag should have a vertical red band on the left, a vertical white band in the middle, and a vertical red band on the right. It also should have a red maple leaf in the middle.

Output 2 : 

Instruct. 3 : [I_0] [O_0] [I_1] [O_1] [I_2] [O_2] The middle thing does not look like a maple leaf, which is not aligned with your description. Revise the image.

Output 3 : 

Figure 8: Example of the Canadian flag drawn by InstructGPT.

C List of Evaluation Datasets

We provide a detailed list of all the datasets used in our experiment on Table 19.

Dataset	Task	Description	Reference	#Test Size	#ChatGPT Eval
National Flag Drawing	IG	National Flag Drawing is a designed synthetic dataset which is used to evaluate the multimodal understanding of LLMs. The instruction for the National Flag Drawing is as follow: given a nation, draw the corresponding national flag and revise it based on the follow-up correction requests.	<i>Curated by authors of this paper</i>	50	50
CNN/DM	SUM	The CNN/DailyMail Dataset is an English-language dataset containing just over 300k unique news articles as written by journalists at CNN and the Daily Mail. The current version supports both extractive and abstractive summarization, though the original version was created for machine-reading and comprehension and abstractive question answering.	Nallapati et al. (2016)	11490	50
SAMSum	SUM	SAMSum dataset contains about 16k messenger-like conversations with summaries. Conversations were created and written down by linguists fluent in English. Linguists were asked to create conversations similar to those they write on a daily basis, reflecting the proportion of topics of their real-life messenger conversations.	Gliwa et al. (2019)	819	50
FLoRes-200	MT	FLoRes is a benchmark dataset for machine translation between English and four low resource languages, Nepali, Sinhala, Khmer and Pashto, based on sentences translated from Wikipedia.	Goyal et al. (2021)	1012 per language (200 languages)	30 per language (12 languages)
NusaX	SA	NusaX is a high-quality multilingual parallel corpus that covers 12 languages, Indonesian, English, and 10 Indonesian local languages, namely Acehnese, Balinese, Banjarese, Buginese, Madurese, Minangkabau, Javanese, Ngaju, Sundanese, and Toba Batak.	Winata et al. (2022)	400	50
bAbI task 15	QA	This basic deduction bAbI tasks is taken from the (20) QA bAbI tasks that a set of proxy tasks that evaluate reading comprehension via question answering. The tasks measure understanding in several ways: whether a system is able to answer questions via simple deduction. The tasks are designed to be prerequisites for any system that aims to be capable of conversing with a human.	Weston et al. (2016b)	1000	30

bAbI task 16	QA	This basic induction bAbI tasks is taken from the (20) QA bAbI tasks that a set of proxy tasks that evaluate reading comprehension via question answering. The tasks measure understanding in several ways: whether a system is able to answer questions via simple induction. The tasks are designed to be prerequisites for any system that aims to be capable of conversing with a human.	Weston et al. (2016b)	1000	30
EntailmentBank	QA	ENTAILMENTBANK, the first dataset of multistep entailment trees for QA, to support entailment-based explanation. ENTAILMENTBANK contains two parts: 1,840 entailment trees, each tree showing how a question-answer pair (QA) is entailed from a small number of relevant sentences (e.g., Figure 1); and a general corpus C, containing those and other sentences of domain-specific and general knowledge relevant to the QA domain.	Dalvi et al. (2021)	340	30
CLUTRR	QA	CLUTRR (Compositional Language Understanding and Text-based Relational Reasoning), a diagnostic benchmark suite, is first introduced in (https://arxiv.org/abs/1908.06177) to test the systematic generalization and inductive reasoning capabilities of NLU systems. The CLUTRR benchmark allows us to test a model's ability for systematic generalization by testing on stories that contain unseen combinations of logical rules, and test for the various forms of model robustness by adding different kinds of superfluous noise facts to the stories.	Sinha et al. (2019)	1146	30
α NLI	QA	Abductive Natural Language Inference (α NLI) is a new commonsense benchmark dataset designed to test an AI system's capability to apply abductive reasoning and common sense to form possible explanations for a given set of observations. Formulated as a binary-classification task, the goal is to pick the most plausible explanatory hypothesis given two observations from narrative contexts.	Bhagavatula et al. (2020)	3059	30
CommonsenseQA	QA	CommonsenseQA is a new multiple-choice question answering dataset that requires different types of commonsense knowledge to predict the correct answers . It contains 12,102 questions with one correct answer and four distractor answers. The dataset is provided in two major training/validation/testing set splits: "Random split" which is the main evaluation split, and "Question token split", see paper for details.	Talmor et al. (2018)	1221	30

HotpotQA	QA	HotpotQA is a new dataset with 113k Wikipedia-based question-answer pairs with four key features: (1) the questions require finding and reasoning over multiple supporting documents to answer; (2) the questions are diverse and not constrained to any pre-existing knowledge bases or knowledge schemas; (3) we provide sentence-level supporting facts required for reasoning, allowing QA systems to reason with strong supervision and explain the predictions; (4) we offer a new type of factoid comparison questions to test QA systems' ability to extract relevant facts and perform necessary comparison.	Yang et al. (2018)	7405	30
PiQA	QA	To apply eyeshadow without a brush, should I use a cotton swab or a toothpick? Questions requiring this kind of physical commonsense pose a challenge to state-of-the-art natural language understanding systems. The PIQA dataset introduces the task of physical commonsense reasoning and a corresponding benchmark dataset Physical Interaction: Question Answering or PIQA. Physical commonsense knowledge is a major challenge on the road to true AI-completeness, including robots that interact with the world and understand natural language. PIQA focuses on everyday situations with a preference for atypical solutions. The dataset is inspired by instructables.com, which provides users with instructions on how to build, craft, bake, or manipulate objects using everyday materials.	Bisk et al. (2020)	1838	30
E-Care	QA	Understanding causality has vital importance for various Natural Language Processing (NLP) applications. Beyond the labeled instances, conceptual explanations of the causality can provide a deep understanding of the causal fact to facilitate the causal reasoning process. We present a human-annotated explainable CAusal REasoning dataset (e-CARE), which contains over 20K causal reasoning questions, together with natural language formed explanations of the causal questions.	Du et al. (2022)	2122	30
Letter string analogy	QA	The letter string analogy domain was introduced in order to evaluate computational models of analogical reasoning. This task is composed of simple alphanumeric characters, but nevertheless require a significant degree of abstraction to identify an analogy.	Webb et al. (2022b)	-	30

SpARTQA	QA	SpARTQA is a textual question answering benchmark for spatial reasoning on natural language text which contains more realistic spatial phenomena not covered by prior datasets and that is challenging for state-of-the-art language models (LM). SPARTQA is built on NLVR’s images containing more objects with richer spatial structures. SPARTQA’s stories are more natural, have more sentences, and richer in spatial relations in each sentence, and the questions require deeper reasoning and have four types: find relation (FR), find blocks (FB), choose object (CO), and yes/no (YN), which allows for more fine-grained analysis of models’ capabilities. The default test set of this dataset is human-annotated test set, which consists of 510 instances.	Mirzaee et al. (2021)	510	64
StepGame	QA	StepGame is a robust multi-hop spatial reasoning dataset in textual format which addresses the limitation from the bAbI dataset task 17 and task 19. In this task, the AI system is required to interpret a story of k spatial relations (e.g left, top-right, 90 degree clockwise) of k+1 entities (k is up to 10) expressed in natural language and answer a question about the relative position of two of those entities. For each k from 1 to 10, there are 30000, 1000, and 30000 samples for train, valid, and test respectively.	Shi et al. (2022b)	30000	30
Pep-3k	QA	Pep-3k is a dataset of physical semantic plausibility judgments of single events. It requires a mixture of commonsense knowledge and conceptual knowledge to solve. Each event consists of a subject, a verb, and an object, i.e it has the simple s-v-o format. For example, the event can be <i>man swallow paintball</i> , with the label 0 (implausible). In total, Pep-3k has 3080 instances with plausible-implausible data balance.	Wang et al. (2018b)	3080	30
COVID-Social	MD	COVID-19-related claims fact-checked by journalists from a website called Politifact.com. Unlike the Covid19-Scientific dataset, it contains non-scientific and socially-related claims, such as “For the coronavirus, the death rate in Texas, per capita of 29 million people, we’re one of the lowest in the country.” Such claims may not be life-and-death matters, but they still have the potential to bring negative sociopolitical effects.	Lee et al. (2021)	340	50
COVID-Scientific	MD	COVID-Scientific is a testset consists of COVID-19-related scientific or medical myths that must be debunked correctly to ensure the safety of the public (e.g., "Drinking a bleach solution will prevent you from getting COVID-19").	Lee et al. (2021)	172	50

TruthfulQA	QA	The dataset comprises questions that span 38 categories, including health, law, finance and politics. The questions are of those that some humans would answer falsely due to a false belief or misconception. (e.g., [conspiracy] “If it’s cold outside what does that tell us about global warming?”; [health] “Can coughing effectively stop a heart attack?”)	Lin et al. (2022)	817	65
MultiWOZ 2.2	TOD	MultiWOZ 2.2 is an improved version of MultiWOZ 2.1, which identifies and fixes dialogue state annotation errors across 17.3% of the utterances on top of MultiWOZ 2.1 and redefines the ontology by disallowing vocabularies of slots with a large number of possible values (e.g., restaurant name, time of booking) and introducing standardized slot span annotations for these slots.	Zang et al. (2020)	5000 dialogue turns	50 dialog turns
OpenDialKG	KGD	OpenDialKG is a dataset of conversations between two crowdsourcing agents engaging in a dialog about a given topic. Each dialog turn is paired with its corresponding “KG paths” that weave together the KG entities and relations that are mentioned in the dialog.	Moon et al. (2019)	2.25k dialogue sessions	50 dialogue turns

Table 19: List of all datasets used in our experiments. IG denotes image generation, SUM denotes summarization, MT denotes machine translation, SA denotes sentiment analysis, QA denotes question answering, MD denotes misinformation detection, TOD denotes task-oriented dialogue, and KGD denotes knowledge-grounded dialogue. Some of the descriptions are directly from the original reference.

D Examples from Machine Translation and Post-Editing

Target	English Text	Label	Translation	Post-Edited Text
Chinese	Although three people were inside the house when the car impacted it, none of them were hurt.	虽然车撞到房子时，房子里面有三个人，但最后并没有人受伤。	尽管有三个人在汽车撞上房子的时候在屋里，但他们都没有受伤。	尽管汽车撞上房子时有三个人在屋里，但他们都没有受伤。
Chinese	34 per cent of those in the poll share this view, wanting Queen Elizabeth II to be Australia's last monarch.	34% 的受访者一致表示，希望英国女王伊丽莎白二世成为澳大利亚的最后一任君主。	这项民意调查显示，34%的人希望伊丽莎白二世是澳大利亚的最后一位君主。	根据民意调查显示，34%的人希望伊丽莎白二世成为澳大利亚的最后一位君主。
Chinese	The qualities that determine a subculture as distinct may be linguistic, aesthetic, religious, political, sexual, geographical, or a combination of factors.	让亚文化变得与众不同的特质，可能是语言、美学、宗教、政治、性别、地理或多种因素的组合。	什么决定一个亚文化作为不同的特质可能是语言的、审美的、宗教的、政治的、性的、地理的，或者这些因素的组合。	决定一个亚文化与众不同的特质可能是语言的、审美的、宗教的、政治的、性的、地理的，或者是这些因素的组合。
Chinese	New Kingdom ancient Egyptians marvelled at their predecessors monuments, which were then well over a thousand year old.	新王国时期的古埃及人惊叹于其前辈的已有一千多年历史的纪念碑。	新王朝古埃及人惊叹于他们的前人的纪念碑，当时已经有一千多年的历史了。	新王朝的古埃及人惊叹于他们前人留下的纪念碑，当时已经有超过一千年的历史。
Chinese	The U.S. Corps of Engineers estimated that 6 inches of rainfall could breach the previously damaged levees.	美国陆军工兵部队估计，6英寸的降雨量就可能冲垮之前损坏的堤坝。	美国工程兵部队估计，6英寸的降雨可能会破坏先前受损的堤坝。	美国工程兵部队估计，6英寸的降雨可能会使先前受损的堤坝再次损坏。
Chinese	Several large television screens were installed in various places in Rome to let the people watch the ceremony.	罗马各地安装了数个大型电视屏幕，供人们观看仪式。	在罗马的若干个地方安装了几台大型电视屏幕，让人们观看仪式。	可以，这是修改后的翻译：在罗马的各个地方安装了几台大型电视屏幕，以便让人们观看仪式。
Chinese	Pests can spoil food, cause irritation, or in a worse case cause allergic reactions, spread venom, or transmit infections.	害虫会破坏食物，引起刺激性反应，在更严重的情况下会引起过敏反应，传播毒液或传染病。	害虫可以毁坏食物，引起刺激，或者在更糟糕的情况下引起过敏反应，传播毒液，或传播感染。	害虫可以毁坏食物，引起不适，甚至会引起过敏反应、传播毒液或传染感染。
Chinese	It is obvious enough that the world has changed much because of humankind's scientific and technological advancements, and problems have become greater because of overpopulation and mankind's extravagant lifestyle.	显而易见，由于人类的科学和技术进步，世界发生了很大的变化，而人口过剩和人类奢侈的生活方式使问题变得更为严重。	很明显，由于人类的科学技术进步，世界已经发生了很大的变化，并且由于人口过剩和人类奢侈的生活方式，问题变得更加严重。	毫无疑问，人类的科学和技术进步使世界发生了巨大的变化，同时人口过剩和人类奢侈的生活方式使问题变得更加严重。

Chinese	The correlation between brain pathology and behaviour supports scientists in their research.	大脑病理和行为之间的相关性可以为科学家们的研究提供支持。	大脑病理学与行为之间的相关性支持科学家们的研究。	研究表明，大脑病理学和行为之间的相关性为科学家提供了支持。
Chinese	Like some other experts, he is skeptical about whether diabetes can be cured, noting that these findings have no relevance to people who already have Type 1 diabetes.	和其他一些专家一样，他对糖尿病能否治愈持怀疑态度。他指出，这些发现与已患有1型糖尿病的人无关。	像其他一些专家一样，他对于糖尿病是否可以被治愈持怀疑态度，并指出这些发现对已经患有1型糖尿病的人没有意义。	和其他一些专家一样，他对糖尿病是否可以被治愈表示怀疑，并指出这些发现对已经患有1型糖尿病的人没有意义。
Korean	Although three people were inside the house when the car impacted it, none of them were hurt.	차가 집에 충돌했을 때 세 명이 안에 있었지만, 그들 중 한 명도 다치지 않았다.	차가 집에 부딪히자 집 안에 세 사람이 있었지만, 아무도 다치지 않았습니다.	차가 집에 부딪히자 집 안에 세 사람이 있었지만, 아무도 상해를 입지 않았습니다.
Korean	34 per cent of those in the poll share this view, wanting Queen Elizabeth II to be Australia's last monarch.	여론 조사에서 34 퍼센트가 엘리자베스 2세가 호주의 마지막 군주이길 바란다는 의견을 보였습니다.	34%가 이 의견을 공감하며, 아스트레일리아의 최후의 왕조가 엘리자베스 2세가 되기를 원한다.	이 조사에서는 34%가 엘리자베스 2세가 아스트레일리아의 최후의 왕조가 되기를 원한다는 의견을 공감한다.
Korean	The qualities that determine a subculture as distinct may be linguistic, aesthetic, religious, political, sexual, geographical, or a combination of factors.	하위문화를 명확하게 구분하는 특징은 언어적, 미적, 종교적, 정치적, 성적, 지리적 요소가 있으며, 이러한 요소들의 결합일 수도 있다.	"다음 문장의 한국어 번역은 무엇입니까? 부문문화를 구별되게 하는 특성은 언어, 예술, 종교, 정치, 성, 지리 요소일 수 있거나 이들 요소의 조합일 수도 있습니다."	부문문화가 구별되는 특성은 언어적, 예술적, 종교적, 정치적, 성적, 지리적 요소 중 하나일 수도 있고, 이들 요소의 조합일 수도 있습니다.
Korean	New Kingdom ancient Egyptians marvelled at their predecessors monuments, which were then well over a thousand year old.	고대 신왕국 이집트인들은 조상의 기념비적인 건축물을 보고 경탄했고 이것은 그 당시 기준으로 1000년은 족히 넘은 건축물이었습니다.	신나라 이집트인들은 전자들이 지금보다 약 1,000년 이상 오래된 고대 유적을 창구로 칭했습니다.	신나라 이집트인들은 전자들이 지금보다 약 1,000년 이상 오래된 고대 유적을 창구로 칭했고, 그들은 그것들을 참화했습니다.
Korean	The U.S. Corps of Engineers estimated that 6 inches of rainfall could breach the previously damaged levees.	미국 공병대는 시간당 6인치의 강우량이 기 파손된 제방을 무너뜨릴 수 있다고 추정했다.	미국 엔지니어 중대는 6인치의 비가 이전에 손상된 제방을 깨고 들어올 수 있다고 추정했습니다.	미국 엔지니어 중대는 6인치의 비가 이전에 손상된 제방을 깨고 가로막을 칠 수 있다고 추정했습니다.
Korean	Several large television screens were installed in various places in Rome to let the people watch the ceremony.	대형 텔레비전 스크린 여러 대가 로마 곳곳에 설치되어 사람들이 장례식을 관람할 수 있도록 했습니다.	로마에서 여러 곳에 거대한 텔레비전 스크린이 설치되어 이들이 의식을 시청할 수 있게 했습니다.	로마에서 여러 곳에 거대한 텔레비전 스크린이 설치되어 이들이 의식을 시청할 수 있게 해주었습니다.

Korean	Pests can spoil food, cause irritation, or in a worse case cause allergic reactions, spread venom, or transmit infections.	해충은 음식이 썩게 만들고 염증을 유발하거나, 더 나쁜 경우 알레르기 반응을 일으키고 독을 퍼뜨리거나 전염병을 옮길 수 있습니다.	당신은 잡초가 식료품을 망치거나, 자극을 유발할 수 있거나, 최악의 경우 알레르기 반응을 유발하거나, 독을 전파하거나, 감염을 전파할 수 있음을 알고 있나요?	잡초는 식료품을 망칠 수 있고, 자극을 유발할 수 있거나, 최악의 경우 알레르기 반응을 유발할 수 있고, 독을 전파할 수 있거나, 감염을 전파할 수 있습니다.
Korean	It is obvious enough that the world has changed much because of humankind's scientific and technological advancements, and problems have become greater because of overpopulation and mankind's extravagant lifestyle.	인류의 과학 기술 발전으로 세상이 많이 변했다는 것은 충분히 명백하고, 또 인구과잉과 인류의 사치스러운 생활방식 때문에 문제가 더 커졌다.	인간의 과학기술 진보와 인간의 과감한 생활방식으로 인해 세계가 상당히 바뀌었고, 인구 초과와 문제가 커졌다는 것은 분명합니다.	인간의 과학기술 진보로 인해 세계가 상당히 바뀌었고, 인구 초과와 인간의 과감한 생활방식으로 인해 문제가 커졌음은 확실합니다.
Korean	The correlation between brain pathology and behaviour supports scientists in their research.	뇌 병리와 행동 사이의 상관관계가 과학자들의 연구를 돕습니다.	뇌의 질환과 행동간의 상관관계가 과학자들의 연구를 지원한다.	뇌의 질환과 행동 사이의 상관관계는 과학자들의 연구를 지원합니다.
Korean	Like some other experts, he is skeptical about whether diabetes can be cured, noting that these findings have no relevance to people who already have Type 1 diabetes.	다른 전문가들과 마찬가지로, 그는 당뇨병의 치료 여부에 회의적이며, 이러한 결과는 제1형 당뇨병 환자에게는 관련이 없음을 지적합니다.	몇몇 전문가들과 마찬가지로, 그는 저주파 치료가 당뇨병을 완전히 치료할 수 있을지에 대해 의문을 갖고 있으며, 이러한 연구 결과는 이미 1형 당뇨병을 가지고 있는 사람들과는 전혀 관계가 없다고 주장한다.	그는 일부 전문가들과 마찬가지로, 이들 연구 결과는 이미 1형 당뇨병을 가진 사람들에게는 전혀 관계가 없다고 주장하며, 당뇨병이 완전히 치료될 수 있을지에 대해 의문을 갖고 있다.
Japanese	Although three people were inside the house when the car impacted it, none of them were hurt.	車の衝突時には3人が家の中にいましたが、怪我人はいませんでした。	この文を日本語に訳すと、「車が家に衝突したときには家の中に3人いましたが、そのうちの誰も傷を負わなかった」となります。	もしよろしければ、この文を元の文の意味を維持しながら、より日本語らしい表現に書き換えることもできます。例えば次のようになります。「車が家に衝突したときには、家の中には3人いましたが、どの人も怪我をしなかった」
Japanese	34 per cent of those in the poll share this view, wanting Queen Elizabeth II to be Australia's last monarch.	世論調査では34%の人がこの意見に賛成しており、エリザベス女王2世がオーストラリア最後の君主であることを望んでいます。	34%の人たちは、この見解を共有して、イングランド女王エリザベス2世が最後の王位を持つことを望む。	34%のアンケートに参加した人たちのうち、この見解を共有しており、イングランド女王エリザベス2世がオーストラリアで最後の王位を持つことを望んでいます。

Japanese	The qualities that determine a subculture as distinct may be linguistic, aesthetic, religious, political, sexual, geographical, or a combination of factors.	サブカルチャーの特徴を決定付ける性質は、言語、美、宗教、政治、性、地理、またはそれらの要素の組み合わせかもしれません。	サブカルチャーを特徴的なものとする品質は、言語的、美学的、宗教的、政治的、性的、地理的、あるいはそれらの要因の組み合わせであることがあります。	そのような品質は、サブカルチャーを特徴的なものとするものである可能性があります。言語的、美学的、宗教的、政治的、性的、地理的、あるいはそれらの要因の組み合わせであることがあります。
Japanese	New Kingdom ancient Egyptians marvelled at their predecessors monuments, which were then well over a thousand year old.	新王国の古代エジプト人たちは、当時千年以上前に建てられた前身のモニュメントに驚嘆しました。	新王国時代のエジプト人は、当時すでに1000年以上も古い前人のモニュメントに驚嘆しました。	新王国時代のエジプト人は、それらの先人の墓苑が当時すでに1000年以上も古いことに驚嘆しました。
Japanese	The U.S. Corps of Engineers estimated that 6 inches of rainfall could breach the previously damaged levees.	米工兵隊は、以前に決壊した堤防が6インチの降水量で再び決壊する可能性があるかと推定しています。	米国エンジニア兵団は、6インチの降雨で以前に損傷した堤防が浸水すると推定しています。	以前に損傷を受けた堤防が、米国エンジニア兵団が推定する6インチの降雨で浸水するかもしれません。
Japanese	Several large television screens were installed in various places in Rome to let the people watch the ceremony.	ローマの各所には、人々に式典を見てもらおうと、大型テレビ画面がいくつも設置されていました。	ローマにはさまざまな場所に大型テレビ画面が設置され、人々が式典を観られるようになりました。	ローマには、さまざまな場所に大型のテレビスクリーンが設置されました。これにより、人々は式典を観ることができるようになりました。
Japanese	Pests can spoil food, cause irritation, or in a worse case cause allergic reactions, spread venom, or transmit infections.	害虫は食物を台無しにしたり、皮膚を刺したり、最悪の場合はアレルギー反応の原因になったり、毒を撒き散らしたり、感染症を媒介することがあります。	害虫は食物を台無しにすることがあり、かゆみを起こすこともあり、さらに悪い場合にはアレルギー反応を引き起こすこともあり、毒を撒いたり、感染を伝播することもあります。	害虫は食物を損なうことがあり、かゆみを起こすこともあります。また、悪い場合にはアレルギー反応を引き起こすこともあり、毒を撒いたり、感染を伝播することもあります。
Japanese	It is obvious enough that the world has changed much because of humankind's scientific and technological advancements, and problems have become greater because of overpopulation and mankind's extravagant lifestyle.	人類の科学技術の進歩によって世界は大きく変わり、過剰な人口と贅沢な生活が原因で問題が拡大したことは十分に明白な事実です。	人類の科学技術の進歩によって世界は大きく変化し、人口の増加や人間の浪費的な生活スタイルによって問題は大きくなっていることは明らかだ。	明らかであるように、人類の科学技術の進歩により、世界は大きく変化し、人口増加や人間の浪費的な生活スタイルにより、問題はより大きくなっている。
Japanese	The correlation between brain pathology and behaviour supports scientists in their research.	脳病理と行動の相関関係は、科学者たちの研究を裏付けるものです。	脳の病理学と行動の関連は、科学者たちの研究を支援する。	脳の病理学と行動の関連は、科学者たちの研究を支援する

Japanese	Like some other experts, he is skeptical about whether diabetes can be cured, noting that these findings have no relevance to people who already have Type 1 diabetes.	他の専門家たちと同様に、彼は糖尿病を完治できるかどうかについては懐疑的であり、これらの調査結果はすでにI型糖尿病を患っている人々には何の関連性もないことを指摘しています。	日本語に翻訳すると、他の専門家のように、彼は糖尿病が治せるかどうかについて懐疑的であり、これらの発見は既に1型の糖尿病を持つ人々には関係がないと述べている。	日本語に翻訳すると、他の専門家のように、彼は糖尿病が完治できるかどうかについて懐疑的であり、これらの発見は既に1型の糖尿病を持つ人々には無関係だと述べている。
French	Although three people were inside the house when the car impacted it, none of them were hurt.	Bien que trois personnes aient été présentes dans la maison quand la voiture l'a percutée, aucune n'a été blessée.	Bien que trois personnes se trouvaient à l'intérieur de la maison lorsque la voiture a percuté celle-ci, aucune d'entre elles n'a été blessée.	Même si trois personnes étaient à l'intérieur de la maison lorsque la voiture a heurté celle-ci, aucune d'entre elles n'a été blessée.
French	34 per cent of those in the poll share this view, wanting Queen Elizabeth II to be Australia's last monarch.	34 % des personnes interrogées partagent ce point de vue, et souhaitent que la reine Elizabeth II soit la dernière monarque à régner sur l'Australie.	34 pour cent de ceux qui ont participé à l'enquête partagent cet avis, souhaitant que la reine Elizabeth II soit la dernière monarchie d'Australie	34 pour cent de ceux qui ont participé à l'enquête souhaitent que la reine Elizabeth II soit la dernière monarchie d'Australie.
French	The qualities that determine a subculture as distinct may be linguistic, aesthetic, religious, political, sexual, geographical, or a combination of factors.	Les critères qui déterminent une sous-culture comme distincte peuvent être linguistiques, esthétiques, religieux, politiques, sexuels, géographiques ou une combinaison de facteurs.	Les qualités qui déterminent une sous-culture comme distincte peuvent être linguistiques, esthétiques, religieuses, politiques, sexuelles, géographiques ou une combinaison de facteurs."	Les caractéristiques qui déterminent une sous-culture comme distincte peuvent être linguistiques, esthétiques, religieuses, politiques, sexuelles, géographiques ou une combinaison de ces éléments.
French	New Kingdom ancient Egyptians marvelled at their predecessors monuments, which were then well over a thousand year old.	Les Égyptiens antiques du Nouvel Empire étaient émerveillés par les monuments de leurs prédécesseurs, qui avaient alors bien plus de mille ans.	Les Égyptiens anciens de la Nouvelle-République admiraient les monuments de leurs prédécesseurs, qui étaient alors bien plus âgés d'un millénaire.	Les Égyptiens anciens de la Nouvelle-République admiraient les monuments de leurs prédécesseurs, qui avaient alors plus de mille ans.
French	The U.S. Corps of Engineers estimated that 6 inches of rainfall could breach the previously damaged levees.	D'après le U.S. Corps of Engineers, une pluviométrie de 6 pouces pourrait ouvrir une brèche dans les digues précédemment endommagées.	Le Corps des ingénieurs des États-Unis a estimé qu'une pluie de 15 cm pourrait affaiblir les digues précédemment endommagées.	Le Corps des ingénieurs des États-Unis a estimé qu'une pluie de 15 cm pourrait mettre à mal les digues précédemment endommagées.
French	Several large television screens were installed in various places in Rome to let the people watch the ceremony.	Plusieurs grands écrans de télévision ont été installés à différents endroits à Rome pour permettre aux gens de regarder la cérémonie.	Plusieurs grands écrans de télévision ont été installés dans différents endroits de Rome pour permettre aux gens de regarder la cérémonie.	Plusieurs grands écrans de télévision ont été installés dans différents endroits de Rome afin de permettre aux gens de regarder la cérémonie.

French	Pests can spoil food, cause irritation, or in a worse case cause allergic reactions, spread venom, or transmit infections.	Les parasites peuvent contaminer la nourriture, provoquer des irritations ou, dans les cas les plus graves, provoquer des réactions allergiques, répandre du venin ou transmettre des infections.	Les nuisibles peuvent gâcher les aliments, provoquer des irritations, ou dans le pire des cas provoquer des réactions allergiques, propager du venin ou transmettre des infections.	Les nuisibles peuvent altérer les aliments, causer des irritations, ou dans les cas les plus graves provoquer des réactions allergiques, répandre du venin ou transmettre des infections.
French	It is obvious enough that the world has changed much because of humankind's scientific and technological advancements, and problems have become greater because of overpopulation and mankind's extravagant lifestyle.	Il est clair que le monde a beaucoup changé grâce aux progrès scientifiques et technologiques réalisés par les êtres humains, et que les problèmes sont devenus plus importants en raison de la surpopulation et du mode de vie extravagant de l'humanité.	Il est assez évident que le monde a beaucoup changé en raison des avancées scientifiques et technologiques de l'humanité et que les problèmes se sont accrus en raison de la surpopulation et du mode de vie extravagant de l'homme.	Il est évident que le monde a beaucoup changé en raison des avancées scientifiques et technologiques de l'humanité, et que les problèmes ont augmenté à cause de la surpopulation et du mode de vie extravagant de l'homme.
French	The correlation between brain pathology and behaviour supports scientists in their research.	La corrélation entre la pathologie du cerveau et le comportement soutient les scientifiques dans leurs recherches.	La corrélation entre la pathologie cérébrale et le comportement soutient les scientifiques dans leur recherche.	La corrélation entre la pathologie cérébrale et le comportement soutient les scientifiques dans leurs recherches.
French	Like some other experts, he is skeptical about whether diabetes can be cured, noting that these findings have no relevance to people who already have Type 1 diabetes.	À l'instar d'autres experts, il se montre sceptique quant à la possibilité de guérir le diabète, faisant remarquer que ces résultats ne sont pas applicables aux personnes qui souffrent déjà de diabète de type 1.	Comme d'autres experts, il est sceptique quant à la possibilité de guérir du diabète, notant que ces résultats n'ont aucune pertinence pour les personnes atteintes de diabète de type 1.	Tout comme certains autres experts, il est sceptique quant à la possibilité de guérir du diabète, et souligne que ces résultats n'ont aucune pertinence pour les personnes atteintes de diabète de type 1.
Indonesian	Although three people were inside the house when the car impacted it, none of them were hurt.	Walaupun ada tiga orang di dalam rumah yang ditabrak mobil, tak satu pun yang cedera.	Meskipun tiga orang ada di dalam rumah saat mobil menabraknya, tidak ada yang terluka.	Meskipun tiga orang ada di dalam rumah saat mobil menabraknya, mereka semua selamat.
Indonesian	34 per cent of those in the poll share this view, wanting Queen Elizabeth II to be Australia's last monarch.	Dalam jajak pendapat mengenai pandangan ini, 34 persen dari mereka menginginkan Ratu Elizabeth II menjadi pemimpin monarki terakhir Australia.	34 persen dari orang yang terlibat dalam polling ini memiliki pandangan yang sama, yaitu ingin Ratu Elizabeth II menjadi ratu terakhir di Australia.	34 persen dari orang yang terlibat dalam polling ini memiliki pandangan yang sama, yaitu ingin Ratu Elizabeth II menjadi ratu terakhir di Australia.
Indonesian	The qualities that determine a subculture as distinct may be linguistic, aesthetic, religious, political, sexual, geographical, or a combination of factors.	Kualitas yang menentukan suatu subkultur berbeda dari yang lain bisa berupa bahasa, estetika, agama, politik, seksual, geografis, atau kombinasi sejumlah faktor.	Kualitas yang menentukan suatu subkultur sebagai distinct dapat berupa linguistik, estetika, agama, politik, seksual, geografis, atau kombinasi faktor-faktor tersebut.	Faktor-faktor yang menentukan sebuah subkultur sebagai distinct bisa berupa faktor linguistik, estetika, agama, politik, seksual, geografis, atau kombinasi dari beberapa faktor tersebut.

Indonesian	New Kingdom ancient Egyptians marvelled at their predecessors monuments, which were then well over a thousand year old.	Kerajaan Baru Mesir kuno mengagumi monumen para pendahulunya, yang saat itu berusia lebih dari seribu tahun.	Kekaisaran Baru orang-orang Mesir kuno terpesona dengan monumen-monumen pendahulunya, yang saat itu sudah lebih dari seribu tahun tua.	Orang-orang Mesir kuno Kekaisaran Baru terpesona dengan monumen-monumen pendahulunya, yang saat itu sudah lebih dari seribu tahun tua.
Indonesian	The U.S. Corps of Engineers estimated that 6 inches of rainfall could breach the previously damaged levees.	U.S. Corps of Engineers memperkirakan bahwa curah hujan sebesar 6 inci ini mampu menerobos tanggul yang sudah rusak.	Korps Engineer AS menghitung bahwa hujan 6 inci dapat merusak parit yang sebelumnya sudah rusak.	Badan Engineer AS memperkirakan bahwa hujan 6 inci dapat menembus parit yang sebelumnya sudah rusak.
Indonesian	Several large television screens were installed in various places in Rome to let the people watch the ceremony.	Beberapa layar televisi besar dipasang di berbagai tempat di Roma supaya orang-orang bisa menyaksikan upacara tersebut.	Beberapa layar televisi besar diinstal di berbagai tempat di Roma agar orang-orang dapat menonton upacara.	Berbagai layar televisi besar diinstal di tempat-tempat yang berbeda di Roma agar masyarakat dapat menonton upacara tersebut.
Indonesian	Pests can spoil food, cause irritation, or in a worse case cause allergic reactions, spread venom, or transmit infections.	Hama dapat merusak makanan, mengakibatkan iritasi, atau dalam kasus yang lebih parah, mengakibatkan reaksi alergi, menyebarkan racun atau menularkan penyakit.	Hama dapat merusak makanan, menyebabkan iritasi, atau dalam keadaan yang lebih buruk menyebabkan reaksi alergi, menyebarkan racun, atau menularkan infeksi.	Hama dapat merusak makanan, menyebabkan iritasi, atau dalam keadaan yang lebih buruk bahkan dapat menyebabkan reaksi alergi, menyebarkan racun, atau menularkan infeksi.
Indonesian	It is obvious enough that the world has changed much because of humankind's scientific and technological advancements, and problems have become greater because of overpopulation and mankind's extravagant lifestyle.	Cukup jelas bahwa dunia berkembang karena ilmu pengetahuan dan teknologi, serta kelebihan jumlah penduduk dan gaya hidup boros manusia menjadi masalah yang cukup besar.	Jelas sekali bahwa dunia telah berubah banyak karena kemajuan ilmiah dan teknologi manusia, dan masalah telah menjadi lebih besar karena kelebihan populasi dan gaya hidup manusia yang berlebihan.	Sudah jelas sekali bahwa dunia telah berubah sangat banyak karena kemajuan ilmiah dan teknologi manusia, dan masalah telah menjadi lebih besar karena kelebihan populasi dan gaya hidup manusia yang berlebihan.
Indonesian	The correlation between brain pathology and behaviour supports scientists in their research.	Korelasi antara patologi otak dan perilaku mendukung para ilmuwan dalam penelitian mereka.	Korelasi antara patologi otak dan perilaku membantu para ilmuwan dalam penelitian mereka.	Korelasi antara patologi otak dan tingkah laku membantu para ilmuwan dalam penelitian mereka.
Indonesian	Like some other experts, he is skeptical about whether diabetes can be cured, noting that these findings have no relevance to people who already have Type 1 diabetes.	Seperti halnya ahli-ahli lain, dia merasa skeptis apakah diabetes dapat disembuhkan, menegaskan bahwa penemuan-penemuan ini tidak relevan dengan orang-orang yang telah mengidap diabetes Jenis 1.	Sama seperti beberapa ahli lainnya, ia ragu-ragu tentang apakah diabetes dapat disembuhkan, mencatat bahwa temuan-temuan ini tidak relevan bagi orang-orang yang sudah menderita diabetes tipe 1.	Seperti beberapa ahli lainnya, ia meragukan apakah diabetes dapat disembuhkan, mencatat bahwa temuan-temuan ini tidak relevan bagi orang-orang yang sudah menderita diabetes tipe 1.

Javanese	Although three people were inside the house when the car impacted it, none of them were hurt.	Sanadyan uwong telu ana ing njero omah nalika mobil kuwi nabrak, nanging kabeh selamat.	Walaupun tiga orang berada di dalam rumah saat mobil menabraknya, tidak ada di antara mereka yang terluka.	Walaupun ana ing tiga wong ing rumah kala mobil mbentak ing kono, ora ana saka ingkang terluka.
Javanese	34 per cent of those in the poll share this view, wanting Queen Elizabeth II to be Australia's last monarch.	34 per sen sing menehi swara sesawangane ngene, kepengin Ratu Elizabeth II dadi raja sing keridhewane ana ning Australia.	Sebagian 34 persen dari mereka yang terlibat dalam polling ini berpendapat demikian, ingin Ratu Elizabeth II menjadi Ratu Australia terakhir.	Sebagian 34 persen dari mereka yang terlibat dalam polling ini memiliki pandangan yang sama, yaitu menginginkan Ratu Elizabeth II menjadi Ratu Australia terakhir.
Javanese	The qualities that determine a subculture as distinct may be linguistic, aesthetic, religious, political, sexual, geographical, or a combination of factors.	Kualitas sing nemtokake sawijining cabang kabudayan mbokmenawa basa, estetika, sing ana gayutane karo agama, politik, seksual, geografis, utawa sawijining kombinasi saka faktor-faktor.	Sifat-sifat yang menentukan subkultur sebagai yang terpisah bisa berupa linguistik, estetika, agama, politik, seksual, geografis, atau kombinasi faktor-faktor.	Kualitas-kualitas sing menentukan subkultur minangka distinct bisa ngandhut linguistik, estetika, agama, politik, seksual, geografis, utawa gabungan faktor-faktor.
Javanese	New Kingdom ancient Egyptians marvelled at their predecessors monuments, which were then well over a thousand year old.	Krajan Anyar Mesir kuno nggumun marang monumen-monumen leluhure, sing nalika iku tetep apik luwih seka sewu taun.	Rakyat Kerajaan Baru dari Mesir Kuno terpesona oleh monumen-monumen pendahulunya, yang saat itu sudah lebih dari seribu tahun tua.	Rakyat Kerajaan Baru dari Mesir Kuno terpesona oleh monumen-monumen pendahulunya, yang saat itu sudah lebih dari seribu tahun tua.
Javanese	The U.S. Corps of Engineers estimated that 6 inches of rainfall could breach the previously damaged levees.	Korps Insinyur AS ngira manawa udan 6 inci bisa nembus tanggul sing rusak sadurunge.	Korps Insinyur Amerika Serikat menghitung bahwa hujan 6 inci bisa merusak parit yang sebelumnya rusak.	Korps Insinyur Amerika Serikat menaksir bahwa hujan 6 inci bisa merusak parit yang telah rusak sebelumnya.
Javanese	Several large television screens were installed in various places in Rome to let the people watch the ceremony.	Sawetara layar TV gedhe dipasang ing sawetara panggonan ing Roma supaya masyarakat bisa nonton upacara kasebut.	Beberapa layar televisi besar dipasang di berbagai tempat di Roma agar orang-orang dapat menonton upacara.	Beberapa layar televisi besar dipasang di beberapa tempat di Roma supaya orang-orang dapat menyaksikan upacara itu.
Javanese	Pests can spoil food, cause irritation, or in a worse case cause allergic reactions, spread venom, or transmit infections.	Ama bisa ngrusak panganan, nyebabake iritasi, utawa ing kasus sing paling parah nyebabake reaksi alergi, nyebar racun, utawa ngirim infeksi.	Hama bisa merusak makanan, menyebabkan iritasi, atau dalam kasus yang lebih buruk menyebabkan reaksi alergi, menyebarkan racun, atau menularkan infeksi.	Hama bisa merusak makanan, menyebabkan iritasi, atau dalam keadaan yang lebih buruk menyebabkan reaksi alergi, menyebarkan racun, atau menularkan infeksi.

Javanese	It is obvious enough that the world has changed much because of humankind's scientific and technological advancements, and problems have become greater because of overpopulation and mankind's extravagant lifestyle.	Cukup jelas manawa donya wis malih akeh amarga majune ilmu lan teknologi manungsa, lan masalah-masalah dadi tambah gedhe amarga kaluwihan populasi lan gaya urip sing boros saka manungsa.	Sawetara bisa dilihat manawa donya wis duwé sak kabèh pangowahan sain-teknologi manungsa lan masalah duwé gedhé luwih amarga tambah-tambahé jumlah penduduk lan gaya hidup manungsa sing luwih mewah.	Sawetara bisa dilihat manawa donya wis duwé sak kabèh pangowahan sain lan teknologi saka manungsa, lan masalah-masalah duwé gedhé luwih amarga tambah-tambahé jumlah penduduk lan gaya hidup manungsa sing luwih mewah
Javanese	The correlation between brain pathology and behaviour supports scientists in their research.	Sesambungan antarane ilmu lelara utek lan tumindak ndhukung para ilmuwan kanggo nлити.	Korelasi antara patologi otak dan perilaku mendukung para ilmuwan dalam penelitian mereka.	Korelasi antara patologi otak dan tingkah laku mendukung para ilmuwan dalam penelitian mereka.
Javanese	Like some other experts, he is skeptical about whether diabetes can be cured, noting that these findings have no relevance to people who already have Type 1 diabetes.	Kaya sawetara ahli liyane, dheweke ora yakin apa diabetes bisa ditambani, nimbang yen panemuan kasebut ora duwe hubungane karo uwong sing uwis ngalami diabetes Jinis 1.	Sebagaimana beberapa ahli lainnya, dia skeptis tentang apakah diabetes dapat disembuhkan, dengan catatan bahwa temuan ini tidak relevan bagi orang-orang yang sudah menderita diabetes tipe 1.	Sekaya karo para ahli liya, dia luwih leres saka wether diabetes bisa dibenerake, nggandeng rembug kuwi ora duwe arti kanggo wong-wong sing wis duwe diabetes tipe 1.
Sundanese	Although three people were inside the house when the car impacted it, none of them were hurt.	Sanaos tilu urang aya di jero bumi pas mobil na nabrak, teu aya anu cidera.	Walaupun tiga orang berada di dalam rumah saat mobil menabraknya, tidak ada yang terluka.	Walaupun ada tiga orang di dalam rumah saat mobil menabraknya, tak ada yang terluka.
Sundanese	34 per cent of those in the poll share this view, wanting Queen Elizabeth II to be Australia's last monarch.	34 persen jalma dina polling di gaduh pendapat ieu, nyaeta hoyong Ratu Elizabeth II janten monarkis pamungkas di Australia.	34 persen dari mereka yang terlibat dalam survei memiliki pandangan ini, yaitu ingin Ratu Elizabeth II menjadi ratu terakhir Australia.	34 persen dari mereka yang terlibat dalam survei memiliki pandangan yang sama, yaitu ingin Ratu Elizabeth II menjadi ratu terakhir di Australia.
Sundanese	The qualities that determine a subculture as distinct may be linguistic, aesthetic, religious, political, sexual, geographical, or a combination of factors.	Kualitas anu nangtukeun subkultur salaku bédana panginten basa, éstétis, religius, politik, séks, géografis, atanapi kombinasi faktor.	Sifat-sifat yang menentukan subkultur sebagai terpisah mungkin bersifat linguistik, estetika, agama, politik, seksual, geografis, atau kombinasi dari faktor-faktor tersebut.	Sifat-sifat nu nentukkeun subkultur salaku nu béda bisa jadi linguistik, estetika, agama, politik, seksual, geografis, atawa kombinasi ti faktor-faktor ieu.
Sundanese	New Kingdom ancient Egyptians marvelled at their predecessors monuments, which were then well over a thousand year old.	Karajaan anyar urang Mesir kuno takjub kana monumén paninggalan karuhunna, anu umurna langkung ti sarebu taun.	Para raja anyar dina Egypt kuno nyogog nara-nara jaman awal, anu dina jaman ayeuna kira-kira ribuan taun.	Para raja anyar dina Egypt kuno terkagum-kagum kana nara-nara candi jaman awal, anu dina jaman ayeuna kira-kira leuweung taun.
Sundanese	The U.S. Corps of Engineers estimated that 6 inches of rainfall could breach the previously damaged levees.	Korps insinyur A.S ngira-ngira curah hujan 6 inci tiasa nembus tanggul anu samemehna geus ruksak.	Korps Engineer AS diperkirakan bahwa 6 inci hujan bisa merusak benteng yang sebelumnya sudah rusak.	Korps Engineer AS diperkirakan bahwa hujan sebesar 6 inci dapat merusak benteng yang sudah rusak sebelumnya.

Sundanese	Several large television screens were installed in various places in Rome to let the people watch the ceremony.	Sababaraha tipi ageung dipasang di sababaraha tempat di Roma supados warga tiasa nyaksian upacara éta.	Beberapa layar televisi besar dipasang di berbagai tempat di Roma agar orang-orang bisa menonton upacara.	Sabagéan layar televisi gaduh dipasang di hiji bilik roma sinaunna pikeun ngajadikeun jalma bisa nonton upacara.
Sundanese	Pests can spoil food, cause irritation, or in a worse case cause allergic reactions, spread venom, or transmit infections.	Hama tiasa ngarusak dahareun, nyieun iritasi, atawa leuwih parah deui tiasa nyieun alérgi, nyebarkeun racun, atawa nularkeun panyakit.	Hama bisa merusak makanan, menyebabkan iritasi, atau dalam kasus yang lebih buruk menyebabkan reaksi alergi, menyebarkan racun, atau menularkan infeksi.	Hama bisa merusak makanan, menyebabkan iritasi, atau dalam kasus yang lebih buruk menyebabkan reaksi alergi, menyebarkan racun, atau menularkan infeksi.
Sundanese	It is obvious enough that the world has changed much because of humankind's scientific and technological advancements, and problems have become greater because of overpopulation and mankind's extravagant lifestyle.	Jelas pisan yén dunya parantos seueur robih kusabab kamajuan élmu manusa sareng téknologi, sareng masalahna janten langkung ageung kusabab seueur penduduk sareng gaya hirup anu boros umat manusa.	Nyaéta nujul bisa ngeuyeuk yén dunya geus robah heubeul dina jaman sains jeung teknologi manusa, jeung masalah geus ngaronjat luyu ku luweung laju penduduk jeung gaya hidup manusa anu boros.	Nyaéta nujul bisa ngeuyeuk yén dunya geus robah heubeul dina jaman sains jeung teknologi manusa, jeung masalah geus ngaronjat luyu ku luweung laju penduduk jeung gaya hidup manusa anu boros.
Sundanese	The correlation between brain pathology and behaviour supports scientists in their research.	Hubungan patologi polo sareng kalakuan ngabantos para élmuwan dina panalungtikanna.	Korelasi antara patologi otak dan perilaku mendukung para ilmuwan dalam penelitian mereka.	Korelasina antara patologi otak jeung ulah-ulahan ngalapkeun dukungan sakurang-kurangna pikeun para ilmuwan dina penelitian maranéhanana.
Sundanese	Like some other experts, he is skeptical about whether diabetes can be cured, noting that these findings have no relevance to people who already have Type 1 diabetes.	Sapertos sababaraha ahli anu sanés, anjeunna henteu percanteun upami diabétés tiasa disembuhkeun, kusabab pamanggihna ieu téh henteu aya hubunganana jeung jalma anu parantos gaduh diabétés tipe 1.	Seperti beberapa ahli lainnya, dia skeptis tentang apakah diabetes bisa disembuhkan, mencatat bahwa temuan ini tidak relevan bagi orang-orang yang sudah memiliki diabetes tipe 1.	Kayaku ngan ahli séjén, dia bérék-bérék ngeunaan jangdi diabetes bisa diobat, ngeunaan yén kajadian ieu teu aya hubunganna jeung jalma anu geus ngalami diabetes tipe 1."

Table 20: Examples of ChatGPT translated and post-edited sentences.

E Evaluation Results for Reasoning

We provide the complete results for reasoning tasks on Table 21.

Categories	Testset	Result
Deductive	ENTAILMENTBANK	28/30
	bAbI (task 15)	28/30 (as is - 19/30)
Inductive	CLUTRR	13/30
	bAbI (task16)	20/30 (as is - 0/30)
Abductive	α NLI	26/30
Mathematical	Math	13/30
Temporal	Timedial	26/30
Spatial	SpartQA (hard)	8/32
	SpartQA (basic)	20/32
	StepGame (hard)	7/30
	StepGame (basic)	19/30
	StepGame (basic-cardinal)	17/20
	StepGame (diagonal)	11/20
	StepGame (clock-direction)	5/20
Commonsense	CommonsenseQA	27/30
	PIQA	25/30
	Pep-3k (Hard)	28/30
Causal	E-Care	24/30
Multi-hop	hotpotQA	8/30
Analogical	Letter string analogy	30/30

Table 21: Composed results for all reasoning tasks.

F Multi-turn for Task-Oriented Dialogue

We provide the example for the modular and unified approaches for Task-Oriented Dialogue in Table 22 and Table 23, respectively.

Task	Key	Text Content
Dialogue State Tracking	Prompt	<p>Give the dialogue state of the last utterance in the following dialogue in the form of 'STATE: Domain-Intent: [Slot, Possible value], ... (for example: STATE: Hotel-Inform: ['area', 'centre']) by using the following pre-defined slots and possible values:</p> <p>Intents: Request, Inform, general-thank, general-bye Domain: hotel, Slots: pricerange, Possible values: ['expensive', 'cheap', 'moderate'] Domain: hotel, Slots: type, Possible values: ['guesthouse', 'hotel'] Domain: hotel, Slots: parking, Possible values: ['free', 'no', 'yes'] Domain: hotel, Slots: bookday, Possible values: ['monday', 'tuesday', 'wednesday', 'thursday', 'friday', 'saturday', 'sunday'] Domain: hotel, Slots: bookpeople, Possible values: ['1', '2', '3', '4', '5', '6', '7', '8'] Domain: hotel, Slots: bookstay, Possible values: ['1', '2', '3', '4', '5', '6', '7', '8'] Domain: hotel, Slots: stars, Possible values: ['0', '1', '2', '3', '4', '5'] Domain: hotel, Slots: internet, Possible values: ['free', 'no', 'yes'] Domain: hotel, Slots: name, Possible values: [] Domain: hotel, Slots: area, Possible values: ['centre', 'east', 'north', 'south', 'west'] Domain: hotel, Slots: address, Possible values: [?] Domain: hotel, Slots: phone, Possible values: [?] Domain: hotel, Slots: postcode, Possible values: [?] Domain: hotel, Slots: ref, Possible values: [?]</p> <p>USER: I need a five stars hotel close to a mall and main restaurants. The hotel should include free wifi in the room. STATE : 'Hotel-Inform': [['internet', 'yes'], ['stars', '5']] SYSTEM: I'm sorry, but there are no 5-star hotels anywhere in Cambridge. Would you like to broaden your criteria? USER: How about something with 4 stars and something expensive. STATE : 'Hotel-Inform': [['pricerange', 'expensive'], ['stars', '4']] SYSTEM: I have 2 options for you - The Huntingdon Marriott on the west side, and the University Arms in the center. Are you interested in one of these? USER: Maybe. Is either one a 4 star hotel? If so, I'd like to book a room for 4 nights. STATE : 'Hotel-Inform': [['bookstay', '4'], ['stars', '4']] SYSTEM: Both are 4 stars. How many days would you be staying, and how many people are in your party? I can see if something is available. USER: I need a reservation for 4 nights starting Saturday.</p>
	Gold	STATE : 'Hotel-Inform': [['bookday', 'saturday'], ['stars', '4']]
	ChatGPT Generation	STATE : 'Hotel-Inform': [['bookday', 'saturday'], ['bookstay', '4'], ['stars', '4']]
	Response Generation using Actions	Prompt
Gold	SYSTEM: As I mentioned it is cheap and the phone number is 01223316074.	
ChatGPT Generation	SYSTEM: The phone number for the 4-star hotel located in North Cambridge is 01223316074 and it is in the cheap price range.	

Table 22: Examples of modular Task-Oriented Dialogue using ChatGPT: dialogue state tracking and response generation

Task	Key	Text Content
Multi-turn unified approach	Prompt	<p>Use the following knowledge base to complete the task of "recommending a restaurant" by continuing the conversation as a task-oriented dialogue system:</p> <p>Restaurant: Mama Julia, Food: French, Price: Expensive, Location: 7th street, Rating: 5 Restaurant: Papa John, Food: Italian, Price: Moderate, Location: 3rd street, Rating: 4 Restaurant: The Crossroad, Food: Morocco, Price: Moderate, Location: Downtown, Rating: 2 Restaurant: Tacos City, Food: Mexian, Price: Cheap, Location: Center, Rating: 1 Restaurant: Golden Rice Bowl, Food: Chinese, Price: Cheap, Location: 3rd district, Rating: 3 Restaurant: Veggie Garden, Food: Chinese, Price: Expensive, Location: Town Hall, Rating: 4 Restaurant: Pizza House, Food: Italian, Price: Moderate, Location: 3rd street, Rating: 2 Restaurant: The Palace, Food: Vietnamese, Price: Expensive, Location: Hotel Grandview, Rating: 5</p> <p>Hi, I'm looking for a restaurant</p>
	ChatGPT	Sure, what type of food are you in the mood for?
	Prompt	I only want restaurants rated 3 or higher. Any suggestions?
	ChatGPT	Here are some restaurants that have a rating of 3 or higher: Mama Julia (French, Expensive, 5 rating), Papa John (Italian, Moderate, 4 rating), Golden Rice Bowl (Chinese, Cheap, 3 rating), Veggie Garden (Chinese, Expensive, 4 rating), and The Palace (Vietnamese, Expensive, 5 rating). Which one would you like to know more about?
	Prompt	What food does The Palace serve?
	ChatGPT	The Palace serves Vietnamese cuisine.
	Prompt	Is Taxos City expensive to eat at?
	ChatGPT	Tacos City is considered as cheap restaurant.
	Prompt	I want to eat Italian food, any recommendations?
	ChatGPT	Based on the information provided in the knowledge base, I would recommend Papa John and Pizza House for Italian food. Papa John is rated 4 out of 5 and is classified as "Moderate" in terms of price, while Pizza House is rated 2 out of 5 and is also classified as "Moderate". Papa John is located on 3rd street, while Pizza House is located on 3rd street as well.
	Prompt	Did you forget I only wanted restaurants rated 3 or higher?
	ChatGPT	I apologize, I must have misunderstood your previous request. Based on the knowledge base provided, The only Italian restaurant that is rated 3 or higher is Papa John which is rated 4 out of 5 and is located on 3rd street. It is classified as "Moderate" in terms of price.
	Prompt	Great thanks for your help.
	ChatGPT	You're welcome! I hope you enjoy your meal at Papa John. If you have any other questions or need further assistance, feel free to ask.

Table 23: Example of multi-turn unified approach for Task-Oriented Dialogue using ChatGPT