

INTENTIONQA: A Benchmark for Evaluating Purchase Intention Comprehension Abilities of Large Language Models in E-commerce

Anonymous ACL submission

Abstract

Enhancing Large Language Models' (LLMs) ability to understand purchase intentions in E-commerce scenarios is crucial for their effective assistance in various downstream tasks. However, previous approaches that distill intentions from LLMs often fail to generate meaningful and human-centric intentions applicable in real-world E-commerce contexts. This raises concerns about the true comprehension and utilization of purchase intentions by LLMs. In this paper, we present INTENTIONQA, a double-task multiple-choice question answering benchmark to evaluate LLMs' comprehension of purchase intentions in E-commerce. Specifically, LLMs are tasked to infer intentions based on purchased products and utilize them to predict additional purchases. INTENTIONQA consists of 4,375 carefully curated problems across three difficulty levels, constructed using an automated pipeline to ensure scalability on large E-commerce platforms. Human evaluations demonstrate the high quality and low false-negative rate of our benchmark. Extensive experiments across 19 language models show that they still struggle with certain scenarios, such as understanding products and intentions accurately, jointly reasoning with products and intentions, and more, in which they fall far behind human performances.

1 Introduction

Understanding customers' purchase intentions and making reasonable inferences accordingly are crucial for revolutionizing E-commerce services, whose benefits have been demonstrated in myriads of downstream tasks, such as product recommendation (Grbovic et al., 2015; Zhao et al., 2014; Li et al., 2020) and query answering (Zhao et al., 2019; Hirsch et al., 2020). However, intention comprehension (Fogassi et al., 2005) is a non-trivial task as it involves reasoning with implicit mental states, which are not typically expressed in text or conversations. Thus, in the context of E-commerce, ex-

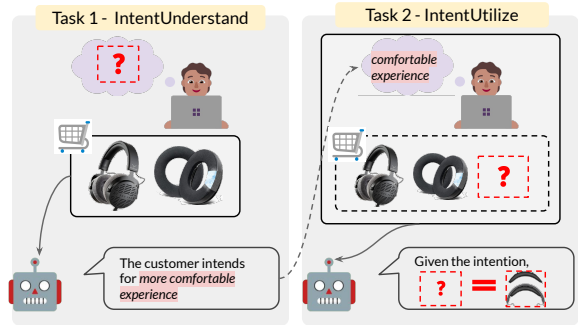


Figure 1: Examples of two tasks in INTENTIONQA. Task 1 requires the language model to determine the customer's intention in purchasing two products, and Task 2 involves recommending a product that fulfills the customer's intention and matches their currently purchased product.

tracting purchase intentions from behaviors without explicit external cues has been challenging (Yang and Tang, 2015).

Recently, Yu et al. (2023, 2024) proposed to distill purchase intentions from LLMs, such as OPT (Zhang et al., 2022b), by leveraging their inherent advantages of generative and commonsense reasoning abilities, as well as being pre-trained on vast textual data including E-commerce knowledge. However, recent analyses by Zhou et al. (2024) show that LLMs struggle to generate meaningful and user-centric intentions. Instead, they are biased by over-focusing on similarities among different products' metadata, such as their properties, and often end up regurgitating information from the provided prompts without truly comprehending the underlying purchase intentions.

Thus, an important yet under-explored question arises: *Can LLMs comprehend the customers' purchase intention and how effective are they in performing such tasks?* To dive into this, we first break down the comprehension of intention into two key aspects, as shown in Figure 1. First, we have **intention understanding**, which evaluates LLMs' capacity to accurately infer customers' purchase in-

tentions based on the products bought. Second, we consider **intention utilization**, which investigates LLMs’ ability to predict additional purchases based on customer’s intentions. Together, they make up the entire process of intention comprehension and play a significant role in enhancing E-commerce search services.

Although LLMs have been extensively used in intention knowledge distillation, their actual performances in this area have not been adequately benchmarked. This is because current methods that leverage LLMs have been adopting an open-ended generation fashion, which is difficult to consistently evaluate (Gu et al., 2021). Additionally, the extensive and constantly growing number of products on E-commerce platforms makes it infeasible and expensive to construct human-curated benchmarks.

To address these challenges and benchmark LLMs on purchase intention comprehension in E-commerce, we introduce INTENTIONQA, a double-task multiple-choice question answering (MCQA) dataset, featuring intention understanding and intention utilization respectively. INTENTIONQA contains 4,375 problems for two tasks and covers varying difficulty levels, allowing for fine-grained evaluation. The MCQA setting enables using consistent evaluation metrics to assess the LLMs’ intention comprehension abilities.

Specifically, we design a pipeline that automatically synthesizes QA pairs by transforming human-annotated intentions from FolkScope (Yu et al., 2023), each involving a pair of co-buy products and the corresponding intention of purchasing them, into questions by masking out the intention or one of the products. To achieve this, we define context-based product similarity and intention similarity metrics. They are computed over ASER (Zhang et al., 2022a), a large-scale eventuality knowledge graph, which we leverage as a reference for our automatic distractor sampling strategy. For each question, we include 3 negative distractors alongside the gold answer through a strict similarity filtering process. We then assign difficulty labels to each QA pair based on the product similarity between the co-buy products in the original intention assertion. These steps are done without human supervision, enabling our benchmark construction pipeline to generalize and accommodate larger-scale product databases and practical applications.

We further conduct human evaluations to demonstrate the high quality and low false-negative rate

of INTENTIONQA, followed by extensive experiments across 19 language models with varying sizes and approaches. Results demonstrate that the existing language models still struggle with certain scenarios, such as understanding products and intentions accurately, jointly reasoning with the products and intentions, and more. In the long run, we hope that our benchmark serves as an important cornerstone toward intention-aware E-commerce services that promote integrating intention reasoning abilities into product recommendations.

2 Related Works

2.1 Intention Discovery with Large Language Models in E-commerce

Understanding intentions with language models have been studied in various domains, such as smoothing chatbox conversations (Ouyang et al., 2022), enhancing web search (Zhang et al., 2019), and more. In the E-commerce domain, understanding customers’ purchase intentions benefits various downstream tasks (Koo and Ju, 2010; Xu et al., 2024), such as automated on-call customer support (Goyal et al., 2022), recommendation systems (Dai et al., 2006; Qian et al., 2023; Jung et al., 2023), product question answering (Deng et al., 2023; Yu and Lam, 2018). While Yu et al. (2023, 2024) proposed leveraging the generation abilities of LLMs to distill purchase intentions from co-buy records, Zhou et al. (2024) showed that LLMs struggle with generating meaningful intentions or understanding user-centric intentions. In this work, we construct INTENTIONQA, a benchmark to evaluate LLMs’ intention comprehension abilities by selecting highly typical intentions in previously available resources and provide insights for human-centric intention comprehension.

2.2 Benchmarking Large Language Models

Since the emergence of LLMs, various studies have explored their capabilities in various domains, including temporal reasoning (Tan et al., 2023), causal reasoning (Chan et al., 2024), commonsense reasoning (Jain et al., 2023), and more (Qin et al., 2023). These benchmarks have made significant contributions to the understanding of large language models, assessing their performance across different parameters and prompting methods. However, in the field of E-commerce, existing works primarily leverage LLMs with explicit instruction-tuning (Li et al., 2024), while neglecting the infeasibility

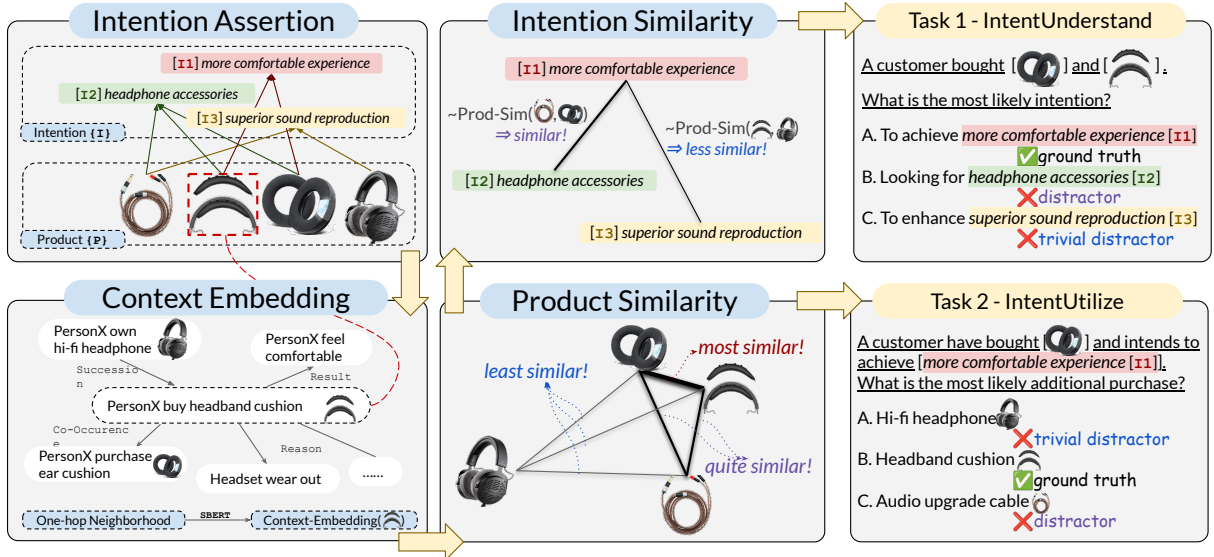


Figure 2: Overview of INTENTIONQA and the construction pipeline. We map products from intention assertions to event nodes in ASER and calculate their context embedding with the one-hop neighborhood. Product and intention similarities are then computed accordingly. Products/intentions with higher similarities are represented closer to each other. Negative distractor sampling for Task 1/2 is based on intention/product similarity respectively.

bility of directly applying LLMs in a generalizable manner. Furthermore, current evaluation benchmarks in E-commerce primarily emphasize product and session comprehension (Jin et al., 2023), which overlooks the important aspect of intention comprehension. In this paper, we step forward by presenting the first benchmark that evaluates the intention comprehension abilities of LLMs.

3 INTENTIONQA

In this section, we introduce INTENTIONQA, a Multiple-Choice Question Answering (MCQA) benchmark consisting of two tasks targeting different aspects of purchase intention comprehension and with progressive difficulties, to evaluate the intention understanding and utilization abilities of LLMs thoroughly.

3.1 Task Definitions

We begin by formally defining two tasks associated with INTENTIONQA.

Task 1: INTENTUNDERSTAND The first task examines whether LLMs can infer the purchase intentions correctly given a real-world record of the products bought. Formally, given a pair of co-buy products p_1, p_2 , LLMs are tasked with selecting the most likely purchase intention i^* from a list of candidate options $\mathcal{I} = [i_1, i_2, \dots, i_{|\mathcal{I}|}]$.

Task 2: INTENTUTILIZE The second task looks further into the capacity of LLMs to utilize pur-

chase intention for the product recommendation process. We approach this by examining their abilities to predict the most likely additional purchase based on customer intention. Specifically, given the purchase intention i^* and one product that has been Bought p^B , the LLMs are tasked with selecting the most likely Additional purchase p^{A*} from a list of candidate options $\mathcal{P}^A = [p_1^A, p_2^A, \dots, p_{|\mathcal{P}^A|}^A]$.

3.2 Source Intention Collection and Context Augmentation

We collect co-buy products and intention assertions from FolkScope (Yu et al., 2023) as our source data. FolkScope is an intention knowledge base that is constructed by distilling knowledge from a pre-trained large language model, OPT (Zhang et al., 2022b). It associates customers' co-purchase behaviors with their purchase intentions, as shown in the upper left part of Figure 2. Two scores are also assigned to each intention, indicating its plausibility and typicality. To accommodate our tasks, we preprocess FolkScope by filtering and retaining plausible assertions with typicality scores above 0.5. This is to minimize the number of overly-general intentions, which may be plausible for most products but are not specifically related to the given products. Including these intentions in INTENTIONQA could lead to many false negative distractors, which harms the quality of our QA pairs.

Since we are aiming for automatic QA pair construction, determining the similarity between differ-

ent intentions and products can serve as powerful hints in selecting appropriate distractors given a correct answer. However, relying solely on product metadata and corresponding purchase behavior falls short of capturing the similarity between intentions, as similar or identical intentions can align with multiple products. To address this limitation and enhance the sampling of distractors while reducing the occurrence of false-negative distractors, we introduce a method to augment customers’ purchase behavior. This is achieved by retrieving additional relevant context from ASER (Zhang et al., 2020, 2022a), a large-scale eventuality knowledge graph that covers billions of commonly seen eventualities. Specifically, we first consider the purchasing event as an eventuality and design heuristic rules to align it with nodes in ASER. Formally, we denote ASER as $\mathcal{G} = \{(h, r, t) | h, t \in \mathcal{V}, r \in \mathcal{R}\}$, where \mathcal{V} and \mathcal{R} are the sets of nodes and relations in ASER. Inspired by recent works in conceptualization (Wang et al., 2023b,a), we simplify the product name p by instructing ChatGPT to conceptualize it into three plausible categories $\mathcal{C}(p) = [c_1, c_2, c_3]$, representing the possible classifications of the product. For example, *iPhone 14* can be conceptualized as a *phone*, *communication device*, and *Apple product*. This augmentation expands the semantic coverage of the purchasing event, increasing the likelihood of finding relevant nodes in ASER.

Next, we design natural language templates (Appendix B) to convert noun phrases of conceptualized product categories into purchasing events $\mathcal{E}(\mathcal{C}) = [e_1, e_2, \dots, e_{|\mathcal{E}|}]$. These events are then matched against nodes in ASER to identify overlapping ones through strict string matching. Formally, we denote ASER as $\mathcal{G} = \{(h, r, t) | h, t \in \mathcal{V}, r \in \mathcal{R}\}$, where \mathcal{V} and \mathcal{R} are the sets of nodes and relations in ASER. We denote the set of matched nodes for p in ASER as $\mathcal{V}_{\mathcal{E}} = \mathcal{E}(\mathcal{C}) \cap \mathcal{V} = [ve_1, ve_2, \dots, ve_{|\mathcal{V}_{\mathcal{E}}|}]$. Next, we compute the sentence embedding of edges in the one-hop neighborhood of each node in $\mathcal{V}_{\mathcal{E}}$ using Sentence-BERT (Reimers and Gurevych, 2019). The context embedding $CE(p)$ is then computed by averaging these embeddings, which serves as the semantic representation of relevant contexts for purchasing the product p .

3.3 Product and Intention Similarity

With the context embeddings of all products computed, they contain valuable background knowl-

edge about purchasing events associated with each product. This includes edges from ASER that capture cause-effect relationships (“Reason” and “Result”), event precedence and succession (“Precedence” and “Succession”), and other relevant information. Intuitively, similar products should have similar contextual information in ASER, and vice versa. Thus, we define the similarity between purchasing events of p_1 and p_2 as follows:

$$Sim^{(p)}(p_1, p_2) = \text{cos_sim}(CE(p_1), CE(p_2))$$

where cos_sim is the cosine similarity between embeddings from Sentence-BERT.

Similarly, we define the similarity between two intentions (i_1, i_2) in FolkScope by leveraging their corresponding purchase events associated with ASER as follows:

$$Sim^{(i)}(i_1, i_2) = \min_{m=1,2;n=1,2} \{Sim^{(p)}(p_m^{(i_1)}, p_n^{(i_2)})\}$$

where $p_m^{(i)}$ is the m th product linked to intention i .

3.4 Distractor Sampling and QA Construction

Finally, we design INTENTIONQA as a multiple-choice QA benchmark and design specific rules to transform intention assertions in FolkScope into question and gold answer pairs. Each gold answer is accompanied by three distractors, i.e., $|\mathcal{I}|, |\mathcal{P}^A| = 4$. For each task, we propose its unique distractor sampling strategy specifically designed for the task objective, based on the similarity scores defined in §3.3.

Task 1: INTENTUNDERSTAND INTENTUNDERSTAND targets LLMs’ ability to accurately infer purchase intentions based on the products bought by a customer. We convert the intention assertions from FolkScope to questions by masking out the intentions. These masked intentions are then treated as gold answers, denoted as i^* . To obtain the distractor intentions $\mathcal{I}^- = [i_1^-, i_2^-, i_3^-]$, we randomly select intentions from FolkScope whose intention-similarity score with i^* ($Sim^{(i)}(i^*, i^-)$) fall within $[0.6, 0.9]$. The lower bound of the range filters out trivial distractors, while the upper bound minimizes the false negative rate in the resulting benchmark.

Task 2: INTENTUTILIZE INTENTUTILIZE evaluates the LLMs’ ability to utilize intentions to predict future purchase behavior. Specifically, we formulate the task as providing LLMs with one

product that the customer has bought and the corresponding intention, and task LLMs with predicting the most likely purchase accompanied by the purchased product. Questions for INTENTUTILIZE are obtained by masking out one of the products (p^{A*}) in each intention assertion of FolkScope. The distractor products ($p_i^-, i = 1, 2, 3$) are randomly selected from those products whose product-similarity score $Sim^p(p^{A*}, p^{A-})$ falls within $[0.7, 0.9]$. Threshold values for both tasks are determined through observations of the distribution and preliminary experiments.

Difficulty Labeling To allow for fine-grained evaluation, we categorize each question into three difficulty levels. Intuitively, intention assertions with high product-similarity scores among co-buy products result in relatively easy problems. This is based on the assumption that understanding just one product is sufficient for comprehending the corresponding intention, without necessitating reasoning about the relationship between the products. Conversely, intention assertions with low product-similarity scores contribute to harder problems as they require comprehending both products and their corresponding intentions, as well as reasoning about the potentially complementary relationship between the products.

Therefore, we categorize the problems based on the product-similarity scores of co-buy products in the original intention assertion. Specifically, problems with a product-similarity score within the range of $[0.85, 1]$ are classified as easy problems, those within the range of $[0.6, 0.85)$ are considered medium, and those within the range of $[0, 0.6)$ are classified as hard problems. These thresholds are determined based on distributions and human observations of problem difficulty.

4 Benchmark Evaluations

4.1 Statistics

We initially construct INTENTIONQA by using 2,315 intention assertions sourced from FolkScope. They are selected by filtering those with high plausibility and typicality scores and whose both products can be aligned with purchasing event nodes of ASER. We then construct 4,375 problems for both tasks in INTENTIONQA, with each problem labeled with difficulty accordingly. The benchmark statistics are reported in Table 1.

Subset	TASK 1		TASK 2	
	#Q	Avg. Sim^p	#Q	Avg. Sim^p
easy	1700	0.972	1625	0.971
medium	423	0.740	385	0.744
hard	118	0.532	133	0.514
Average	2241	0.905	2143	0.902

Table 1: Statistics of the INTENTIONQA. We report the number of questions (#Q) and the average product-similarity scores between the co-buy products among all intentions (Avg. Sim^p) within each difficulty subset.

4.2 Human Evaluations

To evaluate the effectiveness of our benchmark construction pipeline and assess the quality of our constructed QA benchmark, we conduct human annotation to evaluate two aspects: (1) the correctness of product conceptualization by ChatGPT and (2) the quality of the QA pairs in both tasks.

4.2.1 Annotation Setups

We recruit human annotators from the Amazon Mechanical Turk platform. For strict quality control, we only invite workers satisfying the following requirements: 1) at least 1K HITs approved, and 2) at least 95% approval rate. We then host two rounds of qualification rounds using questions sampled from our curated benchmark, with expert-annotated answers. 400 workers are invited in total and around 60 (15%) of them are selected.

For product conceptualization, we randomly sample 2,000 products and task each annotator to label the plausibility of the generated categories. Specifically, we ask the annotators to assess whether all three generated product categories are reasonable according to the original products. Each product is annotated by three annotators and the majority vote is taken as the final label. Results show that 89.4% of products are reasonably conceptualized, demonstrating the strong product understanding ability of ChatGPT and validating the feasibility of leveraging its generative power to aid our benchmark construction process.

We then evaluate the resulting QA pairs from INTENTUNDERSTAND and INTENTUTILIZE. For each task, we randomly sample 300 QA pairs and ask the annotators to assess the quality of these problems. Firstly, they need to annotate the correctness of ground truth options, denoted as *Correct*. Secondly, they assess the false-negativeness of the distractor options by determining whether a distractor option is superior to or equally plausible as the ground truth option, denoted as *F-Neg*. Still, we

Subset	TASK 1		TASK 2	
	Correct	F-Neg	Correct	F-Neg
easy	96.06	6.33	98.20	1.20
medium	94.00	1.33	92.59	4.32
hard	85.71	1.00	100.00	0.00
Average	95.00	2.89	97.33	1.67

Table 2: Annotated correctness (Correct; %) and false-negative rate (F-Neg; %) of 600 randomly sampled QA pairs from two tasks.

collect three votes for each QA pair and take the majority of them.

4.2.2 Results

We report the annotation results in Table 2. We find that INTENTIONQA exhibits high correctness rates among ground truth options. Meanwhile, the low false-negative rates demonstrate the high quality of both tasks. Both statistics validate the reliability of our automatic QA construction pipeline and the quality of the resulting INTENTIONQA benchmark.

5 Experiments and Analysis

5.1 Baseline Selection and Setup

Evaluation Metric We use accuracy as the evaluation metric, which is quantified by the percentage of QA pairs that a language model answers correctly in INTENTIONQA.

Model Selection We evaluate a wide range of (L)LMs in four categories: **(1) PTLM:** We evaluate several pre-trained language models, including RoBERTa (Liu et al., 2019), DeBERTa-v3 (He et al., 2023), T0 (Sanh et al., 2022), T5 (Raffel et al., 2020), and Flan-T5 (Chung et al., 2022). **(2) COMMONSENSE:** We also evaluate PTLMs with commonsense knowledge injected, including HyKAS (Ma et al., 2021), CAR (Wang et al., 2023a), VERA (Liu et al., 2023b), CANDLE (Wang et al., 2024), and VERA-CANDLE (Wang et al., 2024). **(3) OPEN LLM:** We then evaluate representative open-sourced LLMs of varying sizes and versions in zero-shot settings as well as after fine-tuning on intention knowledge (OPEN LLM + MIND, details in §5.5). These models cover LLaMA2 (Touvron et al., 2023), Gemma (Mesnard et al., 2024), Mistral (Jiang et al., 2023), Falcon (Almazrouei et al., 2023), Vicuna (Zheng et al., 2023), Phi-2 (Gunasekar et al., 2023), and Alpaca (Taori et al., 2023; Wang et al., 2023d). **(4) LLM API:** Finally, we adopt Chain-of-Thought prompting (CoT; Wei et al., 2022) and

CoT with Self-Consistency (CoT-SC; Wang et al., 2023c) together with zero-shot prompting to assess ChatGPT (OpenAI, 2022) and GPT-4 (OpenAI, 2023). The sampling temperature τ is set to 0.1 by default. 5 CoT responses are sampled with τ set to 0.7 under CoT-SC. RANDOM and MAJORITY voting are also added as baselines to demonstrate the characteristic of INTENTIONQA. HUMAN performance is calculated based on annotation results of 600 randomly selected QA pairs from both tasks.

5.2 Results

The results of all models are presented in Table 3. From the results, we observe that:

Commonsense knowledge does help in intention comprehension. Models injected with commonsense knowledge showcase comparable performance to significantly larger models. Specifically, CAR and CANDLE (435M) achieve 96.64% of the performance of Flan-T5-xxl (11B) in INTENTUNDERSTAND, despite being 25 times smaller. This demonstrates the effectiveness of incorporating commonsense knowledge in improving intention comprehension in the E-commerce domain.

INTENTUTILIZE is more challenging. For approximately all models, excluding ChatGPT and GPT-4, that exhibit above RANDOM performances in INTENTUNDERSTAND, their performances drop significantly when evaluated on INTENTUTILIZE, with an average accuracy gap of 14.20%. While INTENTUNDERSTAND involves understanding the purchase intention behind a single pair of products, INTENTUTILIZE requires product understanding of all candidate options as well as reasoning with potential intentions behind four pairs of products. This expanded reasoning scope and higher demand for product understanding pose challenges for these models, as their training data may be limited in terms of the variety and quantity of products included. However, ChatGPT and GPT-4 excelled in both tasks, presumably due to their stronger product reasoning abilities.

Intention comprehension abilities of current models are still far from perfect. Although various models perform considerably better than RANDOM guessing, there remains a substantial gap between their performance and that of humans.

5.3 Performances Across Intention Types

To further investigate the reasons why language models fail in intention comprehension, we conduct a more fine-grained analysis by delving into

Methods	Backbone	INTENTUNDERSTAND				INTENTUTILIZE			
		Easy	Medium	Hard	Avg.	Easy	Medium	Hard	Avg.
Random	-	25.00	25.00	25.00	25.00	25.00	25.00	25.00	25.00
Majority Vote	-	26.37	25.24	26.27	26.15	25.97	28.57	28.57	26.60
PTLM	RoBERTa-Large <i>214M</i>	41.46	41.98	38.98	41.43	54.95	35.06	30.08	49.84
	DeBERTa-v3-Large <i>435M</i>	36.40	38.72	37.62	36.90	26.52	29.35	32.33	27.39
	T5-v1.1-xxl <i>11B</i>	24.84	25.47	25.42	24.99	26.71	26.23	25.56	26.55
	Flan-T5-xxl <i>11B</i>	75.98	73.58	63.56	74.88	79.26	81.82	81.95	79.89
	T0-pp <i>11B</i>	71.70	68.87	<u>64.41</u>	70.78	77.11	76.10	78.20	76.99
Commonsense	HyKAS <i>435M</i>	71.81	67.17	46.69	69.61	47.02	45.97	48.12	46.90
	CAR <i>435M</i>	73.69	71.46	54.38	72.20	36.18	43.12	44.36	37.94
	CANDLE <i>435M</i>	<u>74.34</u>	70.75	52.54	72.52	35.94	43.90	43.61	37.84
	VERA <i>11B</i>	69.82	70.52	61.02	69.49	59.20	58.18	64.66	59.36
	VERA-CANDLE <i>11B</i>	70.59	<u>71.33</u>	<u>63.41</u>	<u>70.02</u>	<u>62.18</u>	<u>60.13</u>	<u>66.13</u>	<u>61.81</u>
Open LLM	LLaMA2-7B	22.47	26.24	21.78	23.14	26.42	27.87	29.03	26.84
	LLaMA2-7B-chat	64.98	66.54	53.85	64.61	59.90	54.86	47.37	58.04
	LLaMA2-13B	24.21	27.70	23.23	24.82	27.92	30.59	28.03	28.40
	LLaMA2-13B-chat	69.63	63.96	60.78	68.06	45.53	41.95	39.71	44.52
	Gemma-2B	21.73	23.87	19.81	22.06	30.66	30.63	30.99	30.67
	Gemma-2B-instruct	48.77	47.23	48.21	48.45	39.45	39.15	38.17	39.32
	Gemma-7B	50.94	50.86	42.61	50.48	26.75	30.19	31.20	27.65
	Gemma-7B-instruct	65.55	64.31	52.04	64.61	33.18	36.01	41.51	34.20
	Mistral-7B-instruct-v0.1	53.49	55.04	49.26	53.54	26.18	28.27	28.57	26.70
	Mistral-7B-instruct-v0.2	<u>76.57</u>	<u>74.53</u>	<u>63.56</u>	<u>75.50</u>	<u>59.78</u>	<u>62.60</u>	<u>65.41</u>	<u>60.64</u>
	Falcon-7B	24.19	20.52	23.73	23.47	25.40	25.45	27.82	25.56
	Falcon-7B-instruct	24.54	22.17	28.81	24.32	26.15	28.05	26.32	26.50
	Vicuna-7B-v1.5	57.13	57.08	50.85	56.79	27.88	30.13	23.31	28.00
	Phi-2 <i>3B</i>	33.24	37.97	32.20	34.95	26.71	28.57	28.57	27.16
	Alpaca-LLaMA-7B	48.97	46.93	36.44	47.93	50.15	46.49	37.59	48.72
Open LLM + MIND	LLaMA2-7B-chat	65.78	64.61	55.75	66.15	59.43	57.13	60.03	59.04
	Mistral-7B-instruct-v0.2	<u>78.57</u>	<u>74.31</u>	<u>80.89</u>	<u>76.97</u>	<u>61.14</u>	<u>65.42</u>	<u>62.16</u>	<u>62.02</u>
LLM API	ChatGPT	75.06	73.76	68.64	74.48	80.74	76.62	68.42	79.23
	ChatGPT (CoT)	76.07	74.53	63.56	75.12	78.89	75.32	78.20	78.21
	ChatGPT (CoT-SC)	76.51	73.82	63.56	75.32	85.72	77.14	82.71	83.99
	GPT-4	78.12	75.41	66.10	76.97	86.03	82.34	84.96	85.30
	GPT-4 (CoT)	78.12	75.41	66.10	76.97	86.03	82.34	84.96	85.30
	GPT-4 (CoT-SC)	78.80	72.88	65.25	76.97	84.00	80.78	84.96	83.48
Human	-	89.96	90.00	80.96	89.33	95.50	85.19	100.0	94.00

Table 3: Evaluation results (Accuracy%) of various language models on both tasks of the INTENTIONQA benchmark. The best performances within each category are underlined and the best among all baselines are **bold-faced**.

intentions with different commonsense relations grounded in ConceptNet (Speer et al., 2017). Specifically, we construct a sibling QA set using our proposed pipeline, with the only additional constraint being that the distractor options share the same relation type as the ground truth option. From the results presented in Figure 3, all the evaluated language models are more effective in understanding the product definition, with an average of 70.47% across relations *isA*, *definedAs*, and *relatedTo*. However, a performance decline of 6.69% is observed in relations that require a deeper understanding of the cause and effect behind the purchasing event, such as *capableOf* and *cause*.

5.4 Error Analysis

In this section, we randomly sample 120 questions that GPT-4 answers incorrectly from INTEN-

TIONQA and categorize the errors by asking experts to annotate them manually.

Among 60 annotated error samples from INTENTUNDERSTAND, we found:

- 40.0% errors are caused by failing to identify the most typical intention, e.g., choosing “because the product is of good quality” instead of “because the person wants to build a water cooling system.”
- 13.3% errors are due to overarching inference. The selected options, while seemingly plausible, cannot be deduced from the products provided.
- 8.3% errors are due to selecting implausible options. The model selects an option that is irrelevant to the given products or implausible. Cases where the rationales in CoT responses are irrelevant to selected options are also observed.

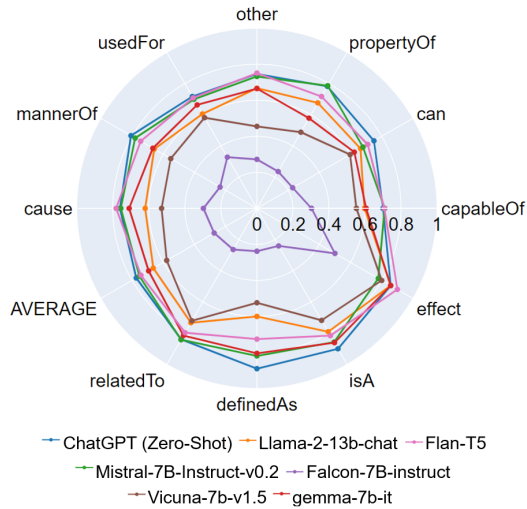


Figure 3: Performances of various language models in comprehending intentions with different relations.

- 10.0% errors are due to inaccurate understanding of the given products.
- 20.0% errors are due to false-negative distractors or incorrect ground truth answers.

Among 60 annotated error samples from INTENTUTILIZE, we found:

- 40% errors are due to inaccurate understanding of the given intention. For example, the model chooses “iPod” when the given intention is “because the customer wanted to use them *with* his/her iPod”.
- 38.3% errors are due to inaccurate understanding of the given products. The reasoning in their response demonstrates inaccurate understanding of the purchased products or those in the options. Or, when the intention is not typical enough to filter out distractors, they fail to rely more on the purchased product to select the best option.
- 21.7% errors are due to false-negative distractors or incorrect ground truth answers.

5.5 Transferring from Other Sources

In this section, we explore whether transferring intention knowledge from other sources can further aid the model’s performance via fine-tuning. Specifically, we use MIND, a knowledge base constructed by Xu et al. (2024), in addition to FolkScope, as a rich source of purchase intentions. MIND is a multi-modal intention knowledge base distilled from LLaVa (Liu et al., 2023a), which includes product images in the knowledge generation process. To ensure the quality of generated intentions, a human-centric intention filtering module

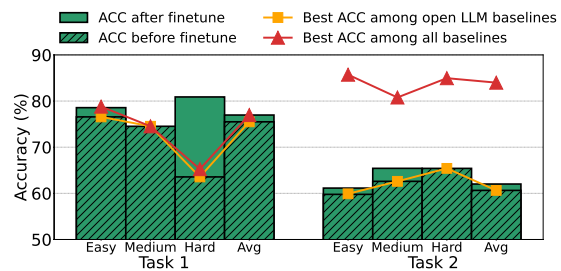


Figure 4: Comparisons between models fine-tuned on intentions from MIND and baseline models achieving top performances.

is developed to eliminate implausible and atypical intentions.

To incorporate MIND’s intention knowledge, we convert 4,059 sets of co-buy records and their corresponding intentions into an instruction-tuning format. We then fine-tune the LLaMA2-7B-chat and Mistral-7B-instruct-v0.2 models on this data using LoRA (Hu et al., 2022). The results, reported in OPEN LLM + MIND of Table 3, reveal an average performance gain of 1.51% and 1.19% for two tasks respectively.

Next, we compare the performance of the fine-tuned Mistral-7B-instruct-v0.2 model with the highest accuracy achieved by all OPEN LLMs and all baselines. The trends are shown in Figure 4. Notably, fine-tuning enables Mistral-7B-instruct-v0.2 to achieve performance comparable to that of GPT-4 in INTENTUNDERSTAND. However, INTENTUTILIZE remains a challenging task even after fine-tuning. This disparity suggests that fine-tuning with intention knowledge facilitates the acquisition of intention understanding abilities, while improving INTENTUTILIZE performance requires more effort. One possible reason is that INTENTUTILIZE places a higher demand on product understanding and reasoning abilities compared to INTENTUNDERSTAND, which cannot be easily improved by simple knowledge injection.

6 Conclusions

In conclusion, this paper presents INTENTIONQA, a double-task MCQA dataset designed to assess the intention comprehension capabilities of LLMs. Extensive experiments and analyses demonstrate that LLMs face significant challenges in certain scenarios, trailing far behind human performance levels, while fine-tuning on external resources brings considerable performance gains. We hope our work sheds light on the limitations of current LLMs in E-commerce intention understanding and facilitates the utilization of LLM in E-commerce scenarios.

606 Limitations

607 We base the negative distractor sampling on sim-
608 ilarity filtering with manually selected thresholds.
609 While these thresholds are decided after multiple
610 rounds of parameter searches and observation of
611 the resulting data quality and have been validated
612 by the human annotation we conduct, automated
613 threshold tuning methods (Xu et al., 2021) could
614 be implemented to facilitate this process.

615 As we build the dataset based on FolkScope,
616 the quality of the latter is upper-bounded by the
617 former. Nevertheless, the construction pipeline in-
618 troduced in this work can be generalized to expand
619 the dataset by incorporating other intention knowl-
620 edge bases. Meanwhile, more advanced LLMs
621 have the potential of curating intention knowledge
622 bases with high quality, further boosting the quality
623 of our QA benchmark.

624 Since LLMs demonstrate strong generative ca-
625 pabilities and commonsense reasoning, it is po-
626 tentially feasible to leverage models such as Chat-
627 GPT to generate contextual information for pur-
628 chase events. However, we rely on the eventuality
629 knowledge graph, ASER, to facilitate the calcula-
630 tion of context embeddings. This offers advantages
631 in terms of cost control and the potential to scale
632 up. Additionally, the human annotation results of
633 our dataset confirm the effectiveness of leveraging
634 ASER for this purpose.

635 Ethics Statement

636 While we adopt LLMs in a generative setting, gen-
637 erating harmful or biased content from them is
638 limited as INTENTIONQA is evaluated in multiple-
639 choice question form. In most cases, the language
640 models generate a single letter representing the
641 option. In CoT, the LLMs generate a short ratio-
642 nale and then output the final answer, where the
643 rationale is closely related to the question itself.
644 All the experiments are conducted using models
645 publicly available via open sources or APIs. The
646 annotators are paid a wage higher than our local
647 law, and the expert annotators are graduate students
648 specializing in natural language processing. They
649 have all agreed to participate voluntarily and are
650 well-instructed about the tasks.

651 References

652 Ebtesam Almazrouei, Hamza Alobeidli, Abdulaziz Al-
653 shamsi, Alessandro Cappelli, Ruxandra Cojocaru,

Mérouane Debbah, Étienne Goffinet, Daniel Hesslow, 654
Julien Launay, Quentin Malartic, Daniele Mazzotta, 655
Badreddine Noune, Baptiste Pannier, and Guilherme 656
Penedo. 2023. [The falcon series of open language 657](#)
[models](#). *CoRR*, abs/2311.16867. 658

Chunkit Chan, Jiayang Cheng, Weiqi Wang, Yuxin 659
Jiang, Tianqing Fang, Xin Liu, and Yangqiu Song. 660
2024. [Exploring the potential of chatgpt on sentence 661](#)
[level relations: A focus on temporal, causal, and 662](#)
[discourse relations](#). In *Findings of the Association 663*
for Computational Linguistics: EACL 2024, St. Ju- 664
lian’s, Malta, March 17-22, 2024, pages 684–721. 665
Association for Computational Linguistics. 666

Hyung Won Chung, Le Hou, Shayne Longpre, Barret 667
Zoph, Yi Tay, William Fedus, Eric Li, Xuezhi Wang, 668
Mostafa Dehghani, Siddhartha Brahma, Albert Web- 669
son, Shixiang Shane Gu, Zhuyun Dai, Mirac Suz- 670
gun, Xinyun Chen, Aakanksha Chowdhery, Sharan 671
Narang, Gaurav Mishra, Adams Yu, Vincent Y. Zhao, 672
Yanping Huang, Andrew M. Dai, Hongkun Yu, Slav 673
Petrov, Ed H. Chi, Jeff Dean, Jacob Devlin, Adam 674
Roberts, Denny Zhou, Quoc V. Le, and Jason Wei. 675
2022. [Scaling instruction-finetuned language models](#). 676
CoRR, abs/2210.11416. 677

Honghua (Kathy) Dai, Lingzhi Zhao, Zaiqing Nie, Ji- 678
Rong Wen, Lee Wang, and Ying Li. 2006. [Detecting 679](#)
[online commercial intention \(OCI\)](#). In *Proceedings 680*
of the 15th international conference on World Wide 681
Web, WWW 2006, Edinburgh, Scotland, UK, May 682
23-26, 2006, pages 829–837. ACM. 683

Yang Deng, Wenxuan Zhang, Qian Yu, and Wai Lam. 684
2023. [Product question answering in e-commerce: A 685](#)
[survey](#). In *Proceedings of the 61st Annual Meeting of 686*
the Association for Computational Linguistics (Vol- 687
ume 1: Long Papers), ACL 2023, Toronto, Canada, 688
July 9-14, 2023, pages 11951–11964. Association for 689
Computational Linguistics. 690

Leonardo Fogassi, Pier Francesco Ferrari, Benno 691
Gesierich, Stefano Rozzi, Fabian Chersi, and Gia- 692
como Rizzolatti. 2005. Parietal lobe: from action 693
organization to intention understanding. *Science*, 694
308(5722):662–667. 695

Abhinav Goyal, Anupam Singh, and Nikesh Garera. 696
2022. [End-to-end speech to intent prediction to 697](#)
[improve e-commerce customer support voicebot in 698](#)
[hindi and english](#). In *Proceedings of the 2022 Confer- 699*
ence on Empirical Methods in Natural Language Pro- 700
cessing: EMNLP 2022 - Industry Track, Abu Dhabi, 701
UAE, December 7 - 11, 2022, pages 579–586. Asso- 702
ciation for Computational Linguistics. 703

Mihajlo Grbovic, Vladan Radosavljevic, Nemanja 704
Djuric, Narayan Bhamidipati, Jaikit Savla, Varun 705
Bhagwan, and Doug Sharp. 2015. [E-commerce in 706](#)
[your inbox: Product recommendations at scale](#). In 707
Proceedings of the 21th ACM SIGKDD International 708
Conference on Knowledge Discovery and Data Min- 709
ing, Sydney, NSW, Australia, August 10-13, 2015, 710
pages 1809–1818. ACM. 711

712	Jing Gu, Qingyang Wu, and Zhou Yu. 2021. Perception score: A learned metric for open-ended text generation evaluation . In <i>Thirty-Fifth AAAI Conference on Artificial Intelligence, AAAI 2021, Thirty-Third Conference on Innovative Applications of Artificial Intelligence, IAAI 2021, The Eleventh Symposium on Educational Advances in Artificial Intelligence, EAAI 2021, Virtual Event, February 2-9, 2021</i> , pages 12902–12910. AAAI Press.	
713		
714		
715		
716		
717		
718		
719		
720		
721	Suriya Gunasekar, Yi Zhang, Jyoti Aneja, Caio César Teodoro Mendes, Allie Del Giorno, Sivakanth Gopi, Mojan Javaheripi, Piero Kauffmann, Gustavo de Rosa, Olli Saarikivi, Adil Salim, Shital Shah, Harkirat Singh Behl, Xin Wang, Sébastien Bubeck, Ronen Eldan, Adam Tauman Kalai, Yin Tat Lee, and Yuanzhi Li. 2023. Textbooks are all you need . <i>CoRR</i> , abs/2306.11644.	
722		
723		
724		
725		
726		
727		
728		
729	Pengcheng He, Jianfeng Gao, and Weizhu Chen. 2023. Debertav3: Improving deberta using electra-style pre-training with gradient-disentangled embedding sharing . In <i>The Eleventh International Conference on Learning Representations, ICLR 2023, Kigali, Rwanda, May 1-5, 2023</i> . OpenReview.net.	
730		
731		
732		
733		
734		
735	Sharon Hirsch, Ido Guy, Alexander Nus, Arnon Dagan, and Oren Kurland. 2020. Query reformulation in e-commerce search . In <i>Proceedings of the 43rd International ACM SIGIR conference on research and development in Information Retrieval, SIGIR 2020, Virtual Event, China, July 25-30, 2020</i> , pages 1319–1328. ACM.	
736		
737		
738		
739		
740		
741		
742	Edward J. Hu, Yelong Shen, Phillip Wallis, Zeyuan Allen-Zhu, Yuanzhi Li, Shean Wang, Lu Wang, and Weizhu Chen. 2022. Lora: Low-rank adaptation of large language models . In <i>The Tenth International Conference on Learning Representations, ICLR 2022, Virtual Event, April 25-29, 2022</i> . OpenReview.net.	
743		
744		
745		
746		
747		
748	Raghav Jain, Daivik Sojitra, Arkadeep Acharya, Sriparna Saha, Adam Jatowt, and Sandipan Dandapat. 2023. Do language models have a common sense regarding time? revisiting temporal commonsense reasoning in the era of large language models . In <i>Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing, EMNLP 2023, Singapore, December 6-10, 2023</i> , pages 6750–6774. Association for Computational Linguistics.	
749		
750		
751		
752		
753		
754		
755		
756		
757	Albert Q. Jiang, Alexandre Sablayrolles, Arthur Mensch, Chris Bamford, Devendra Singh Chaplot, Diego de Las Casas, Florian Bressand, Gianna Lengyel, Guillaume Lample, Lucile Saulnier, Léo Renard Lavaud, Marie-Anne Lachaux, Pierre Stock, Teven Le Scao, Thibaut Lavril, Thomas Wang, Timothée Lacroix, and William El Sayed. 2023. Mistral 7b . <i>CoRR</i> , abs/2310.06825.	
758		
759		
760		
761		
762		
763		
764		
765	Wei Jin, Haitao Mao, Zheng Li, Haoming Jiang, Chen Luo, Hongzhi Wen, Haoyu Han, Hanqing Lu, Zhengyang Wang, Ruirui Li, Zhen Li, Monica Cheng, Rahul Goutam, Haiyang Zhang, Karthik Subbian, Suhang Wang, Yizhou Sun, Jiliang Tang, Bing Yin,	
766		
767		
768		
769		
	and Xianfeng Tang. 2023. Amazon-m2: A multilingual multi-locale shopping session dataset for recommendation and text generation . In <i>Advances in Neural Information Processing Systems 36: Annual Conference on Neural Information Processing Systems 2023, NeurIPS 2023, New Orleans, LA, USA, December 10 - 16, 2023</i> .	770
		771
		772
		773
		774
		775
		776
	Yeongseo Jung, Eunseo Jung, and Lei Chen. 2023. Towards a unified conversational recommendation system: Multi-task learning via contextualized knowledge distillation . In <i>Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing, EMNLP 2023, Singapore, December 6-10, 2023</i> , pages 13625–13637. Association for Computational Linguistics.	777
		778
		779
		780
		781
		782
		783
		784
	Dong-Mo Koo and Seon-Hee Ju. 2010. The interactional effects of atmospherics and perceptual curiosity on emotions and online shopping intention. <i>Computers in Human Behavior</i> , 26(3):377–388.	785
		786
		787
		788
	Lei Li, Yongfeng Zhang, and Li Chen. 2020. Generate neural template explanations for recommendation . In <i>CIKM '20: The 29th ACM International Conference on Information and Knowledge Management, Virtual Event, Ireland, October 19-23, 2020</i> , pages 755–764. ACM.	789
		790
		791
		792
		793
		794
	Yangning Li, Shirong Ma, Xiaobin Wang, Shen Huang, Chengyue Jiang, Haitao Zheng, Pengjun Xie, Fei Huang, and Yong Jiang. 2024. Ecomgpt: Instruction-tuning large language models with chain-of-task tasks for e-commerce . In <i>Thirty-Eighth AAAI Conference on Artificial Intelligence, AAAI 2024, Thirty-Sixth Conference on Innovative Applications of Artificial Intelligence, IAAI 2024, Fourteenth Symposium on Educational Advances in Artificial Intelligence, EAAI 2014, February 20-27, 2024, Vancouver, Canada</i> , pages 18582–18590. AAAI Press.	795
		796
		797
		798
		799
		800
		801
		802
		803
		804
		805
	Haotian Liu, Chunyuan Li, Qingyang Wu, and Yong Jae Lee. 2023a. Visual instruction tuning . In <i>Advances in Neural Information Processing Systems 36: Annual Conference on Neural Information Processing Systems 2023, NeurIPS 2023, New Orleans, LA, USA, December 10 - 16, 2023</i> .	806
		807
		808
		809
		810
		811
	Jiacheng Liu, Wenya Wang, Dianzhuo Wang, Noah A. Smith, Yejin Choi, and Hannaneh Hajishirzi. 2023b. Vera: A general-purpose plausibility estimation model for commonsense statements . In <i>Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing, EMNLP 2023, Singapore, December 6-10, 2023</i> , pages 1264–1287. Association for Computational Linguistics.	812
		813
		814
		815
		816
		817
		818
		819
	Yinhan Liu, Myle Ott, Naman Goyal, Jingfei Du, Mandar Joshi, Danqi Chen, Omer Levy, Mike Lewis, Luke Zettlemoyer, and Veselin Stoyanov. 2019. Roberta: A robustly optimized BERT pretraining approach . <i>CoRR</i> , abs/1907.11692.	820
		821
		822
		823
		824
	Kaixin Ma, Filip Ilievski, Jonathan Francis, Yonatan Bisk, Eric Nyberg, and Alessandro Oltramari. 2021.	825
		826

827	Knowledge-driven data construction for zero-shot evaluation in commonsense question answering. In <i>Thirty-Fifth AAAI Conference on Artificial Intelligence, AAAI 2021, Thirty-Third Conference on Innovative Applications of Artificial Intelligence, IAAI 2021, The Eleventh Symposium on Educational Advances in Artificial Intelligence, EAAI 2021, Virtual Event, February 2-9, 2021</i> , pages 13507–13515. AAAI Press.	
828		
829		
830		
831		
832		
833		
834		
835		
836	Thomas Mesnard, Cassidy Hardin, Robert Dadashi, Surya Bhupatiraju, Shreya Pathak, Laurent Sifre, Morgane Rivière, Mihir Sanjay Kale, Juliette Love, Pouya Tafti, Léonard Hussenot, Aakanksha Chowdhery, Adam Roberts, Aditya Barua, Alex Botev, Alex Castro-Ros, Ambrose Slone, Amélie Héliou, Andrea Tacchetti, Anna Bulanova, Antonia Paterson, Beth Tsai, Bobak Shahriari, Charline Le Lan, Christopher A. Choquette-Choo, Clément Crepy, Daniel Cer, Daphne Ippolito, David Reid, Elena Buchatskaya, Eric Ni, Eric Noland, Geng Yan, George Tucker, George-Christian Muraru, Grigory Rozhdestvenskiy, Henryk Michalewski, Ian Tenney, Ivan Grishchenko, Jacob Austin, James Keeling, Jane Labanowski, Jean-Baptiste Lespiau, Jeff Stanway, Jenny Brennan, Jeremy Chen, Johan Ferret, Justin Chiu, and et al. 2024. Gemma: Open models based on gemini research and technology . <i>CoRR</i> , abs/2403.08295.	
837		
838		
839		
840		
841		
842		
843		
844		
845		
846		
847		
848		
849		
850		
851		
852		
853		
854	OpenAI. 2022. Chatgpt: Optimizing language models for dialogue . <i>OpenAI</i> .	
855		
856	OpenAI. 2023. GPT-4 technical report . <i>CoRR</i> , abs/2303.08774.	
857		
858	Long Ouyang, Jeffrey 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 F. Christiano, Jan Leike, and Ryan Lowe. 2022. Training language models to follow instructions with human feedback .	
859		
860		
861		
862		
863		
864		
865		
866	Mingjie Qian, Yongsen Zheng, Jinghui Qin, and Liang Lin. 2023. Hutcrs: Hierarchical user-interest tracking for conversational recommender system . In <i>Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing, EMNLP 2023, Singapore, December 6-10, 2023</i> , pages 10281–10290. Association for Computational Linguistics.	
867		
868		
869		
870		
871		
872		
873	Chengwei Qin, Aston Zhang, Zhuosheng Zhang, Jiaao Chen, Michihiro Yasunaga, and Diyi Yang. 2023. Is chatgpt a general-purpose natural language processing task solver? In <i>Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing, EMNLP 2023, Singapore, December 6-10, 2023</i> , pages 1339–1384. Association for Computational Linguistics.	
874		
875		
876		
877		
878		
879		
880		
881	Colin Raffel, Noam Shazeer, Adam Roberts, Katherine Lee, Sharan Narang, Michael Matena, Yanqi Zhou, Wei Li, and Peter J. Liu. 2020. Exploring the limits of transfer learning with a unified text-to-text transformer . <i>J. Mach. Learn. Res.</i> , 21:140:1–140:67.	
882		
883		
884		
885		
	Nils Reimers and Iryna Gurevych. 2019. Sentence-bert: Sentence embeddings using siamese bert-networks . In <i>Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing, EMNLP-IJCNLP 2019, Hong Kong, China, November 3-7, 2019</i> , pages 3980–3990. Association for Computational Linguistics.	886 887 888 889 890 891 892 893
	Victor Sanh, Albert Webson, Colin Raffel, Stephen H. Bach, Lintang Sutawika, Zaid Alyafeai, Antoine Chaffin, Arnaud Stiegler, Arun Raja, Manan Dey, M Saiful Bari, Canwen Xu, Urmish Thakker, Shanya Sharma Sharma, Eliza Szczechla, Taewoon Kim, Gunjan Chhablani, Nihal V. Nayak, Debajyoti Datta, Jonathan Chang, Mike Tian-Jian Jiang, Han Wang, Matteo Manica, Sheng Shen, Zheng Xin Yong, Harshit Pandey, Rachel Bawden, Thomas Wang, Trishala Neeraj, Jos Rozen, Abheesht Sharma, Andrea Santilli, Thibault Févry, Jason Alan Fries, Ryan Teehan, Teven Le Scao, Stella Biderman, Leo Gao, Thomas Wolf, and Alexander M. Rush. 2022. Multi-task prompted training enables zero-shot task generalization . In <i>The Tenth International Conference on Learning Representations, ICLR 2022, Virtual Event, April 25-29, 2022</i> . OpenReview.net.	894 895 896 897 898 899 900 901 902 903 904 905 906 907 908 909 910
	Robyn Speer, Joshua Chin, and Catherine Havasi. 2017. Conceptnet 5.5: An open multilingual graph of general knowledge . In <i>Proceedings of the Thirty-First AAAI Conference on Artificial Intelligence, February 4-9, 2017, San Francisco, California, USA</i> , pages 4444–4451. AAAI Press.	911 912 913 914 915 916
	Qingyu Tan, Hwee Tou Ng, and Lidong Bing. 2023. Towards benchmarking and improving the temporal reasoning capability of large language models . In <i>Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers), ACL 2023, Toronto, Canada, July 9-14, 2023</i> , pages 14820–14835. Association for Computational Linguistics.	917 918 919 920 921 922 923 924
	Rohan Taori, Ishaan Gulrajani, Tianyi Zhang, Yann Dubois, Xuechen Li, Carlos Guestrin, Percy Liang, and Tatsunori B. Hashimoto. 2023. Stanford alpaca: An instruction-following llama model . https://github.com/tatsu-lab/stanford_alpaca .	925 926 927 928 929
	Hugo Touvron, Louis Martin, Kevin Stone, Peter Albert, Amjad Almahairi, Yasmine Babaei, Nikolay Bashlykov, Soumya Batra, Prajjwal Bhargava, Shrutu Bhosale, Dan Bikel, Lukas Blecher, Cristian Canton-Ferrer, Moya Chen, Guillem Cucurull, David Esiobu, Jude Fernandes, Jeremy Fu, Wenyin Fu, Brian Fuller, Cynthia Gao, Vedanuj Goswami, Naman Goyal, Anthony Hartshorn, Saghar Hosseini, Rui Hou, Hakan Inan, Marcin Kardas, Viktor Kerkez, Madian Khabsa, Isabel Kloumann, Artem Korenev, Punit Singh Koura, Marie-Anne Lachaux, Thibaut Lavril, Jenya Lee, Diana Liskovich, Yinghai Lu, Yuning Mao, Xavier Martinet, Todor Mihaylov, Pushkar Mishra, Igor Molybog, Yixin Nie, Andrew Poulton, Jeremy Reizenstein, Rashi Rungta, Kalyan Saladi, Alan Schelten,	930 931 932 933 934 935 936 937 938 939 940 941 942 943 944

945	Ruan Silva, Eric Michael Smith, Ranjan Subramanian, Xiaoqing Ellen Tan, Binh Tang, Ross Taylor, Adina Williams, Jian Xiang Kuan, Puxin Xu, Zheng Yan, Iliyan Zarov, Yuchen Zhang, Angela Fan, Melanie Kambadur, Sharan Narang, Aurélien Rodriguez, Robert Stojnic, Sergey Edunov, and Thomas Scialom. 2023. Llama 2: Open foundation and fine-tuned chat models . <i>CoRR</i> , abs/2307.09288.	
946		
947		
948		
949		
950		
951		
952		
953	Wei qi Wang, Tianqing Fang, Wenxuan Ding, Baixuan Xu, Xin Liu, Yangqiu Song, and Antoine Bosselut. 2023a. CAR: conceptualization-augmented reasoner for zero-shot commonsense question answering . In <i>Findings of the Association for Computational Linguistics: EMNLP 2023, Singapore, December 6-10, 2023</i> , pages 13520–13545. Association for Computational Linguistics.	
954		
955		
956		
957		
958		
959		
960		
961	Wei qi Wang, Tianqing Fang, Chunyang Li, Haochen Shi, Wenxuan Ding, Baixuan Xu, Zhaowei Wang, Jiaxin Bai, Xin Liu, Jiayang Cheng, Chunkit Chan, and Yangqiu Song. 2024. CANDLE: iterative conceptualization and instantiation distillation from large language models for commonsense reasoning . <i>CoRR</i> , abs/2401.07286.	
962		
963		
964		
965		
966		
967		
968	Wei qi Wang, Tianqing Fang, Baixuan Xu, Chun Yi Louis Bo, Yangqiu Song, and Lei Chen. 2023b. CAT: A contextualized conceptualization and instantiation framework for commonsense reasoning . In <i>Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers), ACL 2023, Toronto, Canada, July 9-14, 2023</i> , pages 13111–13140. Association for Computational Linguistics.	
969		
970		
971		
972		
973		
974		
975		
976		
977	Xuezhi Wang, Jason Wei, Dale Schuurmans, Quoc V. Le, Ed H. Chi, Sharan Narang, Aakanksha Chowdhery, and Denny Zhou. 2023c. Self-consistency improves chain of thought reasoning in language models . In <i>The Eleventh International Conference on Learning Representations, ICLR 2023, Kigali, Rwanda, May 1-5, 2023</i> . OpenReview.net.	
978		
979		
980		
981		
982		
983		
984	Yizhong Wang, Yeganeh Kordi, Swaroop Mishra, Alisa Liu, Noah A. Smith, Daniel Khashabi, and Hannaneh Hajishirzi. 2023d. Self-instruct: Aligning language models with self-generated instructions . In <i>Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers), ACL 2023, Toronto, Canada, July 9-14, 2023</i> , pages 13484–13508. Association for Computational Linguistics.	
985		
986		
987		
988		
989		
990		
991		
992		
993	Jason Wei, Xuezhi Wang, Dale Schuurmans, Maarten Bosma, Brian Ichter, Fei Xia, Ed H. Chi, Quoc V. Le, and Denny Zhou. 2022. Chain-of-thought prompting elicits reasoning in large language models . In <i>Advances in Neural Information Processing Systems 35: Annual Conference on Neural Information Processing Systems 2022, NeurIPS 2022, New Orleans, LA, USA, November 28 - December 9, 2022</i> .	
994		
995		
996		
997		
998		
999		
1000		
1001	Baixuan Xu, Wei qi Wang, Haochen Shi, Wenxuan Ding, Huihao Jing, Tianqing Fang, Jiaxin Bai, Xin Liu,	
1002		
	Changlong Yu, Zheng Li, Chen Luo, Bing Yin, Long Chen, and Yangqiu Song. 2024. Mind: Multi-modal shopping intention distillation from large vision-language models for e-commerce purchase understanding .	1003 1004 1005 1006 1007
	Yi Xu, Lei Shang, Jinxing Ye, Qi Qian, Yu-Feng Li, Baigui Sun, Hao Li, and Rong Jin. 2021. Dash: Semi-supervised learning with dynamic thresholding . In <i>Proceedings of the 38th International Conference on Machine Learning, ICML 2021, 18-24 July 2021, Virtual Event</i> , volume 139 of <i>Proceedings of Machine Learning Research</i> , pages 11525–11536. PMLR.	1008 1009 1010 1011 1012 1013 1014
	Yang Yang and Jie Tang. 2015. Beyond query: Interactive user intention understanding . In <i>2015 IEEE International Conference on Data Mining, ICDM 2015, Atlantic City, NJ, USA, November 14-17, 2015</i> , pages 519–528. IEEE Computer Society.	1015 1016 1017 1018 1019
	Changlong Yu, Xin Liu, Jefferson Maia, Tianyu Cao, Laurence Yang Li, Yifan Gao, Yangqiu Song, Rahul Goutam, Haiyang Zhang, Bing Yin, et al. 2024. Cosmo: A large-scale e-commerce common sense knowledge generation and serving system at amazon . In <i>Proceedings of the 2024 International Conference on Management of Data, SIGMOD 2024</i> .	1020 1021 1022 1023 1024 1025 1026
	Changlong Yu, Wei qi Wang, Xin Liu, Jiaxin Bai, Yangqiu Song, Zheng Li, Yifan Gao, Tianyu Cao, and Bing Yin. 2023. Folkscope: Intention knowledge graph construction for e-commerce commonsense discovery . In <i>Findings of the Association for Computational Linguistics: ACL 2023, Toronto, Canada, July 9-14, 2023</i> , pages 1173–1191. Association for Computational Linguistics.	1027 1028 1029 1030 1031 1032 1033 1034
	Qian Yu and Wai Lam. 2018. Product question intent detection using indicative clause attention and adversarial learning . In <i>Proceedings of the 2018 ACM SIGIR International Conference on Theory of Information Retrieval, ICTIR 2018, Tianjin, China, September 14-17, 2018</i> , pages 75–82. ACM.	1035 1036 1037 1038 1039 1040
	Hongfei Zhang, Xia Song, Chenyan Xiong, Corby Rosset, Paul N. Bennett, Nick Craswell, and Saurabh Tiwary. 2019. Generic intent representation in web search . In <i>Proceedings of the 42nd International ACM SIGIR Conference on Research and Development in Information Retrieval, SIGIR 2019, Paris, France, July 21-25, 2019</i> , pages 65–74. ACM.	1041 1042 1043 1044 1045 1046 1047
	Hongming Zhang, Xin Liu, Haojie Pan, Haowen Ke, Jiefu Ou, Tianqing Fang, and Yangqiu Song. 2022a. ASER: towards large-scale commonsense knowledge acquisition via higher-order selectional preference over eventualities . <i>Artif. Intell.</i> , 309:103740.	1048 1049 1050 1051 1052
	Hongming Zhang, Xin Liu, Haojie Pan, Yangqiu Song, and Cane Wing-Ki Leung. 2020. ASER: A large-scale eventuality knowledge graph . In <i>WWW '20: The Web Conference 2020, Taipei, Taiwan, April 20-24, 2020</i> , pages 201–211. ACM / IW3C2.	1053 1054 1055 1056 1057

1058	Susan Zhang, Stephen Roller, Naman Goyal, Mikel Artetxe, Moya Chen, Shuohui Chen, Christopher Dewan, Mona T. 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. 2022b. OPT: open pre-trained transformer language models . <i>CoRR</i> , abs/2205.01068.
1066	Jiashu Zhao, Hongshen Chen, and Dawei Yin. 2019. A dynamic product-aware learning model for e-commerce query intent understanding . In <i>Proceedings of the 28th ACM International Conference on Information and Knowledge Management, CIKM 2019, Beijing, China, November 3-7, 2019</i> , pages 1843–1852. ACM.
1073	Wayne Xin Zhao, Yanwei Guo, Yulan He, Han Jiang, Yuxin Wu, and Xiaoming Li. 2014. We know what you want to buy: a demographic-based system for product recommendation on microblogs . In <i>The 20th ACM SIGKDD International Conference on Knowledge Discovery and Data Mining, KDD '14, New York, NY, USA - August 24 - 27, 2014</i> , pages 1935–1944. ACM.
1081	Lianmin Zheng, Wei-Lin Chiang, Ying Sheng, Siyuan Zhuang, Zhanghao Wu, Yonghao Zhuang, Zi Lin, Zhuohan Li, Dacheng Li, Eric P. Xing, Hao Zhang, Joseph E. Gonzalez, and Ion Stoica. 2023. Judging llm-as-a-judge with mt-bench and chatbot arena . In <i>Advances in Neural Information Processing Systems 36: Annual Conference on Neural Information Processing Systems 2023, NeurIPS 2023, New Orleans, LA, USA, December 10 - 16, 2023</i> .
1090	Yaowei Zheng, Richong Zhang, Junhao Zhang, Yanhan Ye, Zheyang Luo, and Yongqiang Ma. 2024. Llamafactory: Unified efficient fine-tuning of 100+ language models . <i>CoRR</i> , abs/2403.13372.
1094	Wendi Zhou, Tianyi Li, Pavlos Vougiouklis, Mark Steedman, and Jeff Z. Pan. 2024. A usage-centric take on intent understanding in e-commerce . <i>CoRR</i> , abs/2402.14901.

Appendices

A Implementation Details

A.1 Hyperparameter Settings

For models in the category of LLM API, we set the `max_tokens` to 10 and 200 respectively for ZERO-SHOT and CoT (CoT+SC) prompting. The temperature τ is set to 0.1 for non-Self-Consistency decoding and 0.7 otherwise.

For models in the category of OPEN LLM, we use the default setting as presented in Hugging Face.

For fine-tuning LLMs, we use the open-sourced library LLaMA-Factory¹ (Zheng et al., 2024) to train all models. All hyperparameters follow the default settings, and a LoRA rank of $\alpha = 64$ is used. We conduct all experiments on a Linux machine with eight NVIDIA V100 GPUs.

B Prompts

We report the prompt used for product simplification with ChatGPT in Tabel 4.

To transfer a conceptualized product into ASER nodes’ format, we utilize natural language templates. These templates consist of assembling each product category with a subject and a verb that carry semantic meanings related to purchasing. Specifically, the subjects we use include: PersonX, PersonY, PeopleX, PeopleY; Similarly, the verbs we employ are: buy, shop, purchase, get, obtain, have, in simple present tense, original form, simple perfect tense, or past tense, with optional articles (a, an, the, 1, 2) added before the conceptualized product name. As a result, when a product such as “iPhone 14” occurs, we transform it into a list of concise yet semantically complete events that can potentially be matched in ASER. For example, one of the transformed events could be “PersonX bought a phone.”

We report the prompts used for INTENTUNDERSTAND and INTENTUTILIZE in Table 5 and Table 6 respectively.

C Case Study

We present example questions that GPT-4 successfully answer or fail with CoT for both tasks in Table 7.

¹<https://github.com/hiyouga/LLaMA-Factory>

Method	Prompt for Product Name Simplification
ZERO-SHOT	Product name: <product>; What is the category of the product? Generate three possible categories, each in 2 words, separated by a comma.

Table 4: Prompt used to instruct ChatGPT to conceptualize the product name.

Method	Prompt for INTENTUNDERSTAND
ZERO-SHOT	A customer buys <product 1> and <product 2>. What is the most likely intention for buying them? A. because <intention 1> B. because <intention 2> C. because <intention 3> D. because <intention 4> Answer A or B or C or D only without any other word.
CoT	A customer buys <product 1> and <product 2>. What is the most likely intention for buying them? A. because <intention 1> B. because <intention 2> C. because <intention 3> D. because <intention 4> Formulate your answer in this way. Step 1: Give a short and brief rationale by thinking step by step. Step 2: Answer A or B or C or D only without any other word.

Table 5: Prompts for INTENTUNDERSTAND with ZERO-SHOT prompting and CoT respectively.

Method	Prompt for INTENTUTILIZE
ZERO-SHOT	A customer buys <product>, because <intention>. What is the customer's most probable additional purchase? A. <product 1> B. <product 2> C. <product 3> D. <product 4> Answer A or B or C or D only without any other word.
CoT	A customer buys <product>, because <intention>. What is the customer's most probable additional purchase? A. <product 1> B. <product 2> C. <product 3> D. <product 4> Formulate your answer in this way. Step 1: Give a short and brief rationale by thinking step by step. Step 2: Answer A or B or C or D only without any other word.

Table 6: Prompts for INTENTUTILIZE with ZERO-SHOT prompting and CoT respectively.

Task	Prompt & Response
Task 1	<p>[Prompt] A customer buys Belkin USB A/A Extension Cable, USB Type-A Female and USB Type-A Male (6 Feet) and Cable Matters 3-Pack USB Extension Cable (USB Extender Cable/USB 2.0 Extension Cable) - 3 Feet - Available 3FT - 10FT in Length. What is the most likely intention for buying them? A. because they both are a type of USB extension cable. B. because they both are a manner of "Audio & Video Accessories". C. because they both are related to the topic of "Army Style". D. because they could both be worn as a bracelet. Formulate your answer in this way. Step 1: Give a short and brief rationale by thinking step by step. Step 2: Answer A or B or C or D only without any other word. [Response from GPT-4] Step 1: The customer is buying two different USB extension cables, which indicates they need to extend the length of their USB connections. Step 2: A [Correct!]</p>
Task 1	<p>[Prompt] A customer buys Barbra Collection Hawaiian Style Kukui Nut Lei Hand Painted Multi-color Turtle 32 Inches and Barbra Collection Hawaiian Style Kukui Nut Lei Hibiscus Hand Painted Turquoise Flower 32 Inches. What is the most likely intention for buying them? A. because the person wants to replace the LCD screen of the laptop. B. because they both are capable of giving him a lot of benefits. C. because they both are a type of toys. D. because they could both be found on Aliexpress.com. Formulate your answer in this way. Step 1: Give a short and brief rationale by thinking step by step. Step 2: Answer A or B or C or D only without any other word. [Response from GPT-4] Step 1: The customer is buying two different styles of Kukui Nut Lei, which are Hawaiian-style accessories. The intention is most likely related to their use as accessories or decorations. Step 2: D [Correct answer: B]</p>
Task 2	<p>[Prompt] A customer buys StarTech.com 60x10mm Replacement Ball Bearing Computer Case Fan w/ TX3 Connector - 3 pin case Fan - TX3 Fan - 60mm Fan, because they both are defined as "Computer Accessories". What is the customer's most probable additional purchase? A. StarTech 8-Inch 24 Pin ATX 2.01 Power Extension Cable (ATX24POWEXT) B. Genuine Power On-Off On/Off On Off Flex Cable Ribbon FOR Apple iPad 2 iPad2 C. WiFi Bluetooth Antenna Flex Cable Replacement Part for iPad 3 and iPad 4 D. Headphone Audio Jack Flex Cable Ribbon Replacement for iPad Air (Black) Formulate your answer in this way. Step 1: Give a short and brief rationale by thinking step by step. Step 2: Answer A or B or C or D only without any other word. [Response from GPT-4] Step 1: The customer has purchased a computer case fan, which indicates that they are likely working on a desktop computer or building one. Therefore, the most probable additional purchase would be related to desktop computer components or accessories. Step 2: A [Correct!]</p>
Task 2	<p>[Prompt] A customer buys Nakamichi NK2000 Over-the-Head Headphones White, because he wanted to use them with his iPod. What is the customer's most probable additional purchase? A. Nakamichi Amplified Stereo Headphones NK2000 Turbulence B. Apple iPod nano 8 GB Silver (4th Generation) (Discontinued by Manufacturer) C. Iwoo Bluetooth Headphone, Blue D. Portable CD Cassette Boombox with AM/FM Radio Formulate your answer in this way. Step 1: Give a short and brief rationale by thinking step by step. [Response from GPT-4] Step 2: Answer A or B or C or D only without any other word. Step 1: The customer already bought headphones for his iPod, so he is unlikely to buy another pair of headphones or a different audio device. He might need an accessory that is compatible with his iPod. Step 2: B [Correct answer: A]</p>

Table 7: Example prompts and responses from GPT-4 with CoT prompting methods.