

# Can Large Language Models Understand Context?

Anonymous ACL submission

## Abstract

Understanding context is key to understanding human language, an ability which Large Language Models (LLMs) have been increasingly seen to demonstrate to an impressive extent. However, though the evaluation of LLMs encompasses various domains within the realm of Natural Language Processing, limited attention has been paid to probing their linguistic capability of understanding contextual features. This paper introduces a context understanding benchmark by adapting existing datasets to suit the evaluation of generative models. This benchmark comprises of four distinct tasks and nine datasets, all featuring prompts designed to assess the models' ability to understand context. First, we evaluate the performance of LLMs under the in-context learning pretraining scenario. Experimental results indicate that pre-trained dense models struggle with understanding more nuanced contextual features when compared to state-of-the-art fine-tuned models. Second, as LLM compression holds growing significance in both research and real-world applications, we assess the context understanding of quantized models under in-context-learning settings. We find that 3-bit quantization leads to varying degrees of performance reduction on our benchmark. We conduct an extensive analysis of these scenarios to substantiate our experimental results.

## 1 Introduction

Discourse understanding, as one of the fundamental problems in NLP, focuses on modeling linguistic features and structures that go beyond individual sentences (Joty et al., 2019). Understanding discourse requires resolving the relations between words/phrases (coreference resolution) and discourse units (discourse parsing and discourse relation classification) in the previous context, identifying carry-over information for the following context (dialogue state tracking), and recognizing discourse-specific phenomena (ellipsis).

LLMs have garnered substantial attention from both academia and the industry due to their remarkable capability in comprehending language and world knowledge. Their unparalleled performance across a diverse range of benchmarks and datasets has firmly established their significance in a relatively short period of time. As LLMs continue to push the boundaries of scale and capability, the evaluation of their multifaceted abilities becomes an equally vital endeavor. Consequently, the development of robust evaluation methodologies to assess specific aspects of LLMs becomes imperative. In addition, these methodologies should focus on helping achieve a comprehensive understanding of their advancement while clearly delineating their limitations. However, recently published LLMs, such as OPT (Zhang et al., 2022), LLaMA (Touvron et al., 2023) and GPT-4 (OpenAI, 2023), are only evaluated on limited benchmarks, and have a significant drawback: they neglect the inclusion of discourse-related datasets for evaluation, thereby limiting the comprehensive assessment of their language understanding capabilities.

To provide a comprehensive evaluation, plenty of benchmarks and datasets address various facets of language understanding, including benchmarks that delve into common sense knowledge (Hendrycks et al., 2021a; Kwiatkowski et al., 2019), as well as linguistic capabilities like sentiment analysis, natural language inference, summarization, text classification, and more (Bang et al., 2023b; Liang et al., 2022). These general benchmarks and specific dataset evaluations exhibit certain limitations. Despite the requirement for contextual information in these benchmarks to effectively tackle tasks (for example, sentiment analysis requires an understanding of polarities within the given text), none of these benchmarks cater to tasks that demand a nuanced comprehension of linguistic features within a provided context.

On the other hand, recent LLMs, by virtue of

043  
044  
045  
046  
047  
048  
049  
050  
051  
052  
053  
054  
055  
056  
057  
058  
059  
060  
061  
062  
063  
064  
065  
066  
067  
068  
069  
070  
071  
072  
073  
074  
075  
076  
077  
078  
079  
080  
081  
082  
083

possessing billions of parameters, have led to an exponential surge in computational and storage costs (Brown et al., 2020b), which hinders the deployment of large models to personal devices and restricts the on-device performance of language understanding tasks. To address this challenge, model compression methods, which can reduce memory and disk requirements of both model training and inference, have gained attention. Existing compression techniques, such as 3-bit quantization (Frantar et al., 2022), have demonstrated the potential to reduce model sizes with only marginal performance trade-offs. However, the evaluation of quantization methods suffers from two deficiencies. Firstly, quantization methods are primarily evaluated on limited benchmarks and datasets, such as Lambada (Paperno et al., 2016), ARC (Boratto et al., 2018), PIQA (Tata and Patel, 2003), BoolQ (Clark et al., 2019), and StoryCloze (Mostafazadeh et al., 2017). Secondly, previous work has not delved into a linguistic analysis to identify where the model efficacy wanes.

Given the above shortcomings, this paper evaluates LLMs on a context understanding benchmark constructed from varied discourse understanding datasets. We conduct an extensive analysis of LLM performance on this benchmark, including models of varying sizes and those subjected to compression techniques, aiming to provide a more comprehensive understanding of context understanding capability of the LLMs. The contributions of this paper can be summarized as follows:

- Our work introduces a contextual understanding benchmark, including four tasks, for the evaluation of LLMs. We also present prompts designed for in-context learning on each task.
- We evaluate LLMs of varying sizes from different model families and provide an analysis on these models' capability for context understanding.
- We evaluate post-training compressed models in ICL settings and conduct an analysis of the reduction in context understanding capability compared to dense models.

## 2 Related Work

### 2.1 In-context Learning Evaluation

The paradigm of ICL (Brown et al., 2020a) is rapidly gaining importance. Studies have demon-

strated that the generalization of LLMs to various downstream NLP tasks, such as MMLU (Hendrycks et al., 2021b), is significantly enhanced when provided with a small number of examples as prompts (Brown et al., 2020a; Chowdhery et al., 2022; Hoffmann et al., 2022; Rae et al., 2022; Anil et al., 2023; Touvron et al., 2023; OpenAI, 2022, 2023). Recent research has extensively evaluated the performance of LLMs across a spectrum of language-related tasks, spanning from text generation to understanding input sequences. This assessment contains a wide array of benchmarks, including SUPER-GLUE (Wang et al., 2019; Laskar et al., 2023), and tasks such as question answering, information retrieval, sentiment analysis (Bang et al., 2023b; Liang et al., 2022), dialogue (Heck et al., 2023), and text classification (Yang and Menczer, 2023).

### 2.2 Model Compression for LLMs

Model compression techniques can be broadly categorized into three main approaches: compression during training, compression associated with fine-tuning, and post-training methods. In terms of quantization during training, this technique enables LLMs to adapt to low-precision representations during the training process (Liu et al., 2023). Model compression with fine-tuning involves quantization awareness into the fine-tuning stage (Kim et al., 2023; Dettmers et al., 2023). Post-training techniques, on the other hand, are applied after the completion of an LLMs training phase and typically involve the use of calibration data. This category comprises two primary approaches: pruning, which removes redundant or non-salient weights to induce sparsity (Frantar and Alistarh, 2023), and quantization, which employs low-precision numeric representations of weights and activations (Nagel et al., 2020; Frantar et al., 2022; Yuan et al., 2023). Prior research shows that quantization outperforms pruning in several settings (Kuzmin et al., 2023), thus in this work, we focus on model quantization and its impact on the selected context-aware tasks.

## 3 Task Selection & Design

Our contextual understanding benchmark includes four tasks with nine datasets, as presented in Table 1. In the following sections, we provide detailed explanations of each task and the corresponding datasets, along with the designed prompts for ICL evaluations.

Type	Task	Dataset	Context
Doc	Coreference	WSC273	Nominal & eventual reference
		OntoNotes	
	Discourse	PDTB-3	Relations between discourse units
Dial.	DST	MultiWoz	Entity carryover within context
	Query Rewrite	MuDoCo	Ellipsis and reference
		QReCC	
		InCar	
GECOR			
CANARD			

Table 1: Tasks and datasets in the context understanding benchmark.

### 3.1 Coreference Resolution

The coreference resolution task (CR) contributes to achieving a coherent understanding of the overall meaning conveyed within the text. Thus, it plays a critical role in diving into language models’ capability to grasp coreference relations as well as contextual nuances within documents. We select two coreference datasets: WSC273 (Levesque et al., 2012) and OntoNotes 5.0 (Pradhan et al., 2013).

WSC273, which contains the first 273 examples from the Winograd Schema Challenge, is a dataset that requires the system to read a sentence with an ambiguous pronoun and select the referent of that pronoun from two choices. OntoNotes is a human-annotated corpus of documents annotated with multiple layers of linguistic information including syntax, propositions, named entities, word sense, and in-document coreference. As it is one of the most frequently used datasets for training coreference models, prior research has achieved significant advancements under the supervised fine-tuning paradigm (Lee et al., 2017; Joshi et al., 2020; Bohnet et al., 2023). However, these model designs cannot be extended to generative models under ICL settings. Recently, Le and Ritter (2023) have leveraged document templates for LLMs; however, their evaluation is confined to prominent models such as InstructGPT (Ouyang et al., 2022), neglecting the fact that smaller models lack the generative capacity required to accomplish such tasks. Due to these limitations, we propose a novel multiple-choice task design. In this design, we provide the mentions and evaluate the model on resolution. Each option represents a potentially markable span.<sup>1</sup> Ta-

<sup>1</sup>Considering the inferior performance of small models on the mention detection task, we utilize gold markable spans coreference linking.

<b>Instruction:</b> Please carefully read the following passages. For each passage and the options, you must identify which option the mention marked in <b>*bold*</b> refers to. If the marked mention does not have any antecedent, please select “no antecedent”.
<b>Context:</b> ... To express <b>*its*</b> determination ... the Chinese securities regulatory department ... this stock reform ...
<b>Choices:</b>
A. no antecedent
B. the Chinese securities regulatory department
C. this stock reform
...
<b>Question:</b> What does <b>*its*</b> refer to?
<b>Answer:</b> B

Table 2: An OntoNotes example of prompt and answer.

ble 2 presents an example of the input to the model<sup>2</sup>. The entire prompt consists of five parts: (1) an instruction that provides guidance to the model for the task, (2) a document containing plain text with a selected mention span highlighted using a bold symbol, (3) a list of choices, which includes all the gold mentions present in the document, (4) a question that directs the model’s attention, and (5) a guiding word *answer* that prompts for the output. We experiment with multiple instructions and prompts and provide the one with the best performance. Linking scores are computed for each question and the results are subsequently aggregated for evaluation. We utilize the official evaluation metrics from the CoNLL-2012 shared task (Pradhan et al., 2012), which employs the CoNLL F1 score, derived from the averaging of three coreference metrics: MUC, B<sup>3</sup>, and CEAF<sub>φ4</sub>.

### 3.2 Dialogue State Tracking

Dialogue state tracking (DST) is an important task in the area of task-oriented dialogue (TOD) modeling (Young et al., 2013), where the dialogue agent tracks the key information provided by the user as the conversation progresses. Table 3 provides an example from MultiWOZ (Budzianowski et al., 2018) where the user expresses the constraints when looking for a restaurant. The output of DST is typically maintained in slot-value pair format.

Previous research has explored ICL capabilities on MultiWOZ and demonstrated promising results compared to fine-tuning models (Hu et al., 2022; Heck et al., 2023). However, these studies either involve partial training or are untested with smaller and quantized models. Here we adopt a straightforward and simplified ICL approach proposed by

<sup>2</sup>Detailed examples for each task design can be found in Appendix A.

---

**Ontology:**  
 {"slots": {"restaurant-pricerange": "price budget for the restaurant", ... },  
 "categorical": {"restaurant-pricerange": ['cheap', 'expensive', 'moderate'], ... } }

**Instruction:** Now consider the following dialogue between two parties called the "system" and "user". Can you tell me which of the "slot" was updated by the "user" in its latest response to the "system"? Present the updates in JSON format. If no "slots" were updated, return an empty JSON list. If you encounter "slot" that was requested by the "user" then fill them with "?". If a user does not seem to care about a discussed "slot" fill it with "dontcare".

**[Previous Dialogue State]**  
**[Conversation]:**  
 "system": ""  
 "user": "I'm looking for a moderately priced place to eat that's in the centre of town."  
**Output:** {"restaurant-pricerange": "moderate", "restaurant-area": "centre"}

---

Table 3: A DST example of prompt and *answer*.

Heck et al. (2023), and test it on MultiWOZ v2.2 (Zang et al., 2020). The prompt to the model consists of domain knowledge from ontology, an instruction, previous dialogue state (the belief state accumulated until the previous user turn) and the conversation proceeding to the current turn. The ontology could be very long if we consider all domains in the dataset; thus, given the input length constraint of LLMs, only the knowledge relevant to the conversation is provided. Following literature, we report joint goal accuracy (JGA) (Mrkšić et al., 2017) for evaluating the performance of DST.

### 3.3 Implicit Discourse Relation Classification

Discourse demonstrates its importance beyond individual sentences, which emphasizes the ways in which different segments of a text interconnect and structure themselves to convey a coherent and meaningful message. The PDTB-3 corpus, as introduced by Webber et al. (2019), annotates implicit discourse relations across elementary discourse units (EDUs)<sup>3</sup>. These relations imply connections between EDUs and may be made explicit by inserting a connective. Within the context of the understanding benchmark, we opt for the implicit discourse relation classification task for two primary reasons. Firstly, the order of the two EDUs is provided, enabling the model to directly utilize this information. Secondly, the connective triggering the relation is implicit, increasing the task’s complexity. In this task (Disc.), two EDUs are fed as input, and the objective of the task is to

<sup>3</sup>EDU refers to the smallest segment of a text that conveys a complete and coherent meaning within larger discourse.

---

**Instruction:** Given two arguments and a list of connective words, please select the most likely connective between two arguments.  
**[Relation Description]**  
**Input:**  
 Arg 1: Amcore, also a bank holding company, has assets of \$1.06 billion.  
 Arg 2: Central’s assets are \$240 million.  
**Question:** What is the connective that best describes the relation between two arguments?  
**Choices:**  
 A. Temporal B. Contingency C. Comparison D. Expansion  
**Answer:** C

---

Table 4: A PDTB example of prompt and *answer*.

correctly identify the relation between them. Due to the nuanced differences between each relation and the demand for annotators with rich linguistic knowledge and extensive annotation training, the classification task poses challenges to fine-tuned classification models.

The PDTB3 corpus classifies discourse relations into four categories - Temporal, Contingency, Comparison, and Expansion. We convert this task into a multiple-choice question and experiment with *classes* as options. In the *classes* scenario, the task offers four options, with each representing a distinct discourse relation class. Table 4 exhibits the components of the prompt. It includes an instruction at the beginning, followed by a concise description of each relation, a context with two arguments, a question along with answer choices, and a trigger word. We evaluate each model’s performance on this dataset using accuracy as the metric.

### 3.4 Query Rewrite

While document-based CR (OntoNotes, Section 3.1) covers various types of coreference relations across multiple genres, it does not allow the ability to evaluate certain aspects which are important to understand context. Firstly, the CR task typically focuses on document-based coreference chains, neglecting mention resolution in dialogues. Secondly, ellipsis, which is the omission of one or more words from a clause while still allowing it to be understood in context, is a crucial linguistic phenomenon frequently encountered in speech and conversation. It is essential for language models to grasp and accurately identify ellipses within context. Incorporating these features into the benchmark is thus pivotal when evaluating context understanding.

Query rewrite (QR) is the task of rewriting the last utterance of a user into a context-free, independent utterance that can be interpreted without dia-

<b>Instruction:</b> Rewrite the last query following interaction into a well-formed, context independent query. Resolve any disfluencies or grammatical errors in the query.
<b>Input:</b>
User: Try to reach Forbes now .
Bot: Forbes at Washington Post ? Or Forbes of Publishing Division ?
User: Publishing Division .
<b>Rewrite:</b> <i>Forbes of Publishing Division</i>

Table 5: A query rewrite example of prompt and *answer*.

log context. The objective of the task is to identify the entity or events references from the previous query, whether through a pronoun or an omitted word/phrase, and then generate a new query that includes the previous context directly.

We incorporate five QR datasets in the proposed benchmark: MuDoCo (Martin et al., 2020), QReCC (Anantha et al., 2021), InCar (Regan et al., 2019), GECOR (Quan et al., 2019), and CANARD (Elgohary et al., 2019). These datasets span multiple genres and domains in dialogues. We experiment with various prompts used for fine-tuning models and present the results with the best selections. Table 5 presents a concise prompt comprising an instruction along with context for each dialogue. To assess the quality of generated queries, we follow the metrics from previous research, particularly BLEU (Papineni et al., 2002) and ROUGE (Lin, 2004).

## 4 Experiments

The evaluation was conducted on a computational infrastructure comprising  $8 \times$  A100 GPUs. We experiment with three model families. For smaller models, we consider OPT (Zhang et al., 2022), ranging from 125M to 2.7B. Although OPT also offers larger models, we opt for LLaMA (Touvron et al., 2023) as the mid-sized LMs, spanning from 7B to 65B parameters, due to showcased superior performance by prior works. For large-scale LMs, we leverage GPT-3.5-turbo<sup>4</sup>. For each model, on every dataset, we assess five different settings: zero-shot, one-shot, 5-shot, 8-shot, and 10-shot. We randomly select the examples from the training set for the few-shot prompting.<sup>5</sup>

<sup>4</sup><https://platform.openai.com/docs/models/gpt-3-5>

<sup>5</sup>WSC273 itself is a test set and it does not have any fine-tuning scores, so we only report the zero-shot results in the table.

### 4.1 Dense Model

Results of the three model families are reported in Table 6, along with results of fine-tuned (FT) models to help better interpret how well the pre-trained models behave with ICL. For each, we present the N-shot setting that yields the highest score (see Appendix B for details). Overall, performance improves as the model size increases and pre-trained models with ICL struggle to catch up with FT models on most tasks.

**Coreference Resolution** Larger models exhibit promising performance on the WSC273 task, indicating that LLMs can effectively handle "simple" coreference relations within limited contexts and mentions. However, when it comes to document-based CR with complex clusters, their performance substantially drops<sup>6</sup>. Even on providing ground-truth mentions, the highest-performing GPT is only on par with rule-based coreference systems (Manning et al., 2014) and is far from the end-to-end fine-tuned SpanBERT (Joshi et al., 2020). The gap between ICL and FT results highlights that under the ICL setting, LLMs struggle to build coreference chains without adequate domain-specific examples. Specifically, models except GPT perform significantly worse on the MUC metric. Error analysis reveals that these models are inclined to create more clusters, including singleton clusters. This implies that pre-trained LLMs encounter difficulties in understanding long-range contextual information.

**DST** A similar trend is observed as CR where OPT and LLaMA models fall behind GPT-3.5 significantly. This suggests that these models fail to extract key information as the conversation proceeds, even with the provision of 5 to 10 demonstrations and the distilled relevant domain ontology in prompt. Our error analysis indicates that most of the errors happen due to the misdetection of slots or the wrong predicted value in a slot-value pair. Only GPT-3.5 reaches the level of FT results which is a fine-tuned T5 base model (Bang et al., 2023a).

**Implicit discourse relation classification** We observe an increase in scores when the model size exceeds 7B. However, even the best-performing

<sup>6</sup>Note that the OntoNotes dataset is substantially larger than the others. We observe that inference on the entire test set becomes extremely time-consuming, particularly with the larger models; further, the cost of running inference on GPT-3.5 starts becoming non-negligible. Consequently, we propose limiting the OntoNotes test set to a 10% sub-sample, which is the setting we consistently adopt.

Task	Dataset	Metrics	OPT				LLaMA			GPT	FT
			125M	350M	1.3B	2.7B	7B	13B	30B	3.5-turbo	
CR	WSC273	Acc	58.24	66.67	76.19	77.66	86.81	89.38	89.01	88.64	N/A
	OntoNotes	MUC	12.66	7.58	13.21	8.29	10.31	31.8	33.56	56.32	77.26
		B <sup>3</sup>	53.8	52.26	53.54	52.41	52.20	58.43	58.66	68.20	73.43
		CEAF <sub><math>\phi_4</math></sub>	31.09	29.49	31.40	30.10	32.63	38.0	39.27	50.72	74.46
	Avg. F1	32.52	29.78	32.72	30.27	31.71	42.74	43.83	58.41	76.03	
DST	MultiWOZ	JGA	11.11	27.96	26.61	28.08	32.30	28.12	42.24	57.40	63.79
Disc.	PDTB-3	Acc	10.04	10.04	10.04	16.15	17.16	26.01	39.77	43.83	76.23
QR	MuDoCo	BLEU	0.46	0.36	7.02	49.2	41.12	61.15	66.51	57.14	80.31
		ROUGE1	1.52	12.18	10.98	65.61	56.07	74.78	77.88	79.37	92.01
	QReCC	BLEU	4.53	31.27	26.35	40.09	28.19	38.64	58.68	55.24	58.67
		ROUGE1	13.91	58.18	53.10	68.32	48.27	56.40	78.74	79.98	81.75
	InCar	BLEU	0.00	7.66	12.71	27.42	28.20	42.13	48.58	63.66	88.45
		ROUGE1	3.41	28.76	30.45	49.63	49.96	56.73	64.18	83.51	95.24
	GECOR	BLEU	0.20	26.40	26.32	49.99	53.27	66.30	73.80	63.34	82.56
		ROUGE1	4.06	42.13	42.57	65.89	69.23	80.99	86.03	79.00	92.63
	CANARD	BLEU	2.61	19.39	24.24	34.66	21.34	29.32	47.24	47.12	57.46
		ROUGE1	9.82	45.63	49.36	62.73	38.17	46.61	69.73	74.61	81.06

Table 6: Few-shot results of two open-sourced models and GPT-3.5 on the context understanding benchmark. The results with the best number of few-shot examples are reported for each task. Fine-tuning (FT) results serves as a reference when evaluating LLMs’ capability under ICL setup.

model, GPT, achieves  $> 30$  points lower than the current SOTA fine-tuned model (Liu and Strube, 2023). We carefully examine the predictions for each model and found that all models tend to predict the same relation class for every example, albeit with their individual preferences for the selected relation. This suggests that the models struggle to distinguish the nuances between different relation classes and fail to correctly identify relations across EDUs within context.

**Query Rewriting** The gap between small and large models is significantly huge, compared to the other tasks. For instance, OPT-125M cannot even complete the rewriting task. Analysis on predictions of small models indicates that the model is not capable of following the instructions or learning patterns from the few-shot examples. We identify a few major error types: (1) generating the next sentence, instead of rewriting; (2) rewriting the wrong user turn from the conversation; (3) copying the last user utterance without any rewriting. These errors get reduced as the model size increases. However, similar to the previous three tasks, the best ICL results achieved by GPT is far from the fine-tuned models.<sup>7</sup> It is worth noting that OPT-2.7B performs on par or notably better than LLaMA-7B, which is somewhat not aligned with the findings in Beeching

<sup>7</sup>In literature, the best FT results come from different models across five QR datasets, where some are not even LLM based. To ensure fair comparison, we fine-tuned a T5 large model on each QR dataset.

et al. (2023) where LLaMA-7B even outperforms OPT-66B in many tasks, including ARC (Clark et al., 2018), HellaSwag (Zellers et al., 2019), and MMLU (Hendrycks et al., 2021b).

All in all, this section presents a holistic comparison of LLMs’ behaviors on the target context understanding tasks. On the tasks with structured outputs such as CR or DST, even small models show a certain level of context understanding and seem to follow the task instruction. Classification tasks such as discourse relation selection are deemed the easiest among all tasks; however, the small models are even worse than a random guess (25%). As for the generative task, the ability to complete query rewriting can be only observed in the case of larger models, as the model has the freedom to generate arbitrary content that does not follow the prompt. Except for DST, FT models demonstrate marked superiority over pre-trained models, highlighting the potential for improving LLMs’ competence on these context understanding tasks.

## 4.2 Model Compression Technique

As we focus on evaluating context understanding of LLMs in an ICL setup, we evaluate models quantized using GPTQ (Frantar et al., 2022), which is an efficient one-shot weight quantization algorithm based on approximate second-order information that compresses the model post-training. It enables a reduction in memory and disk requirements by up to 80%, compared to the pre-quantized model.

Dataset	Metrics	7B-D	30B-Q	30B-D
WSC273	Acc	86.81	87.18	89.01
OntoNotes	MUC	10.31	25.37	33.56
	B <sup>3</sup>	52.20	56.80	58.66
	CEAF <sub>φ4</sub>	32.63	36.93	39.27
	Avg. F1	31.71	39.70	43.83
MultiWOZ	JGA	32.30	41.99	42.24
PDTB-3	Acc	17.16	31.29	39.77
MuDoCo	BLEU	41.12	59.22	66.51
	ROUGE1	56.07	71.38	77.88
QReCC	BLEU	28.19	53.72	58.68
	ROUGE1	48.27	74.13	78.74
InCar	BLEU	28.20	39.69	48.58
	ROUGE1	49.96	56.32	64.18
GECOR	BLEU	53.27	70.41	83.36
	ROUGE1	69.23	73.80	86.03
CANARD	BLEU	21.34	45.07	47.24
	ROUGE1	38.17	67.15	69.73

Table 7: Comparison between dense and quantized models. Dense LLaMA-7B and 3-bit quantized LLaMA-30B share similar memory and disk requirements. **D** represents dense model and **Q** denotes quantized model.

### 4.3 Quantized Model Results

GPTQ (Frantar et al., 2022) has been shown to effectively reduce the model size to 3 bits without incurring substantial performance losses across a range of NLP tasks, such as MMLU, ARC, StoryCloze. However, whether this performance preservation can be extended to contextual understanding was unclear.

Table 7 presents the comparison between the dense and 3-bit quantized LLaMA models. In contrast to previous studies on 3-bit quantization, we observed that quantization leads to fluctuated drops in performance across the four tasks. Specifically, in WSC273, MultiWoz, and CANARD, post-training quantization incurs only a marginal performance drop ( $\sim 1.7$  points). However, in the remaining datasets, quantization results in significant performance drops.

The results further show that the quantized LLaMA-30B model consistently outperforms the dense LLaMA-7B model across all tasks despite being comparable in disk and memory requirements. For CR, the 30B quantized model achieves significantly higher scores on the OntoNotes dataset across all metrics. The MUC metric shows the most substantial improvement, indicating that the quantized 30B model partially overcomes the tendency to create small clusters for mentions. For DST on MultiWOZ, the 30B quantized model show a 30% relative improvement over the 7B model in JGA. On discourse parsing with PDTB-3, the ac-

Dataset	6.7/7B		13B		30B	
	O.	L.	O.	L.	O.	L.
Mudoco	53.1	41.1	55.2	61.1	55.2	66.5
	71.8	56.0	72.1	74.7	71.5	77.8
QReCC	46.6	28.1	43.7	38.6	43.8	58.6
	73.4	48.2	71.6	56.4	71.9	78.7
InCar	40.3	28.2	41.9	42.1	44.6	48.5
	64.8	49.9	62.6	56.7	65.3	64.1
GECOR	58.8	53.2	60.9	66.3	58.2	73.8
	75.7	69.2	78.3	80.9	76.1	86.0
CANARD	43.8	21.3	37.5	29.3	41.3	47.2
	72.0	38.1	66.0	46.6	69.3	69.7

Table 8: Comparison between OPT (O.) and LLaMA (L.) across five query rewrite datasets

curacy of quantized 30B model is almost double, 17.16% vs 31.29%. Across all QR datasets, the quantized 30B model substantially improves NLG scores compared to the dense 7B model, with relative gains ranging from 15-50%. The largest gap is observed on GECOR.

In general, we show that the quantized 30B LLaMA model consistently and significantly outperforms the dense 7B model as a result of the increased scale, despite using 3-bit quantization. The benefits of greater model scale thus outweigh the impacts of quantization in understanding discourse. We believe this finding would be beneficial when deploying LLMs in real-world applications with disk and runtime constraints.

## 5 Case Study: Query Rewrite

In this section, we provide in-depth analysis by comparing the two open-sourced model families OPT and LLaMA, and the impact of quantization, using query rewrite as the target task.

### 5.1 OPT vs. LLaMA

Prior works (Beeching et al., 2023) have consistently shown that, under the same model size, LLaMA outperforms OPT. However, their performance on QR, as shown in Table 8, does not follow this pattern. When the model size is around 7B, OPT consistently performs better than LLaMA by a significant margin across the five QR datasets. The two models perform on par with each other at 13B. The superiority of LLaMA is only obvious with 30B model size. From another perspective, although we expect performance to improve as model size increases, we observe this trend on LLaMA, but not on OPT. These results suggest that it may not be correct to conclude the overall superiority between two model families by only comparing on

Context	
User:	what is the name of india pakistan border line
Bot:	The Radcliffe Line was the boundary demarcation line between the Indian and Pakistani portions of the Punjab and Bengal provinces of British India.
User:	who created the radcliffe line
Bot:	The Radcliffe Line was named after its architect, Sir Cyril Radcliffe, who was the joint chairman of the two boundary commissions for the two provinces.
User:	when was the line published
<b>Gold answer:</b> when was the <u>radcliffe</u> line published	
<b>Prediction 1 (repeat the last query):</b> when was the line published	
<b>Prediction 2 (language modeling):</b> 1947	

Table 9: An example of two major types of errors found in the query rewrite task.

a certain range of model sizes or on a certain set of tasks.

## 5.2 Dense vs. Quantized

We conduct a quantitative analysis on the error types of query rewriting to investigate the performance gap between dense and quantized models with comparable computing requirements. Across the five datasets, we identify two main error types that account for nearly 80% of the total errors. First, the model *repeats* the last query without rewriting the referred entity. In this case, the model seems to understand the instruction but fails to comprehend the referred entity within the context. This type of error can be primarily associated with the model’s context understanding capability. Second, the model treats the task as a language modeling (*LM*) task, as shown in Table 9, where it provides a response to the last query. In this scenario, the model appears to struggle to understand the task instruction, even with several few-shot examples. We classify this error type as more related to the model’s ICL ability.

We perform manual error annotations on the five QR datasets<sup>8</sup>. Table 10 illustrates the number of errors for the three selected model settings in LLaMA for each dataset. A consistent trend is observed across all QR datasets. In terms of *repeat* errors, the 30B dense model exhibits fewer errors, around half, compared to the 7B dense model (297 vs. 469). However, 3-bit GPTQ quantization leads to an increase in this type of error, reaching a similar error count to the 7B dense model (458 vs. 469). This suggests that 3-bit quantization reduces the model’s ability to comprehend the context. Regarding *LM* errors, the 30B dense model

<sup>8</sup>10% test data on QReCC and CANARD was graded.

Type	Dataset	7B D	30B Q	30B D
Repeat	MuDoCo	260	247	194
	QReCC	86	90	26
	InCar	17	15	8
	GECOR	59	62	37
	CANARD	47	44	32
	Total	469	458	297
LM	MuDoCo	71	29	16
	QReCC	80	28	16
	InCar	19	20	15
	GECOR	6	1	0
	CANARD	127	76	59
	Total	232	125	106

Table 10: Number of the major two types errors on three LLaMA models (7B dense, 30B quantized, and 30B dense) found in Query rewrite. *Repeat* stands for repeat-the-last-query error and *LM* denotes language modeling error.

also significantly outperforms the 7B dense model, with 106 errors compared to 232. It is to be noted that the quantized model generates 125 *LM* errors, slightly more than the 30B dense model. However, it generates significantly fewer (around 50%) errors compared to the 7B dense model. This indicates that 3-bit quantization maintains the ICL capability when evaluated on our benchmark.

## 6 Conclusion

This paper introduces a contextual understanding benchmark designed to assess the performance of LLMs. We collect nine existing datasets spanning four tasks, each carefully tailored to suit generative models. This benchmark encompasses essential elements for assessing linguistic comprehension within context, including both document and dialog based contextual understanding. Experimental results under an in-context learning setting reveal that LLMs struggle with nuanced linguistic features within this challenging benchmark, exhibiting inconsistencies with other benchmarks that emphasize other aspects of language. To the best of knowledge, we are also the first to compare dense models and post-training quantization models in contextual understanding tasks. This comparison highlights that 3-bit post-training quantization reduces the general understanding capacity of context, particularly in complex references and task-oriented dialogue states. Our proposed contextual comprehension benchmark thus provides a unique perspective on the contextual dimension of language understanding and offers a valuable addition to existing LLM evaluations.



## 589 Limitations

590 This work provides an evaluation of various pre-  
591 trained LLMs, including OPT, LLaMA, and GPT,  
592 on our understanding benchmark. However, we  
593 have not evaluated other LLMs designed for longer  
594 input scenarios, such as LongLLaMA (Tworkowski  
595 et al., 2023).

596 In our evaluation, we focus on the GPTQ quan-  
597 tization method, analyzing its performance on our  
598 benchmark. We do not include other post-training  
599 quantization techniques, such as RPTQ (Yuan et al.,  
600 2023), in this work.

601 Our evaluation concentrates on English datasets,  
602 primarily utilizing LLMs pre-trained with English  
603 data. All of the four tasks on our benchmark have  
604 datasets from other languages. The coreference  
605 dataset OntoNotes 5.0 contains annotations of Ara-  
606 bic and Chinese. In addition, recent releases such  
607 as CorefUD (Nedoluzhko et al., 2022) promote  
608 standardization of multilingual coreference anno-  
609 tations. In DST, CrossWOZ (Zhu et al., 2020) is a  
610 cross-domain wizard-of-oz task-oriented dataset.  
611 Long et al. (2020) develop TED-CDB, a Chi-  
612 nese discourse relation dataset. The query rewrite  
613 task also has datasets in other languages, such as  
614 REWRITE (Su et al., 2019) and Restoration-200K  
615 (Pan et al., 2019). Finally, specific LLMs opti-  
616 mized for individual languages, such as ChatGLM  
617 (Du et al., 2022), exist and are not a part of our  
618 evaluation.

## 619 References

620 Raviteja Anantha, Svitlana Vakulenko, Zhucheng Tu,  
621 Shayne Longpre, Stephen Pulman, and Srinivas  
622 Chappidi. 2021. [Open-domain question answering  
623 goes conversational via question rewriting](#). In *Pro-  
624 ceedings of the 2021 Conference of the North Amer-  
625 ican Chapter of the Association for Computational  
626 Linguistics: Human Language Technologies*, pages  
627 520–534, Online. Association for Computational Lin-  
628 guistics.

629 Rohan Anil, Andrew M. Dai, Orhan Firat, Melvin John-  
630 son, Dmitry Lepikhin, Alexandre Passos, Siamak  
631 Shakeri, Emanuel Taropa, Paige Bailey, Zhifeng  
632 Chen, Eric Chu, Jonathan H. Clark, Laurent El  
633 Shafey, Yanping Huang, Kathy Meier-Hellstern, Gau-  
634 rav Mishra, Erica Moreira, Mark Omernick, Kevin  
635 Robinson, Sebastian Ruder, Yi Tay, Kefan Xiao,  
636 Yuanzhong Xu, Yujing Zhang, Gustavo Hernandez  
637 Abrego, Junwhan Ahn, Jacob Austin, Paul Barham,  
638 Jan Botha, James Bradbury, Siddhartha Brahma,  
639 Kevin Brooks, Michele Catasta, Yong Cheng, Colin  
640 Cherry, Christopher A. Choquette-Choo, Aakanksha  
641 Chowdhery, Clément Crepy, Shachi Dave, Mostafa

Dehghani, Sunipa Dev, Jacob Devlin, Mark Díaz,  
Nan Du, Ethan Dyer, Vlad Feinberg, Fangxiaoyu  
Feng, Vlad Fienber, Markus Freitag, Xavier Gar-  
cia, Sebastian Gehrmann, Lucas Gonzalez, Guy Gur-  
Ari, Steven Hand, Hadi Hashemi, Le Hou, Joshua  
Howland, Andrea Hu, Jeffrey Hui, Jeremy Hur-  
witz, Michael Isard, Abe Ittycheriah, Matthew Jagiel-  
ski, Wenhao Jia, Kathleen Kenealy, Maxim Krikun,  
Sneha Kudugunta, Chang Lan, Katherine Lee, Ben-  
jamin Lee, Eric Li, Music Li, Wei Li, YaGuang Li,  
Jian Li, Hyeontaek Lim, Hanzhao Lin, Zhongtao Liu,  
Frederick Liu, Marcello Maggioni, Aroma Mahendru,  
Joshua Maynez, Vedant Misra, Maysam Moussalem,  
Zachary Nado, John Nham, Eric Ni, Andrew Nys-  
trom, Alicia Parrish, Marie Pellat, Martin Polacek,  
Alex Polozov, Reiner Pope, Siyuan Qiao, Emily Reif,  
Bryan Richter, Parker Riley, Alex Castro Ros, Au-  
rko Roy, Brennan Saeta, Rajkumar Samuel, Renee  
Shelby, Ambrose Slone, Daniel Smilkov, David R.  
So, Daniel Sohn, Simon Tokumine, Dasha Valter,  
Vijay Vasudevan, Kiran Vodrahalli, Xuezhi Wang,  
Pidong Wang, Zirui Wang, Tao Wang, John Wiet-  
ing, Yuhuai Wu, Kelvin Xu, Yunhan Xu, Linting  
Xue, Pengcheng Yin, Jiahui Yu, Qiao Zhang, Steven  
Zheng, Ce Zheng, Weikang Zhou, Denny Zhou, Slav  
Petrov, and Yonghui Wu. 2023. [Palm 2 technical  
report](#).

Namo Bang, Jeehyun Lee, and Myoung-Wan Koo.  
2023a. [Task-optimized adapters for an end-to-end  
task-oriented dialogue system](#). In *Findings of the As-  
sociation for Computational Linguistics: ACL 2023*,  
pages 7355–7369, Toronto, Canada. Association for  
Computational Linguistics.

Yejin Bang, Samuel Cahyawijaya, Nayeon Lee, Wen-  
liang Dai, Dan Su, Bryan Wilie, Holy Lovenia, Ziwei  
Ji, Tiezheng Yu, Willy Chung, Quyet V. Do, Yan  
Xu, and Pascale Fung. 2023b. [A multitask, multilin-  
gual, multimodal evaluation of chatgpt on reasoning,  
hallucination, and interactivity](#).

Edward Beeching, Clémentine Fourrier, Nathan Habib,  
Sheon Han, Nathan Lambert, Nazneen Rajani, Omar  
Sanseviero, Lewis Tunstall, and Thomas Wolf. 2023.  
Open llm leaderboard. [https://huggingface.co/  
spaces/HuggingFaceH4/open\\_llm\\_leaderboard](https://huggingface.co/spaces/HuggingFaceH4/open_llm_leaderboard).

Bernd Bohnet, Chris Alberti, and Michael Collins. 2023.  
[Coreference resolution through a seq2seq transition-  
based system](#). *Transactions of the Association for  
Computational Linguistics*, 11:212–226.

Michael Boratko, Harshit Padigela, Divyendra Mikki-  
lineni, Pritish Yuvraj, Rajarshi Das, Andrew McCal-  
lum, Maria Chang, Achille Fokoue-Nkoutche, Pavan  
Kapanipathi, Nicholas Mattei, Ryan Musa, Kartik  
Talamadupula, and Michael Witbrock. 2018. [A sys-  
tematic classification of knowledge, reasoning, and  
context within the ARC dataset](#). In *Proceedings of  
the Workshop on Machine Reading for Question An-  
swering*, pages 60–70, Melbourne, Australia. Associ-  
ation for Computational Linguistics.

Tom Brown, Benjamin Mann, Nick Ryder, Melanie  
Subbiah, Jared D Kaplan, Prafulla Dhariwal, Arvind

702	Neelakantan, Pranav Shyam, Girish Sastry, Amanda	<i>Short Papers</i> ), pages 2924–2936, Minneapolis, Min-	763
703	Askill, Sandhini Agarwal, Ariel Herbert-Voss,	nesota. Association for Computational Linguistics.	764
704	Gretchen Krueger, Tom Henighan, Rewon Child,		
705	Aditya Ramesh, Daniel Ziegler, Jeffrey Wu, Clemens	Peter Clark, Isaac Cowhey, Oren Etzioni, Tushar Khot,	765
706	Winter, Chris Hesse, Mark Chen, Eric Sigler, Ma-	Ashish Sabharwal, Carissa Schoenick, and Oyvind	766
707	teusz Litwin, Scott Gray, Benjamin Chess, Jack	Tafjord. 2018. <a href="#">Think you have solved question an-</a>	767
708	Clark, Christopher Berner, Sam McCandlish, Alec	<a href="#">swering? try arc, the ai2 reasoning challenge.</a>	768
709	Radford, Ilya Sutskever, and Dario Amodei. 2020a.		
710	<a href="#">Language models are few-shot learners.</a> In <i>Ad-</i>	Tim Dettmers, Artidoro Pagnoni, Ari Holtzman, and	769
711	<i>Advances in Neural Information Processing Systems,</i>	Luke Zettlemoyer. 2023. Qlora: Efficient finetuning	770
712	volume 33, pages 1877–1901. Curran Associates,	of quantized llms. <i>arXiv preprint arXiv:2305.14314.</i>	771
713	Inc.		
714	Tom B. Brown, Benjamin Mann, Nick Ryder, Melanie	Zhengxiao Du, Yujie Qian, Xiao Liu, Ming Ding,	772
715	Subbiah, Jared Kaplan, Prafulla Dhariwal, Arvind	Jiezhong Qiu, Zhilin Yang, and Jie Tang. 2022. Glm:	773
716	Neelakantan, Pranav Shyam, Girish Sastry, Amanda	General language model pretraining with autoregres-	774
717	Askill, Sandhini Agarwal, Ariel Herbert-Voss,	sive blank infilling. In <i>Proceedings of the 60th An-</i>	775
718	Gretchen Krueger, Tom Henighan, Rewon Child,	<i>annual Meeting of the Association for Computational</i>	776
719	Aditya Ramesh, Daniel M. Ziegler, Jeffrey Wu,	<i>Linguistics (Volume 1: Long Papers)</i> , pages 320–335.	777
720	Clemens Winter, Christopher Hesse, Mark Chen, Eric		
721	Sigler, Mateusz Litwin, Scott Gray, Benjamin Chess,	Ahmed Elgohary, Denis Peskov, and Jordan Boyd-	778
722	Jack Clark, Christopher Berner, Sam McCandlish,	Graber. 2019. <a href="#">Can you unpack that? learning to</a>	779
723	Alec Radford, Ilya Sutskever, and Dario Amodei.	<a href="#">rewrite questions-in-context.</a> In <i>Proceedings of the</i>	780
724	2020b. <a href="#">Language models are few-shot learners.</a>	<i>2019 Conference on Empirical Methods in Natu-</i>	781
		<i>ral Language Processing and the 9th International</i>	782
		<i>Joint Conference on Natural Language Processing</i>	783
		<i>(EMNLP-IJCNLP)</i> , pages 5918–5924, Hong Kong,	784
		China. Association for Computational Linguistics.	785
725	Paweł Budzianowski, Tsung-Hsien Wen, Bo-Hsiang		
726	Tseng, Iñigo Casanueva, Stefan Ultes, Osman Ra-	Elias Frantar and Dan Alistarh. 2023. SparseGPT: Mas-	786
727	madan, and Milica Gašić. 2018. <a href="#">MultiWOZ - a large-</a>	sive language models can be accurately pruned in	787
728	<a href="#">scale multi-domain Wizard-of-Oz dataset for task-</a>	one-shot. <i>arXiv preprint arXiv:2301.00774.</i>	788
729	<a href="#">oriented dialogue modelling.</a> In <i>Proceedings of the</i>		
730	<i>2018 Conference on Empirical Methods in Natural</i>	Elias Frantar, Saleh Ashkboos, Torsten Hoefler, and	789
731	<i>Language Processing</i> , pages 5016–5026, Brussels,	Dan Alistarh. 2022. GPTQ: Accurate post-training	790
732	Belgium. Association for Computational Linguistics.	compression for generative pretrained transformers.	791
		<i>arXiv preprint arXiv:2210.17323.</i>	792
733	Aakanksha Chowdhery, Sharan Narang, Jacob Devlin,		
734	Maarten Bosma, Gaurav Mishra, Adam Roberts,	Michael Heck, Nurul Lubis, Benjamin Ruppik, Renato	793
735	Paul Barham, Hyung Won Chung, Charles Sutton,	Vukovic, Shutong Feng, Christian Geisshauer, Hsien-	794
736	Sebastian Gehrmann, Parker Schuh, Kensen Shi,	chin Lin, Carel van Niekerk, and Milica Gasic. 2023.	795
737	Sasha Tsvyashchenko, Joshua Maynez, Abhishek	<a href="#">ChatGPT for zero-shot dialogue state tracking: A</a>	796
738	Rao, Parker Barnes, Yi Tay, Noam Shazeer, Vin-	<a href="#">solution or an opportunity?</a> In <i>Proceedings of the</i>	797
739	odkumar Prabhakaran, Emily Reif, Nan Du, Ben	<i>61st Annual Meeting of the Association for Compu-</i>	798
740	Hutchinson, Reiner Pope, James Bradbury, Jacob	<i>tational Linguistics (Volume 2: Short Papers)</i> , pages	799
741	Austin, Michael Isard, Guy Gur-Ari, Pengcheng Yin,	936–950, Toronto, Canada. Association for Compu-	800
742	Toju Duke, Anselm Levskaya, Sanjay Ghemawat,	tational Linguistics.	801
743	Sunipa Dev, Henryk Michalewski, Xavier Garcia,		
744	Vedant Misra, Kevin Robinson, Liam Fedus, Denny	Dan Hendrycks, Collin Burns, Steven Basart, Andy	802
745	Zhou, Daphne Ippolito, David Luan, Hyeontaek Lim,	Zou, Mantas Mazeika, Dawn Song, and Jacob Stein-	803
746	Barret Zoph, Alexander Spiridonov, Ryan Sepassi,	hardt. 2021a. Measuring massive multitask language	804
747	David Dohan, Shivani Agrawal, Mark Omernick, An-	understanding. <i>Proceedings of the International Con-</i>	805
748	drew M. Dai, Thanumalayan Sankaranarayanan Pil-	<i>ference on Learning Representations (ICLR).</i>	806
749	lai, Marie Pellat, Aitor Lewkowycz, Erica Moreira,		
750	Rewon Child, Oleksandr Polozov, Katherine Lee,	Dan Hendrycks, Collin Burns, Steven Basart, Andy Zou,	807
751	Zongwei Zhou, Xuezhi Wang, Brennan Saeta, Mark	Mantas Mazeika, Dawn Song, and Jacob Steinhardt.	808
752	Diaz, Orhan Firat, Michele Catasta, Jason Wei, Kathy	2021b. <a href="#">Measuring massive multitask language un-</a>	809
753	Meier-Hellstern, Douglas Eck, Jeff Dean, Slav Petrov,	<a href="#">derstanding.</a>	810
754	and Noah Fiedel. 2022. <a href="#">Palm: Scaling language mod-</a>		
755	<a href="#">eling with pathways.</a>	Jordan Hoffmann, Sebastian Borgeaud, Arthur Mensch,	811
		Elena Buchatskaya, Trevor Cai, Eliza Rutherford,	812
756	Christopher Clark, Kenton Lee, Ming-Wei Chang,	Diego de Las Casas, Lisa Anne Hendricks, Johannes	813
757	Tom Kwiatkowski, Michael Collins, and Kristina	Welbl, Aidan Clark, Tom Hennigan, Eric Noland,	814
758	Toutanova. 2019. <a href="#">BoolQ: Exploring the surprising</a>	Katie Millican, George van den Driessche, Bogdan	815
759	<a href="#">difficulty of natural yes/no questions.</a> In <i>Proceedings</i>	Damoc, Aurelia Guy, Simon Osindero, Karen Si-	816
760	<i>of the 2019 Conference of the North American Chap-</i>	mony, Erich Elsen, Jack W. Rae, Oriol Vinyals,	817
761	<i>ter of the Association for Computational Linguistics:</i>	and Laurent Sifre. 2022. <a href="#">Training compute-optimal</a>	818
762	<i>Human Language Technologies, Volume 1 (Long and</i>	<a href="#">large language models.</a>	819

820	Yushi Hu, Chia-Hsuan Lee, Tianbao Xie, Tao Yu, Noah A. Smith, and Mari Ostendorf. 2022. <a href="#">In-context learning for few-shot dialogue state tracking</a> . In <i>Findings of the Association for Computational Linguistics: EMNLP 2022</i> , pages 2627–2643, Abu Dhabi, United Arab Emirates. Association for Computational Linguistics.	878
821		879
822		880
823		881
824		882
825		883
826		884
827	Mandar Joshi, Danqi Chen, Yinhan Liu, Daniel S. Weld, Luke Zettlemoyer, and Omer Levy. 2020. <a href="#">SpanBERT: Improving pre-training by representing and predicting spans</a> . <i>Transactions of the Association for Computational Linguistics</i> , 8:64–77.	885
828		886
829		887
830		888
831		889
832	Shafiq Joty, Giuseppe Carenini, Raymond Ng, and Gabriel Murray. 2019. <a href="#">Discourse analysis and its applications</a> . In <i>Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics: Tutorial Abstracts</i> , pages 12–17, Florence, Italy. Association for Computational Linguistics.	890
833		891
834		892
835		893
836		894
837		895
838	Jeonghoon Kim, Jung Hyun Lee, Sungdong Kim, Joon-suk Park, Kang Min Yoo, Se Jung Kwon, and Dong-soo Lee. 2023. <a href="#">Memory-efficient fine-tuning of compressed large language models via sub-4-bit integer quantization</a> .	896
839		897
840		898
841		899
842		900
843	Andrey Kuzmin, Markus Nagel, Mart van Baalen, Arash Behboodi, and Tijmen Blankevoort. 2023. <a href="#">Pruning vs quantization: Which is better?</a>	901
844		902
845		903
846	Tom Kwiatkowski, Jennimaria Palomaki, Olivia Redfield, Michael Collins, Ankur Parikh, Chris Alberti, Danielle Epstein, Illia Polosukhin, Jacob Devlin, Kenton Lee, Kristina Toutanova, Llion Jones, Matthew Kelcey, Ming-Wei Chang, Andrew M. Dai, Jakob Uszkoreit, Quoc Le, and Slav Petrov. 2019. <a href="#">Natural questions: A benchmark for question answering research</a> . <i>Transactions of the Association for Computational Linguistics</i> , 7:452–466.	904
847		905
848		906
849		907
850		908
851		909
852		910
853		911
854		912
855	Md Tahmid Rahman Laskar, M Saiful Bari, Mizanur Rahman, Md Amran Hossen Bhuiyan, Shafiq Joty, and Jimmy Huang. 2023. <a href="#">A systematic study and comprehensive evaluation of ChatGPT on benchmark datasets</a> . In <i>Findings of the Association for Computational Linguistics: ACL 2023</i> , pages 431–469, Toronto, Canada. Association for Computational Linguistics.	913
856		914
857		915
858		916
859		917
860		918
861		919
862		920
863	Nghia T. Le and Alan Ritter. 2023. <a href="#">Are large language models robust zero-shot coreference resolvers?</a>	921
864		922
865	Kenton Lee, Luheng He, Mike Lewis, and Luke Zettlemoyer. 2017. <a href="#">End-to-end neural coreference resolution</a> . In <i>Proceedings of the 2017 Conference on Empirical Methods in Natural Language Processing</i> , pages 188–197, Copenhagen, Denmark. Association for Computational Linguistics.	923
866		924
867		925
868		926
869		927
870		928
871	Hector J. Levesque, Ernest Davis, and Leora Morgenstern. 2012. The winograd schema challenge. In <i>13th International Conference on the Principles of Knowledge Representation and Reasoning, KR 2012</i> , Proceedings of the International Conference on Knowledge Representation and Reasoning, pages 552–561. Institute of Electrical and Electronics Engineers Inc.	929
872		930
873		931
874		932
875		933
876		934
877		
	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 Yuksekogul, 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. <a href="#">Holistic evaluation of language models</a> .	
	Chin-Yew Lin. 2004. <a href="#">ROUGE: A package for automatic evaluation of summaries</a> . In <i>Text Summarization Branches Out</i> , pages 74–81, Barcelona, Spain. Association for Computational Linguistics.	
	Wei Liu and Michael Strube. 2023. <a href="#">Annotation-inspired implicit discourse relation classification with auxiliary discourse connective generation</a> . In <i>Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)</i> , pages 15696–15712, Toronto, Canada. Association for Computational Linguistics.	
	Zechun Liu, Barlas Oguz, Changsheng Zhao, Ernie Chang, Pierre Stock, Yashar Mehdad, Yangyang Shi, Raghuraman Krishnamoorthi, and Vikas Chandra. 2023. <a href="#">Llm-qat: Data-free quantization aware training for large language models</a> .	
	Wanqiu Long, Bonnie Webber, and Deyi Xiong. 2020. <a href="#">TED-CDB: A large-scale Chinese discourse relation dataset on TED talks</a> . In <i>Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP)</i> , pages 2793–2803, Online. Association for Computational Linguistics.	
	Christopher D. Manning, Mihai Surdeanu, John Bauer, Jenny Finkel, Steven J. Bethard, and David McClosky. 2014. <a href="#">The Stanford CoreNLP natural language processing toolkit</a> . In <i>ACL 2014 System Demonstrations</i> , pages 55–60.	
	Scott Martin, Shivani Poddar, and Kartikeya Upasani. 2020. <a href="#">MuDoCo: Corpus for multidomain coreference resolution and referring expression generation</a> . In <i>Proceedings of the 12th Language Resources and Evaluation Conference</i> , pages 104–111, Marseille, France. European Language Resources Association.	
	Nasrin Mostafazadeh, Michael Roth, Annie Louis, Nathanael Chambers, and James Allen. 2017. <a href="#">LSDSem 2017 shared task: The story cloze test</a> . In <i>Proceedings of the 2nd Workshop on Linking Models of Lexical, Sentential and Discourse-level Semantics</i> , pages 46–51, Valencia, Spain. Association for Computational Linguistics.	



1050	pages 22–31, Florence, Italy. Association for Computational Linguistics.	1106
1051		1107
1052	S. Tata and J.M. Patel. 2003. <a href="#">Piqa: an algebra for querying protein data sets</a> . In <i>15th International Conference on Scientific and Statistical Database Management, 2003.</i> , pages 141–150.	1108
1053		1109
1054		1110
1055		
1056	Hugo Touvron, Thibaut Lavril, Gautier Izacard, Xavier Martinet, Marie-Anne Lachaux, Timothée Lacroix, Baptiste Rozière, Naman Goyal, Eric Hambro, Faisal Azhar, Aurelien Rodriguez, Armand Joulin, Edouard Grave, and Guillaume Lample. 2023. <a href="#">Llama: Open and efficient foundation language models</a> .	1111
1057		
1058		
1059		
1060		
1061		
1062	Szymon Tworowski, Konrad Staniszewski, Mikołaj Patek, Yuhuai Wu, Henryk Michalewski, and Piotr Miłoś. 2023. <a href="#">Focused transformer: Contrastive training for context scaling</a> .	1112
1063		
1064		
1065		
1066	Alex Wang, Yada Pruksachatkun, Nikita Nangia, Amanpreet Singh, Julian Michael, Felix Hill, Omer Levy, and Samuel Bowman. 2019. <a href="#">Superglue: A stickier benchmark for general-purpose language understanding systems</a> . In <i>Advances in Neural Information Processing Systems</i> , volume 32. Curran Associates, Inc.	1113
1067		
1068		
1069		
1070		
1071		
1072		
1073	Bonnie Webber, Rashmi Prasad, Alan Lee, and Aravind Joshi. 2019. The penn discourse treebank 3.0 annotation manual. <a href="https://catalog.ldc.upenn.edu/docs/LDC2019T05/PDTB3-Annotation-Manual.pdf">https://catalog.ldc.upenn.edu/docs/LDC2019T05/PDTB3-Annotation-Manual.pdf</a> .	1114
1074		
1075		
1076		
1077		
1078	Kai-Cheng Yang and Filippo Menczer. 2023. <a href="#">Large language models can rate news outlet credibility</a> .	1115
1079		
1080	Steve Young, Milica Gašić, Blaise Thomson, and Jason D Williams. 2013. Pomdp-based statistical spoken dialog systems: A review. <i>Proceedings of the IEEE</i> , 101(5):1160–1179.	1116
1081		
1082		
1083		
1084	Zhihang Yuan, Lin Niu, Jiawei Liu, Wenyu Liu, Xinggang Wang, Yuzhang Shang, Guangyu Sun, Qiang Wu, Jiaxiang Wu, and Bingzhe Wu. 2023. <a href="#">Rptq: Reorder-based post-training quantization for large language models</a> .	1117
1085		
1086		
1087		
1088		
1089	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. In <i>Proceedings of the 2nd Workshop on Natural Language Processing for Conversational AI, ACL 2020</i> , pages 109–117.	1118
1090		
1091		
1092		
1093		
1094		
1095		
1096	Rowan Zellers, Ari Holtzman, Yonatan Bisk, Ali Farhadi, and Yejin Choi. 2019. <a href="#">Hellaswag: Can a machine really finish your sentence?</a>	1119
1097		
1098		
1099	Susan Zhang, Stephen Roller, Naman Goyal, Mikel Artetxe, Moya Chen, Shuohui Chen, Christopher Dewan, Mona Diab, Xian Li, Xi Victoria Lin, Todor Mihaylov, Myle Ott, Sam Shleifer, Kurt Shuster, Daniel Simig, Punit Singh Koura, Anjali Sridhar, Tianlu Wang, and Luke Zettlemoyer. 2022. <a href="#">Opt: Open pre-trained transformer language models</a> .	1120
1100		
1101		
1102		
1103		
1104		
1105		
	Qi Zhu, Kaili Huang, Zheng Zhang, Xiaoyan Zhu, and Minlie Huang. 2020. <a href="#">CrossWOZ: A large-scale Chinese cross-domain task-oriented dialogue dataset</a> . <i>Transactions of the Association for Computational Linguistics</i> , 8:281–295.	1121
		1122
		1123
		1124
		1125
		1126
	<b>A Task Design Examples</b>	1111
	Table 11 presents the input example for each task. For CR, we only show examples from OntoNotes.	1112
		1113
	<b>B Few-shot Settings</b>	1114
	Table 12 shows the number of examples for each dataset that yields the best scores. All datasets except WSC273 and PDTB3 use randomly selected examples from the training set. Since WSC273 does not include any train or validation set, we use the zero-shot setting, as scores presented in Table 6. For each class in PDTB3, we randomly select two examples from the training set for prompting. For some particular datasets, such as OntoNotes, experiments are only performed in the zero-shot and one-shot settings due to the limitation on input length.	1115
		1116
		1117
		1118
		1119
		1120
		1121
		1122
		1123
		1124
		1125
		1126

### Coreference Resolution

Instructions: Please carefully read the following passages. For each passage and the options, you must identify which option the mention marked in \*bold\* refers to. If the marked mention does not have any antecedent, please select "no antecedent".

[Few-shot examples]

Context: — basically , it was unanimously agreed upon by the various relevant parties . To express \*its\* determination , the Chinese securities regulatory department compares this stock reform to a die that has been cast . It takes time to prove whether the stock reform can really meet expectations , and whether any deviations that arise during the stock reform can be promptly corrected . Dear viewers , the China News program will end here . This is Xu Li . Thank you everyone for watching . Coming up is the Focus Today program hosted by Wang Shilin . Good-bye , dear viewers .

Choice:

- A. the Chinese securities regulatory department
- B. this stock reform
- C. the stock reform
- D. you
- E. everyone
- F. no antecedent

Question: What does \*its\* refers to?

Answer: A

---

### Dialogue State Tracking

Consider the following list of concepts, called "slots" provided to you as a json list.

```
"slots": { "restaurant-pricerange": "price budget for the restaurant",
  "restaurant-area": "area or place of the restaurant",
  "restaurant-food": "the cuisine of the restaurant you are looking for",
  ...
  "hotel-postcode": "postal code of the hotel",
  "hotel-ref": "reference number of the hotel booking"
}
```

Some "slots" can only take a value from predefined list:

```
"categorical": { "restaurant-pricerange": [ 'cheap', 'expensive', 'moderate' ],
  "restaurant-area": [ 'centre', 'east', 'north', 'south', 'west' ],
  "restaurant-bookday": [ 'monday', 'tuesday', 'wednesday', 'thursday', 'friday', 'saturday', 'sunday' ],
  ...
  "hotel-internet": [ 'free', 'no', 'yes' ], "hotel-area": [ 'centre', 'east', 'north', 'south', 'west' ]
}
```

Now consider the following dialogue between two parties called the "system" and "user". Can you tell me which of the "slot" was updated by the "user" in its latest response to the "system"? Present the updates in JSON format. If no "slots" were updated, return an empty JSON list. If you encounter "slot" that was requested by the "user" then fill them with "?". If a user does not seem to care about a discussed "slot" fill it with "dontcare".

Input:

Previous state: {}

"system": ""

"user": "I'm looking for a moderately priced place to eat that's in the centre of town."

Output: { "restaurant-pricerange": "moderate", "restaurant-area": "centre" }

---

### Implicit Discourse Relation Classification

Instructions: Given two arguments and a list of connective words, please select the most likely connective between two arguments.

Below are the descriptions of four discourse relation labels. Please find the correct label for each example.

Temporal: The tag temporal is used when the situations described in the arguments are intended to be related temporally.

Contingency: The tag Contingency is used when the situation described by one argument provides the reason, explanation or justification for the situation described by the other.

Comparison: The tag Comparison is used when the discourse relation between two arguments highlights their differences or similarities, including differences between expected consequences and actual ones.

Expansion: The label Expansion is used for relations that expand the discourse and move its narrative or exposition forward.

[Few-shot examples]

Input:

Arg 1: Amcore, also a bank holding company, has assets of \$1.06 billion.

Arg 2: Central's assets are \$240 million.

Question: What is the connective that best describes the relation between two arguments?

- A. Temporal
- B. Contingency
- C. Comparison
- D. Expansion

Answer: C

---

### Query Rewrite

Instructions: Rewrite the last query following interaction into a well-formed, context independent query. Resolve any disfluencies or grammatical errors in the query.

[Few-shot examples]

Input:

User: Try to reach Forbes now .

Bot: Forbes at Washington Post ? Or Forbes of Publishing Division ?

User: Publishing Division .

Rewrite: *Forbes of Publishing Division*

Table 11: Examples of task design for each task in the context understanding benchmark.

Task	Coreference		DST	Discourse	Query Rewrite				
	WSC273	OntoNotes	MultiWOZ	PDTB3	MuDoCo	QReCC	InCar	GECOR	CANARD
N-example	Zero-shot	One-shot	5-shot	8-shot	10-shot	5-shot	10-shot	10-shot	5-shot

Table 12: N-shot settings for each task & dataset that yields the highest scores. For each task and model, we use consistent N-shot settings for comparison.