

SESSIONINTENTBENCH: A Multi-task Inter-session Intention-shift Modeling Benchmark for E-commerce Customer Behavior Understanding

Anonymous ACL submission

Abstract

Session history is a common way of recording user interacting behaviors throughout a browsing activity with multiple products. For example, if a user clicks a product webpage and then leaves, it might be because there are certain features that don't satisfy the user, which serve as an important indicator of on-the-spot user preferences. However, all prior works fail to capture and model customer intention effectively because insufficient information exploitation and only apparent information like descriptions and titles are used. There is also a lack of data and corresponding benchmark for explicitly modeling intention in E-commerce product purchase sessions. To address these issues, we introduce the concept of an *intention tree* and propose a dataset curation pipeline. Together, we construct a sibling multimodal benchmark, SESSIONINTENTBENCH, that evaluates L(V)LMs' capability on understanding inter-session intention shift with four subtasks. With 1,952,177 intention entries, 1,132,145 session intention trajectories, and 13,003,664 available tasks mined using 10,905 sessions, we provide a scalable way to exploit the existing session data for customer intention understanding. We conduct human annotations to collect ground-truth label for a subset of collected data to form an evaluation gold set. Extensive experiments on the annotated data further confirm that current L(V)LMs fail to capture and utilize the intention across the complex session setting. Further analysis show injecting intention enhances LLMs' performances.

1 Introduction

Modeling and analyzing customer intention is of great importance in the E-commerce domain (Dai et al., 2006; Jammalamadaka et al., 2009; Li et al., 2020). This enables us to give better product recommendations and provide more personalized services (Hu et al., 2008; Zhao et al., 2015; Zhu et al., 2024). Conventional ways of understanding user



Figure 1: An example of customer intention-shift in the session. At each session step, the customer interacts with a new product, may change his purchase intent, and then looks for items with desired features.

intention always rely on analyzing user profiles or purchasing records, but such information is not easily retrievable or even missing in real world applications. Therefore, we need a data source with better accessibility and applicability, such as the product purchase sessions, which concludes the user behavior throughout a series of sequential browsing activities. By analyzing the interaction history in this short period of time, we are able to infer the user intention and how it changes over time. The shifting intent behind product searches and inspections can further affect future user interactions. For example, in Figure 1, the customer exposes his intention when he switches from flashy red shoes to plain white ones. After that, browsing for shoes at a much lower price shows customers' need for cheap and cheerful products. By modeling customer session intention and adjusting inferred results when needed, we can provide more customized services in an accurate and timely manner.

Existing work either covers session or intention, but not collectively. There has been an experiment focusing on exploiting the product information within one session and using it to make direct predictions (Jin et al., 2023b), which assembles useful information based on specific product attributes like titles and prices. While some other works explicitly model the user intention behind the single purchase or co-buy behaviors (Xu et al., 2024; Ding et al., 2024). They leverage the most recent user actions for intention understanding and inference, covering only one or two products, but fall short of exploring user preference shifts over a longer horizon, such as sessions. However, Jin et al. (2023b) have shown that session information and fine-grained attribute analysis would help LLMs to give better next-product recommendations. Considering these aspects, it is essential to formulate a method to explicitly model intention over a session period.

But when modeling intention dynamically in more complex purchase contexts, such as sessions, several gaps remain. Firstly, current works only use short-term information and focus on single or co-buy purchases. This approach overlooks the potential motivational intention embedded in earlier user interactions, therefore hindering the models’ capability of making reasonable inferences. Furthermore, among various attributes, only product titles and images are used as product inference hints, which omits important dimensions of product information and results in a waste of information from the collected knowledge base. Last but not least, we lack an automated pipeline to streamline the construction of such intention data, there hasn’t been any formulation of such tasks or benchmark data to evaluate L(V)LM systems.

To combat this, we first propose SESSIONINTENTBENCH tasks, consisting of four sequential subtasks tailored to systematically evaluate L(V)LMs’ capability in understanding customer intention within session browsing records. Then, we design an automated framework to streamline the collection of detailed product metadata, customer intention, and intention shift within the session by prompting L(V)LM in a multi-step manner.

By applying our method to Amazon-M2 (Jin et al., 2023b), we first filter and collect 10,905 sessions with complete textual and visual data. We enrich the original session with intention entries and obtain 1,132,145 possible intention pathways. After that, we further conduct human annotations

to 8,980 sampled intention trajectories to form an evaluation benchmark. Then, we carry out extensive experiments over more than 20 L(V)LMs by applying different evaluation settings and prompting techniques, along with extra fine-tunings. Our findings indicate that current L(V)LMs struggle with the proposed tasks. Further analyses reveal potential underlying causes behind the observed low model accuracy and introduce intention injection as a possible way of assisting models’ understanding of session intent and improving performances. We will make our code, data, and models publicly available after acceptance.

2 Related Works

2.1 Intention Understanding

Intention is the internal mental state that affects people’s decision-making (Alford and Biswas, 2002). By analyzing the inner intention states of the users, service providers are able to present more personalized products (Dai et al., 2006) and give back more accurate responses (Zhang et al., 2016). In E-commerce, customer intention is crucial in understanding their purchase behaviors and preferences (Shim et al., 2001). There has been ongoing research trying to decode how to model shopping intention. For example, using history information like tags (Wang et al., 2025a) and co-buy behaviors (Yu et al., 2023; Xu et al., 2024; Wang et al., 2025b). Recently, studies show that LLMs are struggling to connect the dots between intended products and user intention (Ding et al., 2024). However, figuring out the items the user wanted is even more difficult when it comes to more complex settings like session histories (Jin et al., 2023b). To bridge the gap between understanding intention and providing more precise shopping aids, we formulate SESSIONINTENTBENCH tasking L(V)LMs to infer intent by leveraging session metadata from multiple angles.

2.2 Purchase Session in E-commerce

Purchase session is a record of customer interaction history, which has been becoming an increasingly hot area of research (Alves Gomes et al., 2022; Jia et al., 2023; Wang et al., 2024c). Various methods are proposed trying to exploit the abundant information contained here, such as using deep reinforcement learning models (Bharadwaj et al., 2022), leveraging graph neural networks (Jin et al., 2023a), and carrying out complex logical reasoning tech-

niques (Liu et al., 2023b). While Jin et al. (2023b) systematically introduces session information as an important factor for understanding sequential interacting behavior, Liu et al. (2023b) points out that product attributes play a pivotal role in enhancing user intent capture. This shows that a more fine-grained framework of session intention evaluation is needed. Furthermore, recognizing that multiple intentions can coexist within a session, researchers have explored various approaches to enhance product recommendations. Sun et al. (2024) iteratively updates an intention ranking prompt to optimize recommendations, while Choi et al. (2024) train a neural network to learn intention embedding representations and refine selections accordingly. While these works aim to provide more precise product recommendation, our research focuses on improving language models’ intention understanding and reasoning ability using semantic intention representation. Using the summarization and generation ability of L(V)LMs, in SESSIONINTENTBENCH, we extract and incorporate session intent metadata from multiple aspects for more comprehensive intention capturing.

3 Problem Definition

3.1 SESSIONINTENTBENCH Task Definitions

First of all, we give definitions to the tasks in SESSIONINTENTBENCH. We propose to model the intention shift from four aspects as a comprehensive formulation, as outlined in Figure 2, to facilitate the creation of a L(V)LM shopping agent that is able to: (i) Detect the attribute that is decisive in the intention shift. (ii) Model intention trajectories with mined attributes and leverage them to give better predictions on future interactions. (iii) Compare between the most recently viewed product with previously interacted ones and use this comparison to validate the plausibility of the inferred intent. (iv) Leverage modeled intention trajectories to predict future product interaction preferences.

To this end, we propose tasks that each emphasize a different angle of analysis. Assume we have collected the customer interaction history over time steps $t = 1, 2, \dots, T$, i.e., the interacted products P_1, P_2, \dots, P_T and attributes that affect customer decision-making at each step A_1, A_2, \dots, A_T . Then the history information up till time step t can be summarized as $\mathcal{H}_t = \{(P_j, A_j)\}_{j=1}^t$. We denote inferred customer intention as I_1, I_2, \dots, I_T , and comparisons between

interacted items and internal intent of the current step and previous step C_1, C_2, \dots, C_T . Further discussions on the theory, intuition, task clarifications, and additional details can be found in the Appendix B, C.

TASK 1: Intent-Based Purchasing Likelihood

Estimation: The first task asks the model to verify whether the last proposed intention is a good alignment with the new product we are going to interact with. The model will be given history information \mathcal{H}_{t-1} , the proposed intention I_{t-1} , and new product P_t . It is asked to output a likelihood estimation score $\mathcal{S}_1(P_t, I_{t-1}) \mid_{\mathcal{H}_{t-1}} \in \{0, 1, 2, 3\}$ for the customer to interact with P_t , where 3 means the most likely and 0 means the least probable.

TASK 2: Purchasing Likelihood Inference via Valued Attributes Regularization

The second task requires the model to verify whether the proposed valued attributes of the user are an essential element of the actual unseen product. The model is provided with the history information \mathcal{H}_{t-1} , the proposed valued attribute A_{t-1} , and the new unseen product P_t . The model is required to output an estimated interaction likelihood score $\mathcal{S}_2(P_t, A_{t-1}) \mid_{\mathcal{H}_{t-1}} \in \{0, 1, 2, 3\}$ for the user to interact with P_t under the assumption that the user values the product feature A_{t-1} , where 3 means the most likely and 0 means the least probable.

TASK 3: Intention Justification via Comparison

To ensure that the proposed intent is reasonable and to verify against potential hallucinations, the third task asks the model to justify whether the proposed C_t provides a reasonable justification for the user to interact with P_t after seeing P_{t-1} . Formally, the model is tasked to output a score $\mathcal{S}_3(C_t, P_{t-1}, I_{t-1}, P_t, I_t) \mid_{\mathcal{H}_{t-1}} \in \{0, 1, 2, 3\}$ indicating the plausibility of the generated comparison.

TASK 4: Intention Evolution Modeling

The final task we proposed aims to test the model’s ability to help the recommendation systems decide whether to further recommend similar products or not. Providing the model with all the historical information and inferred purchasing intent, we ask it to choose from exposing the user to (a) Similar products under the same category, (b) Products with different features but still under the same category, (c) Products under different category (exploring more to figure out user preferences). If we map the choices to numerical score $\{1, 2, 3\}$, then we formalize the task as questioning for $\mathcal{S}_4(exploration, I_t) \mid_{\mathcal{H}_t}$

264 $\in \{1, 2, 3\}$. Note that the degree of exploitation
 265 decreases and exploration increases as the score
 266 increases.

267 3.2 Dataset

268 We obtain products in series of sequential interac-
 269 tions from Amazon-M2 dataset (Jin et al., 2023b)
 270 and product image information from Amazon Re-
 271 view Dataset (Hou et al., 2024). We leverage the
 272 abundant textual information (such as titles, price,
 273 color, material, etc.) mentioned in Amazon-M2
 274 and retrieve corresponding product images from
 275 Amazon Review Dataset to curate the dataset. Af-
 276 ter filtering out products whose links are missing or
 277 not accessible, we obtained 10,905 sessions with
 278 complete textual and visual components.

279 4 SESSIONINTENTBENCH Construction

280 In this section, we present our methodology of con-
 281 structing the intention tree from the source data
 282 we collected and how we curated the SESSIONIN-
 283 TENTBENCH. An overview is presented in Fig-
 284 ure 2. Our framework consists of four steps: (i)
 285 Extract attributes for session products to provide
 286 aids for model inference in later steps. (ii) For
 287 each time step, prompt the models to mimic cus-
 288 tomer behavior and infer multiple intentions from
 289 previous interactions. (iii) Enrich intention tree
 290 structure with more nuanced inter-session intention
 291 metadata analyses, which is taken from multiple
 292 perspectives on how and why intention shifts over
 293 time. (iv) Conduct human annotations for the con-
 294 structed tree

295 4.1 Multi-modal Attribute Extraction

296 The first step aims to extract product attributes that
 297 can better assist LVLMS in analyzing user intention
 298 shift in later stages. To achieve this, we use GPT-
 299 4o-mini (OpenAI, 2024a) as the extraction tool
 300 and provide it with ensembled textual and visual
 301 information of session products. The LVLMS is
 302 then asked to output a general classification of the
 303 product itself, and to categorize then instantiate
 304 the extractable features of the product, for example
 305 (e.g., *color: white, size: 7.5 inches*).

306 4.2 Customer Intention Generation

307 To build up the intention tree based on the product
 308 purchase session, we first fill up the tree bones
 309 with predicted user intention using L(V)LMs. The
 310 intentions are inferred at each time step following
 311 the session time frame. Starting with the first item

312 in the session, we ask the model to infer a list of
 313 possible intention $\langle I_{t1}, I_{t2}, I_{t3}, \dots \rangle \big|_{t=1}$ based
 314 on textual and visual information of the product
 315 user interacted, where the prompt is demonstrated
 316 below. Then, repeat the inference every step as we
 317 add the next new session product into the visible
 318 list of items to the model.

319 To make the intention instantiation successional,
 320 we add the intention information of the previous
 321 time step $\{I_i\}_{i=1}^{t-1}$ (**<Prev Intent>**) to facilitate
 322 the model to do the reasoning. And at each time
 323 we do the inference, we will only use one inten-
 324 tion chosen from the previous step intention, to
 325 ensure the coherent intention trajectory sampling.
 326 More specifically, the model is constraint to output
 327 the five most possible user intentions, denoted as
 328 $\{\langle \text{New Intent } i \rangle\}_{i=1}^5$, prior to the fifth product at
 329 each iteration. This process is referred to as *branch-*
 330 *ing*, as it resembles the growth of a tree, wherein
 331 each new intention branches out from the initial
 332 concept, akin to twigs dividing into finer branches.
 333 And starting from the fifth product, we only infer
 334 one possible intention at a time to control the ex-
 335 ponential growth of the tree size (by setting $|\langle \text{New}$
 336 **Intent>|=1**).

```
337 <TASK-PROMPT>
338 <INPUT:>
339 <Prev Intent><Prev Products><New Product>
340 <OUTPUT:>
341 <New Intent 1><Attr 1><Rationale 1><Comp 1>
342 <New Intent 2><Attr 2><Rationale 2><Comp 2>
343 ...
344 <New Intent 5><Attr 5><Rationale 5><Comp 5>
345 <INPUT:>
346 <Prev Intent><Prev Products><New Product>
347 <OUTPUT:>
```

348 4.3 Inter-Sessions Intention Shift Analysis

349 Following this, we want to investigate the specific
 350 reasons behind each intention shift before and af-
 351 ter the customer sees each product and how that
 352 might influence the further decision-making of the
 353 customer. The prompt we used for generation is
 354 given above. To ground the reasoning on the ac-
 355 tual product metadata, we require the model to
 356 point out the most likely feature **<Attr>** A_t that
 357 affects the user choices. Furthermore, we ask for
 358 a more comprehensive comparison **<Comp>** C_t be-
 359 tween the last product P_t the previous one P_{t-1} ,
 360 so that it provides logical support for the modeled
 361 intention pathways. In order to help models reason
 362 better, we require the model to provide rationales
 363

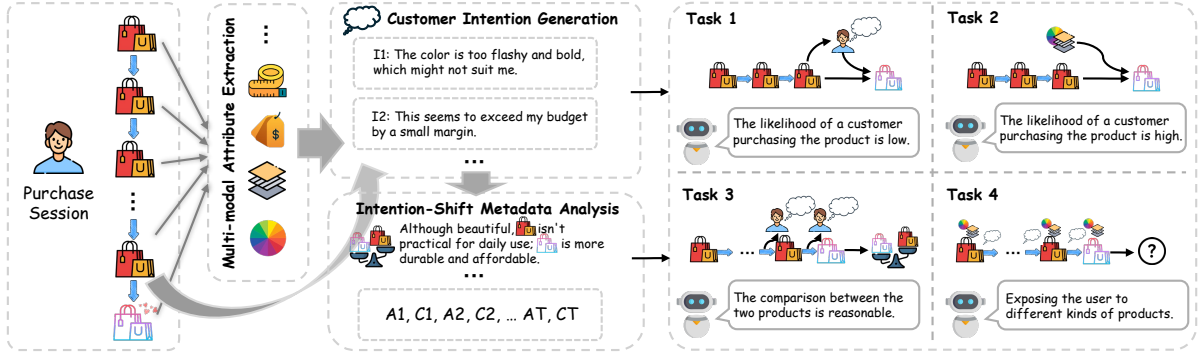


Figure 2: Overview of SESSIONINTENTBENCH and the construction pipeline. Multi-modal attribute extraction is conducted first as an aid for further step intention generation. Metadata analyses are conducted afterward to provide a more fine-grained and detailed inspection of intention shifts in the session interaction. Here, A_i, C_i stands for attributes and comparisons at i -th step. Different task is associated with different collections of metadata.

Genre	Property	Train	Test
Basic Info	# Sessions (uni.)	8963	5306
	# Sampled Tasks	28736	7184
	Avg. # Products	3.4163	3.4123
	Avg. # Intention	3.4163	3.4123
Session Len	# Len = 3	18956	4752
	# Len = 4	7598	1902
	# Len = 5	2182	530
Task Num	# TASK 1	7153	1827
	# TASK 2	7171	1809
	# TASK 3	7154	1826
	# TASK 4	7258	1722

Table 1: Statistics of the sampled and annotated data for the SESSIONINTENTBENCH benchmark. *uni.* means unique sessions are included. Note that the dataset is randomly shuffled and then sampled by 13,003,664 tasks or 1,132,145 intention trajectories, not by sessions.

(**<Rationale>**) behind the generations as part of the output. We collect this analysis metadata in the format of one general categorization plus one detailed instantiation, e.g., *book type: fiction, price: \$20*.

4.4 Human Annotation

We hire Amazon Mechanical Turk annotators to label a randomly sampled subset of our data to balance cost and quality. We ask the workers to annotate emphasizing the following perspectives: (1) The alignment of proposed intention I_t and session products P_{t+1} . (2) The consistency between the inferred valued attribute A_t and actual interacted products P_{t+1} . (3) The plausibility of the generated intention comparison C_t . (4) Predictions on further intention pathways based on historical information. In this way, the session intention could not only provide insights into the thinking process of customers but also meaningful references for when to explore and when to exploit product recommendation systems. To simplify the annotation process,

the annotators are only asked to assign a likelihood score or plausibility score for each task in a format roughly similar to *yes, maybe yes, maybe no, no* (corresponds to $S = 3, 2, 1, 0$). We carried out multiple rounds of annotation worker selections with different criteria to ensure high annotation quality. More details are in Appendix D.

5 Evaluations and Analyses

5.1 Intrinsic Evaluations

We present our detailed statistics in Table 1. By filling up the tree with intention on 10,905 sessions, we obtain more than 1,950,000 intention entries and 1,100,000 intention trajectories. The majority of these sessions contains less than four products, though long sessions also exist with up to 18 products. To sample a subset of sessions to form the SESSIONINTENTBENCH, we first retrieve candidate sessions with lengths three to five. We then sample 2,000 sessions with 2 trajectories per session and later add another disjoint 1,445 sessions with 4 trajectories per session. That gives 9,780 trajectories in total. To grant the model with full information available, we only query the tasks at the end of each session time step, that is, using all the products available and masking the last product when querying the TASK 1, 2.

5.2 Baselines and Model Selections

Evaluation Metric We use accuracy and Macro-F1 score as evaluation metrics. Accuracy is defined as the percentage of questions that are correctly answered. We regard scoring 0,1 in TASK 1-3 as the true positive label and scoring 0 as the one for TASK 4. To start with, we include the Random Selection and Majority Vote score of each task.

Models	Intent-Based Inference		Valued Attributes Reg.		Comparison Just.		Evolution Modeling	
	Acc	Ma-F1	Acc	Ma-F1	Acc	Ma-F1	Acc	Ma-F1
Random	50.00	50.00	50.00	50.00	50.00	50.00	54.38	35.00
Majority	62.30	76.77	54.35	NaN	71.80	83.58	63.15	NaN
LLM (Zero-Shot)								
Meta-Llama-3.1-8B	56.87	70.98	49.36	55.10	71.30	<u>83.24</u>	39.26	53.01
Gemma-2-9B	57.03	69.37	52.18	49.44	41.68	44.19	<u>53.77</u>	34.54
Mistral-7B-v0.3	62.17	<u>76.52</u>	47.65	64.08	71.30	<u>83.24</u>	39.61	53.53
Ministral-8B	56.98	69.33	51.58	50.48	68.02	80.48	38.27	54.08
Mistral-Nemo-12B	53.09	63.82	51.63	35.04	56.79	69.71	47.15	45.11
Falcon-3-7B	57.31	71.74	<u>52.24</u>	49.17	67.36	79.41	44.36	49.68
Falcon-3-10B	54.95	66.93	51.35	48.59	65.49	78.24	43.84	45.89
Qwen-2.5-3B	54.19	64.42	51.96	41.87	68.62	81.01	37.63	<u>53.98</u>
Qwen-2.5-7B	58.62	71.92	51.02	56.18	70.59	82.61	40.07	51.86
LVLm (Zero-Shot)								
LLaVA-v1.6-mistral-7b	58.29	71.90	47.48	62.27	62.94	75.11	37.62	54.20
LLaVA-v1.6-vicuna-7b	<u>62.01</u>	<u>76.55</u>	46.93	63.88	<u>71.27</u>	<u>83.22</u>	37.21	54.24
Qwen-2-VL-7B	58.73	71.48	<u>50.63</u>	56.37	70.61	<u>82.73</u>	<u>37.67</u>	53.95
Meta-Llama-3.2-11B-V	45.10	61.38	38.41	52.35	42.11	59.20	36.33	53.23
L(V)Lm (Few-shots)								
Mistral-7B-v0.3	<u>60.43</u>	74.60	<u>50.64</u>	61.39	<u>67.09</u>	79.08	43.44	49.85
Qwen-2-VL-2B	58.02	73.46	40.63	58.40	66.70	79.92	36.99	53.45
LLaVA-v1.6-vicuna-7b	51.06	77.26	22.61	<u>62.92</u>	66.81	<u>82.99</u>	27.99	<u>54.07</u>
L(V)Lm (Fine-tuned)								
Meta-Llama-3.1-8B	52.82	63.84	51.46	46.27	70.76	82.82	51.92	33.01
Meta-Llama-3.2-3B	55.67	66.80	51.80	46.70	69.61	81.93	51.63	32.66
Mistral-7B-v0.3	57.47	68.56	50.64	44.64	67.69	79.88	55.69	31.69
Ministral-8B	<u>58.35</u>	69.55	51.24	45.01	66.54	79.10	55.57	35.11
Mistral-Nemo-12B	56.10	66.80	52.02	46.68	67.74	79.81	<u>55.81</u>	32.95
Qwen-2.5-7B	54.02	65.63	52.02	46.75	69.50	81.66	54.47	31.59
Falcon-3-7B	55.77	65.02	52.85	<u>48.46</u>	71.41	83.30	54.65	<u>36.86</u>
L(V)Lm (Proprietary API)								
GPT4o-mini	57.44	69.34	51.95	43.81	<u>71.19</u>	<u>83.13</u>	38.39	<u>53.90</u>
GPT4o-mini (5-shots)	58.83	71.86	49.32	<u>53.01</u>	65.25	78.11	46.51	46.96
GPT4o-mini (COT)	57.26	69.02	51.87	43.33	68.86	81.22	42.81	49.42
GPT4o	55.05	65.33	49.75	36.27	56.30	67.51	41.64	52.39
GPT4o (5-shots)	53.10	63.58	44.20	38.61	54.94	65.01	43.44	48.41
GPT4o (COT)	53.30	61.91	<u>52.00</u>	36.08	49.50	50.87	58.42	13.73

Table 2: Evaluation results (%) of various (L)LMs on the annotated testing sets of SESSIONINTENTBENCH. The best performances within each method are underlined, and the best among all methods are **bold-faced**.

Model Selections Then, we test out a diverse set of L(V)LMs on SESSIONINTENTBENCH. Since all the tasks we proposed belong to classification setups, we choose accuracy and Macro-F1 score as evaluation metrics. The models we selected, as given in Table 2, can be classified into three genres: (I) **OPEN L(V)LMs WITH ZERO-SHOT**: Firstly, we select a vast collection of models from different companies or organizations. Text-to-text models includes Llama3.1, Llama3.2 (Grattafiori et al., 2024), Gemma2 (Team et al., 2024), Mistral (Jiang et al., 2023), Falcon (Almazrouei et al., 2023), and Qwen2.5 (Qwen et al., 2025). Image-text-to-text models includes LLaVA (Liu et al., 2023a), Qwen2-VL (Wang et al., 2024a), and Llama with Vision (Grattafiori et al., 2024). Models under this category are prompted using zero-shot. (II) **FINE-TUNED L(V)LMs WITH ZERO-SHOT**: Following that, we fine-tuned Llama3.1, Llama3.2, Mistral, Falcon3, Qwen2.5 on partitioned training set and evaluate them on the testing set. (III) **PRO-**

RIETARY L(V)LM API WITH SEVERAL DIFFERENT PROMPTING TECHNIQUES: Lastly, we test out GPT-4o and GPT-4o-mini (OpenAI et al., 2024; OpenAI, 2024a) using zero-shot prompting, 5-shots prompting and Chain-of-Thought prompting (Wei et al., 2023).

5.3 Main Evaluation Results

INTENTION EVOLUTION MODELING (TASK 4) is the most challenging task. Our experiments show that the average accuracy of the zero-shot models on TASK 4 is 42.34%. Compared to the second hardest task (*Purchasing Likelihood Inference via Valued Attributes Regularization*), which models scored 49.63%, there is a great gap of 7.29% on TASK 4. After being fine-tuned, all open models are able to achieve a minimum accuracy of 51.92%, while the top performing one (Mistral-Nemo-12B) scores 55.81%, just above the RANDOM vote accuracy. It is worth noticing that GPT-4o with Chain-of-Thought prompting is able to achieve the high-

Training Data	Backbone	Intent-Based Inference		Valued Attributes Reg.		Comparison Just.		Evolution Modeling	
		Acc	Ma-F1	Acc	Ma-F1	Acc	Ma-F1	Acc	Ma-F1
Zero-shot	Llama-3.1-8B	56.87	70.98	49.36	55.10	<u>71.30</u>	<u>83.24</u>	39.26	53.01
	Llama-3.2-3B	54.68	63.97	52.02	43.48	33.13	49.48	<u>51.34</u>	36.61
	Mistral-7B-v0.3	62.17	76.52	47.65	64.08	<u>71.30</u>	<u>83.24</u>	39.61	53.53
	Ministral-8B	56.98	69.33	51.58	50.48	68.02	80.48	38.27	54.08
	Falcon-3-7B	57.31	71.74	<u>52.24</u>	49.17	67.36	79.41	44.36	49.68
	Qwen-2.5-7B	58.62	71.92	51.02	56.18	70.59	82.61	40.07	51.86
SIB	Llama-3.1-8B	52.82	63.84	51.46	46.27	70.76	82.82	51.92	33.01
	Llama-3.2-3B	55.67	66.80	51.80	46.70	69.61	81.93	51.63	32.66
	Mistral-7B-v0.3	57.47	68.56	50.64	44.64	67.69	79.88	55.69	31.69
	Ministral-8B	<u>58.35</u>	<u>69.55</u>	51.24	45.01	66.54	79.10	55.57	35.11
	Qwen-2.5-7B	54.02	65.63	52.02	46.75	69.50	81.66	54.47	31.59
	Falcon-3-7B	55.77	65.02	<u>52.85</u>	<u>48.46</u>	71.41	83.30	54.65	<u>36.86</u>
MIND + SIB	Llama-3.1-8B	<u>60.10</u>	68.81	55.33	48.67	70.54	82.54	57.72	39.74
	Llama-3.2-3B	59.88	67.92	55.28	50.15	64.02	75.48	58.54	40.50
	Mistral-7B-v0.3	60.04	<u>69.96</u>	52.90	45.87	67.69	79.56	59.93	37.16
	Ministral-8B	58.24	67.33	53.95	47.44	65.44	77.01	58.77	40.93
	Qwen-2.5-7B	59.00	67.65	53.95	48.62	63.09	74.98	57.84	39.30
	Falcon-3-7B	58.57	68.42	55.94	50.22	<u>71.30</u>	<u>83.25</u>	58.36	40.00

Table 3: Evaluation results (%) of transferring knowledge from MIND to aid SESSIONINTENTBENCH. The best performances among each method are underlined, and the best ones among all methods are **bold-faced**. We abbreviate SESSIONINTENTBENCH as SIB.

est rate of 58.42% among all models and methods. This might be because the larger model size and the trick of enabling reasoning at run time could help the model to better mimic the thinking process of a real-life customer. This result shows that more works need to be done to level up the model’s capability of capturing long-term user intention trends.

Fine-tuning can greatly improve the poor performing models, but struggle to help the mediocre ones. Poor performing models, which we referred to as the ones that receive a low score compared to models under the same category in some evaluation tasks, can quickly pick up relevant capabilities by being fine-tuned on the training set before testing. For example, LLAMA-3.2-3B shows poor performance on TASK 3 (*Intention Justification via Comparison*), but after being fine-tuned on SESSIONINTENTBENCH, it shows a leap of performance by 36.5% and demonstrates comparable outcome with other larger 7B or 8B models. The mediocre performing models, which we referred to as the ones that score near the highest among the models but still struggle to surpass the top accuracy records. Among the proposed tasks, the largest maximum accuracy raise from zero-shot to fine-tuned happens at TASK 4, with a lift of 2.04% in the highest score. As a result of these two factors, the variance between different models shrinks after fine-tuning. See the Appendix F.2 for additional discussions on fine-tuning.

LVLMS struggle to make good usage of visual signals. In comparison to LLMs, which only use textual signals as the input, LVLMS can refer to image information to facilitate their question-

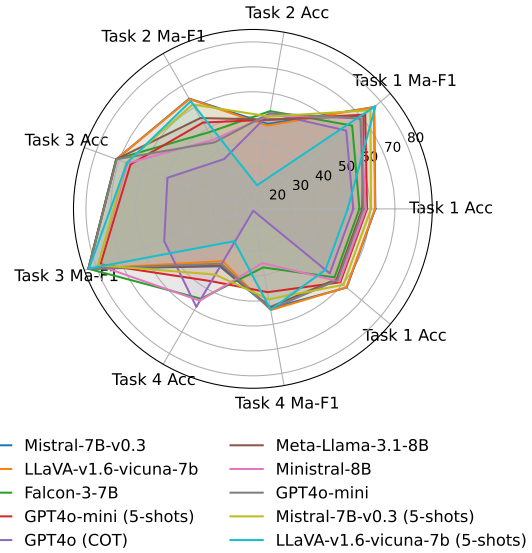


Figure 3: Radar chart of models that have the best performances in multiple tasks within each method. No single model can produce a boundary that encompasses all the data points from other models.

answering and inference reasoning. However, as shown in Table 2, the highest accuracy scores of LVLMS all lag behind compared to the highest ones of LLMs. When evaluated on TASK 4 using direct zero-shot, the best LVLMS outcome is even behind the best LLM by a huge gap of 11.27%. Possible issues could be the low signal-noise ratio of the images collected, and sellers usually include more comprehensive and concise features of products in text format.

No model dominates. The overall scoring result is quite close, especially after fine-tuning, where the variance between models shrinks. The open models that achieve the best accuracy within their category are Mistral-7B-v0.3 (zero-shot LLM), LLaVA-

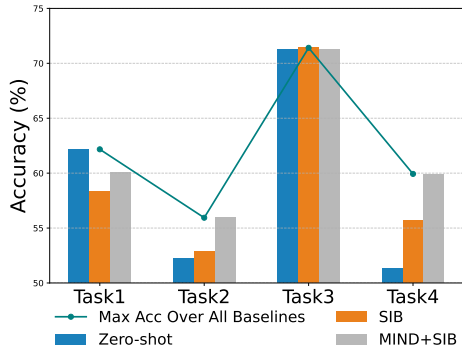


Figure 4: Comparison between best performances across different methods on different tasks. All baseline’s max accuracy line is consistent with the max accuracy line over open models of all methods (i.e., Zero-shot, fine-tuning with SESSIONINTENTBENCH (SIB), and sequential fine-tuning with MIND then SIB). Paying for proprietary API does not add extra value in this case.

v1.6-vicuna-7b (zero-shot LVLML), and Falcon-3-7B (fine-tuned LLM). They belong to different companies or organizations, and none of them achieves the best accuracy in more than two tasks, as shown in Figure 3. It is an indicator that no specific model can dominate the session intention modeling game. Although the GPT-4o with Chain-of-Thought techniques gives the best Intention Evolution Modeling results, surpassing the best fine-tuned one (Mistral-Nemo-12B) with 2.59%, it might not be cost-effective due to expensive pricing and large token consumption required.

5.4 The Impact of Intention Injection

Observing Table 2, we find L(V)LMs struggling with directly leveraging intention for next product inference (*Intent-Based Inference*) and mastering long term trend of shifting intention from session history (*Intention Evolution Modeling*). However, this can be caused by many factors, such as failing to capture diverse user characteristics and preferences. Given this, we hypothesize that fine-tuning models on intention knowledge bases beforehand might enhance their ability to adapt according to different background setups and help them generalize better. Therefore, we tried to first fine-tune models on MIND (Xu et al., 2024), a large-scale intention knowledge base grounded on co-buy behaviors, as shown in Table 3.

After sequentially fine-tuning first on MIND and then on SESSIONINTENTBENCH, we find a leap in performance in the proposed tasks. By injecting intention from MIND to aid SESSIONINTENTBENCH, we improved TASK 1 by 1.75%, TASK 2 by 3.09%, TASK 4 by 4.24% (another great

improvement compared to zero-shot baseline), as demonstrated in Figure 4. This demonstrates that intention injection can be an effective technique to improve the model’s ability to identify user intention from a short yet complex series of interactions. See Appendix F for more discussions on intention injection.

5.5 Error Analyses

We randomly sample 200 tasks where GPT-4o with Chain-of-Thought commits an error. And we recruit experts to analyze the causes behind them manually. Our results show that:

- 47.5% errors are caused by incorrect understanding of the provided metadata. This may be because the model fails to incorporate past product information for deeper comprehension.
- 24% errors are caused by incorrect ground-truth labels. For objective factors, this might be due to internal conflict of session products and metadata from the intention tree or incorporating complicated metadata in the label instruction.
- 7% errors are due to models’ failure to capture important product features contained in the session products, which might be aligned with or different from the metadata described in the problem assumption.
- 6.5% errors are due to irrelevant reasoning or model hallucinations, where the model is often heading towards a different reasoning direction due to some misleading, unimportant features.
- 15% the errors are due to models’ inability to capture the overall intention of the customer when the provided metadata is vague or not decisive when estimating the likelihood.

6 Conclusions

In conclusion, we propose an automated pipeline to construct a large-scale knowledge base and further construct a sample dataset SESSIONINTENTBENCH for L(V)LMs evaluations. Extensive experiments show that current models struggle to understand and infer customers’ intentions while injecting intention from other knowledge bases can level up the performance. We hope our work can bridge the gap between intention understanding in simplified research cases like co-buy intention and more complex yet practical scenarios like session history. We hope this framework can benefit the community by providing better services with future models.

582 Limitations

583 We implemented our intention tree construction
584 pipeline using GPT-4o-mini as the metadata gener-
585 ator. As LLM space advances, more advanced
586 models like GPT-o1 (OpenAI, 2024b), GPT-o3-
587 mini (OpenAI, 2025) will become more accessi-
588 ble to researches, which would potentially better
589 mimic customer thinking process and behavior and
590 generate intention and metadata in higher standard.
591 This would enable our knowledge base and dataset
592 generation with even higher quality.

593 Our current intention modeling process does not
594 incorporate additional personalized factors such
595 as past purchases, user characteristics, and social
596 relationships with other customers. Incorporating
597 these variables, which can be precomputed, could
598 enhance model reasoning during inference, thereby
599 providing more accurate modeling of session intent
600 for specific customers.

601 The modeling setting we proposed contains mul-
602 tiple perspectives of session intent metadata, includ-
603 ing attributes, intention, and comparisons. How-
604 ever, more metadata mined from the session can
605 possibly be added for further knowledge integra-
606 tions and better utilization of available information.
607 More work can be done to explore what other inter-
608 nal factors can be incorporated within the session
609 itself.

610 Ethics Statement

611 **Offensive Content Inspection** We leverage the
612 generation capability of L(V)LMs to construct a
613 knowledge base and carry out experiments. The
614 generated intention at the dataset construction step
615 is closely related to the session product information
616 itself. The remaining metadata is based on the rea-
617 soning and comparison within products and related
618 intentions. As the experiment setting, we only ask
619 models to give out specific scores of likelihood or
620 generate content with constraint reasoning, which
621 is also closely related to sessions and products.

622 **Annotation Wage** The annotators are paid a
623 wage in compliance with the local law, on an aver-
624 age of 15 USD per hour. They have all agreed to
625 participate in annotation voluntarily.

626 **Licenses** Amazon-M2 dataset are released under
627 the license of Apache 2.0. This grants our access
628 to the dataset for free. Our code and data will be
629 shared under the MIT license. It will allow the free
630 distribution of assets we proposed and curated. All

associated licenses permit user access for research
purposes, and we have agreed to follow all terms
of use.

References

- Bruce L Alford and Abhijit Biswas. 2002. [The effects of discount level, price consciousness and sale proneness on consumers' price perception and behavioral intention](#). *Journal of Business Research*, 55(9):775–783.
- Ebtesam Almazrouei, Hamza Alobeidli, Abdulaziz Alshamsi, Alessandro Cappelli, Ruxandra Cojocaru, Mérouane Debbah, Étienne Goffinet, Daniel Hesslow, Julien Launay, Quentin Malartic, Daniele Mazzotta, Badreddine Noune, Baptiste Pannier, and Guilherme Penedo. 2023. [The falcon series of open language models](#). *Preprint*, arXiv:2311.16867.
- Miguel Alves Gomes, Richard Meyes, Philipp Meisen, and Tobias Meisen. 2022. [Will this online shopping session succeed? predicting customer's purchase intention using embeddings](#). In *Proceedings of the 31st ACM International Conference on Information & Knowledge Management, CIKM '22*, page 2873–2882, New York, NY, USA. Association for Computing Machinery.
- Diddigi Raghu Ram Bharadwaj, Lakshya Kumar, Saif Jawaid, and Sreekanth Vempati. 2022. [Fine-grained session recommendations in e-commerce using deep reinforcement learning](#). *Preprint*, arXiv:2210.15451.
- Minjin Choi, Hye-young Kim, Hyunsouk Cho, and Jongwuk Lee. 2024. [Multi-intent-aware session-based recommendation](#). In *Proceedings of the 47th International ACM SIGIR Conference on Research and Development in Information Retrieval, SIGIR '24*, page 2532–2536, New York, NY, USA. Association for Computing Machinery.
- Honghua (Kathy) Dai, Lingzhi Zhao, Zaiqing Nie, Ji-Rong Wen, Lee Wang, and Ying Li. 2006. [Detecting online commercial intention \(oci\)](#). In *Proceedings of the 15th International Conference on World Wide Web, WWW '06*, page 829–837, New York, NY, USA. Association for Computing Machinery.
- Wenxuan Ding, Weiqi Wang, Sze Heng Douglas Kwok, Minghao Liu, Tianqing Fang, Jiaxin Bai, Xin Liu, Changlong Yu, Zheng Li, Chen Luo, Qingyu Yin, Bing Yin, Junxian He, and Yangqiu Song. 2024. [IntentionQA: A benchmark for evaluating purchase intention comprehension abilities of language models in E-commerce](#). In *Findings of the Association for Computational Linguistics: EMNLP 2024*, pages 2247–2266, Miami, Florida, USA. Association for Computational Linguistics.
- Aaron Grattafiori, Abhimanyu Dubey, Abhinav Jauhri, Abhinav Pandey, Abhishek Kadian, Ahmad Al-Dahle, Aiesha Letman, Akhil Mathur, Alan Schelten, Alex Vaughan, Amy Yang, Angela Fan, Anirudh

686	Goyal, Anthony Hartshorn, Aobo Yang, Archi Mitra, Archie Sravankumar, Artem Korenev, Arthur Hinsvark, Arun Rao, Aston Zhang, Aurelien Rodriguez, Austen Gregerson, Ava Spataru, Baptiste Roziere, Bethany Biron, Binh Tang, Bobbie Chern, Charlotte Caucheteux, Chaya Nayak, Chloe Bi, Chris Marra, Chris McConnell, Christian Keller, Christophe Touret, Chunyang Wu, Corinne Wong, Cristian Canton Ferrer, Cyrus Nikolaidis, Damien Alonsius, Daniel Song, Danielle Pintz, Danny Livshits, Danny Wyatt, David Esiobu, Dhruv Choudhary, Dhruv Mahajan, Diego Garcia-Olano, Diego Perino, Dieuwke Hupkes, Egor Lakomkin, Ehab AlBadawy, Elina Lobanova, Emily Dinan, Eric Michael Smith, Filip Radenovic, Francisco Guzmán, Frank Zhang, Gabriel Synnaeve, Gabrielle Lee, Georgia Lewis Anderson, Govind Thattai, Graeme Nail, Gregoire Mialon, Guan Pang, Guillem Cucurell, Hailey Nguyen, Hannah Korevaar, Hu Xu, Hugo Touvron, Iliyan Zarov, Imanol Arrieta Ibarra, Isabel Kloumann, Ishan Misra, Ivan Evtimov, Jack Zhang, Jade Copet, Jaewon Lee, Jan Geffert, Jana Vranes, Jason Park, Jay Mahadeokar, Jeet Shah, Jelmer van der Linde, Jennifer Billock, Jenny Hong, Jenya Lee, Jeremy Fu, Jianfeng Chi, Jianyu Huang, Jiawen Liu, Jie Wang, Jiecao Yu, Joanna Bitton, Joe Spisak, Jongsoo Park, Joseph Rocca, Joshua Johnstun, Joshua Saxe, Junteng Jia, Kalyan Vasuden Alwala, Karthik Prasad, Kartikeya Upasani, Kate Plawiak, Ke Li, Kenneth Heafield, Kevin Stone, Khalid El-Arini, Krithika Iyer, Kshitiz Malik, Kuenley Chiu, Kunal Bhalla, Kushal Lakhotia, Lauren Rantala-Yearly, Laurens van der Maaten, Lawrence Chen, Liang Tan, Liz Jenkins, Louis Martin, Lovish Madaan, Lubo Malo, Lukas Blecher, Lukas Landzaat, Luke de Oliveira, Madeline Muzzi, Mahesh Pasupuleti, Mannat Singh, Manohar Paluri, Marcin Kardas, Maria Tsimpoukelli, Mathew Oldham, Mathieu Rita, Maya Pavlova, Melanie Kambadur, Mike Lewis, Min Si, Mitesh Kumar Singh, Mona Hassan, Naman Goyal, Narjes Torabi, Nikolay Bashlykov, Nikolay Bogoychev, Niladri Chatterji, Ning Zhang, Olivier Duchenne, Onur Çelebi, Patrick Alrassy, Pengchuan Zhang, Pengwei Li, Petar Vasic, Peter Weng, Prajjwal Bhargava, Pratik Dubal, Praveen Krishnan, Punit Singh Koura, Puxin Xu, Qing He, Qingxiao Dong, Ragavan Srinivasan, Raj Ganapathy, Ramon Calderer, Ricardo Silveira Cabral, Robert Stojnic, Roberta Raileanu, Rohan Maheswari, Rohit Girdhar, Rohit Patel, Romain Sauvestre, Ronnie Polidoro, Roshan Sumbaly, Ross Taylor, Ruan Silva, Rui Hou, Rui Wang, Saghar Hosseini, Sahana Chennabasappa, Sanjay Singh, Sean Bell, Seohyun Sonia Kim, Sergey Edunov, Shaoliang Nie, Sharan Narang, Sharath Rapparthi, Sheng Shen, Shengye Wan, Shruti Bhosale, Shun Zhang, Simon Vandenhende, Soumya Batra, Spencer Whitman, Sten Sootla, Stephane Collot, Suchin Gururangan, Sydney Borodinsky, Tamar Herman, Tara Fowler, Tarek Sheasha, Thomas Georgiou, Thomas Scialom, Tobias Speckbacher, Todor Mihaylov, Tong Xiao, Ujjwal Karn, Vedanuj Goswami, Vibhor Gupta, Vignesh Ramanathan, Viktor Kerkez, Vincent Gouget, Virginie Do, Vish Vogeti, Vitor Albiero, Vladan Petrovic, Weiwei Chu, Wenhan Xiong, Wenyin Fu, Whitney Meers, Xavier Martinet, Xiaodong Wang, Xiaofang Wang, Xiaoqing Ellen Tan, Xide Xia, Xinfeng Xie, Xuchao Jia, Xuwei Wang, Yaelle Goldschlag, Yashesh Gaur, Yasmine Babaei, Yi Wen, Yiwen Song, Yuchen Zhang, Yue Li, Yuning Mao, Zacharie Delpierre Coudert, Zheng Yan, Zhengxing Chen, Zoe Papakipos, Aaditya Singh, Aayushi Srivastava, Abha Jain, Adam Kelsey, Adam Shajnfeld, Adithya Gangidi, Adolfo Victoria, Ahuva Goldstand, Ajay Menon, Ajay Sharma, Alex Boesenberg, Alexei Baevski, Allie Feinstein, Amanda Kallet, Amit Sangani, Amos Teo, Anam Yunus, Andrei Lupu, Andres Alvarado, Andrew Caples, Andrew Gu, Andrew Ho, Andrew Poulton, Andrew Ryan, Ankit Ramchandani, Annie Dong, Annie Franco, Anuj Goyal, Aparajita Saraf, Arkabandhu Chowdhury, Ashley Gabriel, Ashwin Bharambe, Assaf Eisenman, Azadeh Yazdan, Beau James, Ben Maurer, Benjamin Leonhardi, Bernie Huang, Beth Loyd, Beto De Paola, Bhargavi Paranjape, Bing Liu, Bo Wu, Boyu Ni, Braden Hancock, Bram Wasti, Brandon Spence, Brani Stojkovic, Brian Gamido, Britt Montalvo, Carl Parker, Carly Burton, Catalina Mejia, Ce Liu, Changan Wang, Changkyu Kim, Chao Zhou, Chester Hu, Ching-Hsiang Chu, Chris Cai, Chris Tindal, Christoph Feichtenhofer, Cynthia Gao, Damon Civin, Dana Beaty, Daniel Kreymer, Daniel Li, David Adkins, David Xu, Davide Testuggine, Delia David, Devi Parikh, Diana Liskovich, Didem Foss, DingKang Wang, Duc Le, Dustin Holland, Edward Dowling, Eissa Jamil, Elaine Montgomery, Eleonora Presani, Emily Hahn, Emily Wood, Eric-Tuan Le, Erik Brinkman, Esteban Arcaute, Evan Dunbar, Evan Smothers, Fei Sun, Felix Kreuk, Feng Tian, Filippos Kokkinos, Firat Ozgenel, Francesco Caggioni, Frank Kanayet, Frank Seide, Gabriela Medina Florez, Gabriella Schwarz, Gada Badeer, Georgia Swee, Gil Halpern, Grant Herman, Grigory Sizov, Guangyi, Zhang, Guna Lakshminarayanan, Hakan Inan, Hamid Shojanazeri, Han Zou, Hannah Wang, Hanwen Zha, Haroun Habeeb, Harrison Rudolph, Helen Suk, Henry Aspegren, Hunter Goldman, Hongyuan Zhan, Ibrahim Damlaj, Igor Molybog, Igor Tufanov, Ilias Leontiadis, Irina-Elena Veliche, Itai Gat, Jake Weissman, James Geboski, James Kohli, Janice Lam, Japhet Asher, Jean-Baptiste Gaya, Jeff Marcus, Jeff Tang, Jennifer Chan, Jenny Zhen, Jeremy Reizenstein, Jeremy Teboul, Jessica Zhong, Jian Jin, Jingyi Yang, Joe Cummings, Jon Carvill, Jon Shepard, Jonathan McPhie, Jonathan Torres, Josh Ginsburg, Junjie Wang, Kai Wu, Kam Hou U, Karan Saxena, Kartikay Khandelwal, Katayoun Zand, Kathy Matosich, Kaushik Veeraraghavan, Kelly Michelena, Keqian Li, Kiran Jagadeesh, Kun Huang, Kunal Chawla, Kyle Huang, Lailin Chen, Lakshya Garg, Lavender A, Leandro Silva, Lee Bell, Lei Zhang, Liangpeng Guo, Licheng Yu, Liron Moshkovich, Luca Wehrstedt, Madian Khabsa, Manav Avalani, Manish Bhatt, Martynas Mankus, Matan Hasson, Matthew Lennie, Matthias Reso, Maxim Groshev, Maxim Naumov, Maya Lathi, Meghan Keneally, Miao Liu, Michael L. Seltzer, Michal Valko, Michelle Restrepo, Mihir Patel, Mik Vyatskov, Mikayel Samvelyan, Mike Clark, Mike Macey, Mike Wang, Miquel Jubert Hermoso,	750 751 752 753 754 755 756 757 758 759 760 761 762 763 764 765 766 767 768 769 770 771 772 773 774 775 776 777 778 779 780 781 782 783 784 785 786 787 788 789 790 791 792 793 794 795 796 797 798 799 800 801 802 803 804 805 806 807 808 809 810 811 812 813
-----	--	--

814	Mo Metanat, Mohammad Rastegari, Munish Bansal,	Dong, Liang Zhang, Lei Cheng, and Linjian Mo.	875
815	Nandhini Santhanam, Natascha Parks, Natasha	2023. Smone: A session-based recommendation	876
816	White, Navyata Bawa, Nayan Singhal, Nick Egebo,	model based on neighbor sessions with similar proba-	877
817	Nicolas Usunier, Nikhil Mehta, Nikolay Pavlovich	bilistic intentions . <i>ACM Trans. Knowl. Discov. Data</i> ,	878
818	Laptev, Ning Dong, Norman Cheng, Oleg Chernoguz,	17(8).	879
819	Olivia Hart, Omkar Salpekar, Ozlem Kalinli, Parkin		
820	Kent, Parth Parekh, Paul Saab, Pavan Balaji, Pe-	Albert Q. Jiang, Alexandre Sablayrolles, Arthur Men-	880
821	dro Rittner, Philip Bontrager, Pierre Roux, Piotr	sch, Chris Bamford, Devendra Singh Chaplot, Diego	881
822	Dollar, Polina Zvyagina, Prashant Ratanchandani,	de las Casas, Florian Bressand, Gianna Lengyel, Guil-	882
823	Pritish Yuvraj, Qian Liang, Rachad Alao, Rachel	laume Lample, Lucile Saulnier, L�lio Renard Lavaud,	883
824	Rodriguez, Rafi Ayub, Raghotham Murthy, Raghu	Marie-Anne Lachaux, Pierre Stock, Teven Le Scao,	884
825	Nayani, Rahul Mitra, Rangaprabhu Parthasarathy,	Thibaut Lavril, Thomas Wang, Timoth�e Lacroix,	885
826	Raymond Li, Rebekkah Hogan, Robin Battey, Rocky	and William El Sayed. 2023. Mistral 7b . <i>Preprint</i> ,	886
827	Wang, Russ Howes, Ruty Rinott, Sachin Mehta,	arXiv:2310.06825.	887
828	Sachin Siby, Sai Jayesh Bondu, Samyak Datta, Sara		
829	Chugh, Sara Hunt, Sargun Dhillon, Sasha Sidorov,	Di Jin, Luzhi Wang, Yizhen Zheng, Guojie Song, Fei	888
830	Satadru Pan, Saurabh Mahajan, Saurabh Verma,	Jiang, Xiang Li, Wei Lin, and Shirui Pan. 2023a.	889
831	Seiji Yamamoto, Sharadh Ramaswamy, Shaun Lind-	Dual intent enhanced graph neural network for	890
832	say, Shaun Lindsay, Sheng Feng, Shenghao Lin,	session-based new item recommendation . In <i>Pro-</i>	891
833	Shengxin Cindy Zha, Shishir Patil, Shiva Shankar,	<i>ceedings of the ACM Web Conference 2023, WWW</i>	892
834	Shuqiang Zhang, Shuqiang Zhang, Sinong Wang,	'23, page 684–693, New York, NY, USA. Association	893
835	Sneha Agarwal, Soji Sajuyigbe, Soumith Chintala,	for Computing Machinery.	894
836	Stephanie Max, Stephen Chen, Steve Kehoe, Steve		
837	Satterfield, Sudarshan Govindaprasad, Sumit Gupta,	Wei Jin, Haitao Mao, Zheng Li, Haoming Jiang,	895
838	Summer Deng, Sungmin Cho, Sunny Virk, Suraj	Chen Luo, Hongzhi Wen, Haoyu Han, Hanqing Lu,	896
839	Subramanian, Sy Choudhury, Sydney Goldman, Tal	Zhengyang Wang, Ruirui Li, Zhen Li, Monica Xiao	897
840	Remez, Tamar Glaser, Tamara Best, Thilo Koehler,	Cheng, Rahul Goutam, Haiyang Zhang, Karthik Sub-	898
841	Thomas Robinson, Tianhe Li, Tianjun Zhang, Tim	bian, Suhang Wang, Yizhou Sun, Jiliang Tang, Bing	899
842	Matthews, Timothy Chou, Tzook Shaked, Varun	Yin, and Xianfeng Tang. 2023b. Amazon-m2: A	900
843	Vontimitta, Victoria Ajayi, Victoria Montanez, Vijai	multilingual multi-locale shopping session dataset for	901
844	Mohan, Vinay Satish Kumar, Vishal Mangla, Vlad	recommendation and text generation . In <i>Advances in</i>	902
845	Ionescu, Vlad Poenaru, Vlad Tiberiu Mihailescu,	<i>Neural Information Processing Systems 36: Annual</i>	903
846	Vladimir Ivanov, Wei Li, Wenchen Wang, Wen-	<i>Conference on Neural Information Processing Sys-</i>	904
847	wen Jiang, Wes Bouaziz, Will Constable, Xiaocheng	<i>tems 2023, NeurIPS 2023, New Orleans, LA, USA,</i>	905
848	Tang, Xiaojian Wu, Xiaolan Wang, Xilun Wu, Xinbo	<i>December 10 - 16, 2023</i> .	906
849	Gao, Yaniv Kleinman, Yanjun Chen, Ye Hu, Ye Jia,		
850	Ye Qi, Yenda Li, Yilin Zhang, Ying Zhang, Yossi Adi,	Lei Li, Yongfeng Zhang, and Li Chen. 2020. Gener-	907
851	Youngjin Nam, Yu, Wang, Yu Zhao, Yuchen Hao,	ate neural template explanations for recommendation .	908
852	Yundi Qian, Yunlu Li, Yuzi He, Zach Rait, Zachary	In <i>Proceedings of the 29th ACM International Con-</i>	909
853	DeVito, Zef Rosnbrick, Zhaoduo Wen, Zhenyu Yang,	<i>ference on Information & Knowledge Management,</i>	910
854	Zhiwei Zhao, and Zhiyu Ma. 2024. The llama 3 herd	<i>CIKM '20</i> , page 755–764, New York, NY, USA. As-	911
855	of models . <i>Preprint</i> , arXiv:2407.21783.	sociation for Computing Machinery.	912
856	Yupeng Hou, Jiacheng Li, Zhankui He, An Yan, Xiusi	Haotian Liu, Chunyuan Li, Qingyang Wu, and Yong Jae	913
857	Chen, and Julian McAuley. 2024. Bridging lan-	Lee. 2023a. Visual instruction tuning . <i>Preprint</i> ,	914
858	guage and items for retrieval and recommendation .	arXiv:2304.08485.	915
859	<i>Preprint</i> , arXiv:2403.03952.		
860	Derek Hao Hu, Qiang Yang, and Ying Li. 2008. An	Xin Liu, Zheng Li, Yifan Gao, Jingfeng Yang, Tianyu	916
861	algorithm for analyzing personalized online commer-	Cao, Zhengyang Wang, Bing Yin, and Yangqiu	917
862	cial intention . In <i>Proceedings of the 2nd Interna-</i>	Song. 2023b. Enhancing user intent capture in	918
863	<i>tional Workshop on Data Mining and Audience In-</i>	session-based recommendation with attribute pat-	919
864	<i>telligence for Advertising</i> , ADKDD '08, page 27–36,	terns . <i>Preprint</i> , arXiv:2312.16199.	920
865	New York, NY, USA. Association for Computing		
866	Machinery.	OpenAI. 2024a. Gpt-4o mini: advancing cost-efficient	921
		intelligence . <i>OpenAI</i> .	922
867	Ravi Chandra Jammalamadaka, Naren Chittar, and San-	OpenAI. 2024b. Introducing openai o1 . <i>OpenAI</i> .	923
868	jay Ghatore. 2009. Mining product intention rules		
869	from transaction logs of an ecommerce portal . In	OpenAI. 2025. Openai o3-mini . <i>OpenAI</i> .	924
870	<i>Proceedings of the 2009 International Database En-</i>		
871	<i>gineering & Applications Symposium</i> , IDEAS '09,	OpenAI, Josh Achiam, Steven Adler, Sandhini Agarwal,	925
872	page 311–314, New York, NY, USA. Association for	Lama Ahmad, Ilge Akkaya, Florencia Leoni Ale-	926
873	Computing Machinery.	man, Diogo Almeida, Janko Altschmidt, Sam Alt-	927
		man, Shyamal Anadkat, Red Avila, Igor Babuschkin,	928
874	Bohan Jia, Jian Cao, Shiyu Qian, Nengjun Zhu, Xin		

929	Suchir Balaji, Valerie Balcom, Paul Baltescu, Haiming Bao, Mohammad Bavarian, Jeff Belgum, Irwan Bello, Jake Berdine, Gabriel Bernadett-Shapiro, Christopher Berner, Lenny Bogdonoff, Oleg Boiko, Madelaine Boyd, Anna-Luisa Brakman, Greg Brockman, Tim Brooks, Miles Brundage, Kevin Button, Trevor Cai, Rosie Campbell, Andrew Cann, Brittany Carey, Chelsea Carlson, Rory Carmichael, Brooke Chan, Che Chang, Fotis Chantzis, Derek Chen, Sully Chen, Ruby Chen, Jason Chen, Mark Chen, Ben Chess, Chester Cho, Casey Chu, Hyung Won Chung, Dave Cummings, Jeremiah Currier, Yunxing Dai, Cory Decareaux, Thomas Degry, Noah Deutsch, Damien Deville, Arka Dhar, David Dohan, Steve Dowling, Sheila Dunning, Adrien Ecoffet, Atty Eleti, Tyna Eloundou, David Farhi, Liam Fedus, Niko Felix, Simón Posada Fishman, Juston Forte, Isabella Fulford, Leo Gao, Elie Georges, Christian Gibson, Vik Goel, Tarun Gogineni, Gabriel Goh, Rapha Gontijo-Lopes, Jonathan Gordon, Morgan Grafstein, Scott Gray, Ryan Greene, Joshua Gross, Shixiang Shane Gu, Yufei Guo, Chris Hallacy, Jesse Han, Jeff Harris, Yuchen He, Mike Heaton, Johannes Heidecke, Chris Hesse, Alan Hickey, Wade Hickey, Peter Hoeschele, Brandon Houghton, Kenny Hsu, Shengli Hu, Xin Hu, Joost Huizinga, Shantanu Jain, Shawn Jain, Joanne Jang, Angela Jiang, Roger Jiang, Haozhun Jin, Denny Jin, Shino Jomoto, Billie Jonn, Heewoo Jun, Tomer Kaftan, Łukasz Kaiser, Ali Kamali, Ingmar Kanitscheider, Nitish Shirish Keskar, Tabarak Khan, Logan Kilpatrick, Jong Wook Kim, Christina Kim, Yongjik Kim, Jan Hendrik Kirchner, Jamie Kiros, Matt Knight, Daniel Kokotajlo, Łukasz Kondraciuk, Andrew Kondrich, Aris Konstantinidis, Kyle Kosic, Gretchen Krueger, Vishal Kuo, Michael Lampe, Ikai Lan, Teddy Lee, Jan Leike, Jade Leung, Daniel Levy, Chak Ming Li, Rachel Lim, Molly Lin, Stephanie Lin, Mateusz Litwin, Theresa Lopez, Ryan Lowe, Patricia Lue, Anna Makanju, Kim Malfacini, Sam Manning, Todor Markov, Yaniv Markovski, Bianca Martin, Katie Mayer, Andrew Mayne, Bob McGrew, Scott Mayer McKinney, Christine McLeavey, Paul McMillan, Jake McNeil, David Medina, Aalok Mehta, Jacob Menick, Luke Metz, Andrey Mishchenko, Pamela Mishkin, Vinnie Monaco, Evan Morikawa, Daniel Mossing, Tong Mu, Mira Murati, Oleg Murk, David Mély, Ashvin Nair, Reiichiro Nakano, Rajeef Nayak, Arvind Neelakantan, Richard Ngo, Hyeonwoo Noh, Long Ouyang, Cullen O’Keefe, Jakub Pachocki, Alex Paino, Joe Palermo, Ashley Pantuliano, Giambattista Parascandolo, Joel Parish, Emy Parparita, Alex Passos, Mikhail Pavlov, Andrew Peng, Adam Peralman, Filipe de Avila Belbute Peres, Michael Petrov, Henrique Ponde de Oliveira Pinto, Michael, Pokorny, Michelle Pokrass, Vitchyr H. Pong, Tolly Powell, Alethea Power, Boris Power, Elizabeth Proehl, Raul Puri, Alec Radford, Jack Rae, Aditya Ramesh, Cameron Raymond, Francis Real, Kendra Rimbach, Carl Ross, Bob Rotsted, Henri Roussez, Nick Ryder, Mario Saltarelli, Ted Sanders, Shibani Santurkar, Girish Sastry, Heather Schmidt, David Schnurr, John Schulman, Daniel Selsam, Kyla Sheppard, Toki Sherbakov, Jessica Shieh, Sarah Shoker, Pranav	
	Shyam, Szymon Sidor, Eric Sigler, Maddie Simens, Jordan Sitkin, Katarina Slama, Ian Sohl, Benjamin Sokolowsky, Yang Song, Natalie Staudacher, Felipe Petroski Such, Natalie Summers, Ilya Sutskever, Jie Tang, Nikolas Tezak, Madeleine B. Thompson, Phil Tillet, Amin Tootoonchian, Elizabeth Tseng, Preston Tuggle, Nick Turley, Jerry Tworek, Juan Felipe Cerón Uribe, Andrea Vallone, Arun Vijayvergiya, Chelsea Voss, Carroll Wainwright, Justin Jay Wang, Alvin Wang, Ben Wang, Jonathan Ward, Jason Wei, CJ Weinmann, Akila Welihinda, Peter Welinder, Jiayi Weng, Lilian Weng, Matt Wiethoff, Dave Willner, Clemens Winter, Samuel Wolrich, Hannah Wong, Lauren Workman, Sherwin Wu, Jeff Wu, Michael Wu, Kai Xiao, Tao Xu, Sarah Yoo, Kevin Yu, Qiming Yuan, Wojciech Zaremba, Rowan Zellers, Chong Zhang, Marvin Zhang, Shengjia Zhao, Tianhao Zheng, Juntang Zhuang, William Zhuk, and Barret Zoph. 2024. Gpt-4 technical report . <i>Preprint</i> , arXiv:2303.08774.	993 994 995 996 997 998 999 1000 1001 1002 1003 1004 1005 1006 1007 1008 1009 1010 1011 1012
	Qwen, :, An Yang, Baosong Yang, Beichen Zhang, Binyuan Hui, Bo Zheng, Bowen Yu, Chengyuan Li, Dayiheng Liu, Fei Huang, Haoran Wei, Huan Lin, Jian Yang, Jianhong Tu, Jianwei Zhang, Jianxin Yang, Jiayi Yang, Jingren Zhou, Junyang Lin, Kai Dang, Keming Lu, Keqin Bao, Kexin Yang, Le Yu, Mei Li, Mingfeng Xue, Pei Zhang, Qin Zhu, Rui Men, Runji Lin, Tianhao Li, Tianyi Tang, Tingyu Xia, Xingzhang Ren, Xuancheng Ren, Yang Fan, Yang Su, Yichang Zhang, Yu Wan, Yuqiong Liu, Zeyu Cui, Zhenru Zhang, and Zihan Qiu. 2025. Qwen2.5 technical report . <i>Preprint</i> , arXiv:2412.15115.	1013 1014 1015 1016 1017 1018 1019 1020 1021 1022 1023 1024
	Soyeon Shim, Mary Ann Eastlick, Sherry L Lotz, and Patricia Warrington. 2001. An online prepurchase intentions model: The role of intention to search . <i>Journal of Retailing</i> , 77(3):397–416.	1025 1026 1027 1028
	Zhu Sun, Hongyang Liu, Xinghua Qu, Kaidong Feng, Yan Wang, and Yew Soon Ong. 2024. Large language models for intent-driven session recommendations . In <i>Proceedings of the 47th International ACM SIGIR Conference on Research and Development in Information Retrieval</i> , SIGIR ’24, page 324–334, New York, NY, USA. Association for Computing Machinery.	1029 1030 1031 1032 1033 1034 1035 1036
	Gemma Team, Morgane Riviere, Shreya Pathak, Pier Giuseppe Sessa, Cassidy Hardin, Surya Bhupatiraju, Léonard Hussenot, Thomas Mesnard, Bobak Shahriari, Alexandre Ramé, Johan Ferret, Peter Liu, Pouya Tafti, Abe Friesen, Michelle Casbon, Sabela Ramos, Ravin Kumar, Charline Le Lan, Sammy Jerome, Anton Tsitsulin, Nino Vieillard, Piotr Stanczyk, Sertan Girgin, Nikola Momchev, Matt Hoffman, Shantanu Thakoor, Jean-Bastien Grill, Behnam Neyshabur, Olivier Bachem, Alanna Walton, Aliaksei Severyn, Alicia Parrish, Aliya Ahmad, Allen Hutchison, Alvin Abdagic, Amanda Carl, Amy Shen, Andy Brock, Andy Coenen, Anthony Laforge, Antonia Paterson, Ben Bastian, Bilal Piot, Bo Wu, Brandon Royal, Charlie Chen, Chintu Kumar, Chris Perry, Chris Welty, Christopher A.	1037 1038 1039 1040 1041 1042 1043 1044 1045 1046 1047 1048 1049 1050 1051 1052

1173	Denny Zhou. 2023. Chain-of-thought prompting elicits reasoning in large language models . <i>Preprint</i> , arXiv:2201.11903.
1174	
1175	
1176	Baixuan Xu, Weiqi Wang, Haochen Shi, Wenxuan Ding, Huihao Jing, Tianqing Fang, Jiaxin Bai, Xin Liu, Changlong Yu, Zheng Li, Chen Luo, Qingyu Yin, Bing Yin, Long Chen, and Yangqiu Song. 2024. MIND: Multimodal shopping intention distillation from large vision-language models for E-commerce purchase understanding . In <i>Proceedings of the 2024 Conference on Empirical Methods in Natural Language Processing</i> , pages 7800–7815, Miami, Florida, USA. Association for Computational Linguistics.
1177	
1178	
1179	
1180	
1181	
1182	
1183	
1184	
1185	
1186	Changlong Yu, Weiqi 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 common-sense discovery . <i>Preprint</i> , arXiv:2211.08316.
1187	
1188	
1189	
1190	
1191	Chenwei Zhang, Wei Fan, Nan Du, and Philip S. Yu. 2016. Mining user intentions from medical queries: A neural network based heterogeneous jointly modeling approach . In <i>Proceedings of the 25th International Conference on World Wide Web, WWW '16</i> , page 1373–1384, Republic and Canton of Geneva, CHE. International World Wide Web Conferences Steering Committee.
1192	
1193	
1194	
1195	
1196	
1197	
1198	
1199	Qi Zhao, Yi Zhang, Daniel Friedman, and Fangfang Tan. 2015. E-commerce recommendation with personalized promotion . In <i>Proceedings of the 9th ACM Conference on Recommender Systems, RecSys '15</i> , page 219–226, New York, NY, USA. Association for Computing Machinery.
1200	
1201	
1202	
1203	
1204	
1205	Yaowei Zheng, Richong Zhang, Junhao Zhang, Yanhan Ye, Zheyang Luo, Zhangchi Feng, and Yongqiang Ma. 2024. Llamafactory: Unified efficient fine-tuning of 100+ language models . <i>Preprint</i> , arXiv:2403.13372.
1206	
1207	
1208	
1209	Xi Zhu, Fake Lin, Ziwei Zhao, Tong Xu, Xiangyu Zhao, Zikai Yin, Xueying Li, and Enhong Chen. 2024. Multi-behavior recommendation with personalized directed acyclic behavior graphs . <i>ACM Trans. Inf. Syst.</i> , 43(1).
1210	
1211	
1212	
1213	

Appendices

1214
1215
1216
1217
1218
1219
1220
1221
1222
1223
1224
1225
1226
1227
1228
1229
1230
1231
1232
1233
1234
1235
1236
1237
1238
1239
1240
1241
1242
1243
1244
1245
1246
1247
1248
1249
1250
1251
1252
1253
1254
1255
1256
1257
1258
1259
1260
1261
1262

A Implementation Details

A.1 Attribute Extraction

To extract attributes from the given session products using GPT-4o-mini, we use the following 3-shot prompt as a general template:

Your goal is to extract the attribute type and attribute values of the product.

You will be provided with the product names and their corresponding product images, and you will output for the product:

Category: general category name of the product. Keep the category name simple and within 3 words.

Attributes: attribute(s) of the product. You can infer new ones from the image. Keep the attribute simple and within 3 words each. Separate different attributes by |. Generate in the format of attribute: value

Below are three examples:

...

Input:

Product Name: Adidas Ultraboost 21 Women's Running Shoes on sale, White/Pink special, Size 8 only, best for daily runs!

Output:

Category: Clothing

Attributes: brand: Adidas | model: Ultraboost 21 | gender: Women's | type: Running Shoes | color: White/Pink | size: 8

Input:

Product Name: Lightweight and powerful Dell XPS 13 Laptop, with newly released Intel i7, 16GB RAM, enhanced 512GB SSD, Silver version

Output:

Category: Electronics

Attributes: brand: Dell | model: XPS 13 | processor: Intel i7 | RAM: 16GB | storage: 512GB | color: Silver

Input:

Product Name: baking enthusiasts' good friend - KitchenAid Artisan Series 5-Quart Stand Mixer, Empire Red

1263	Output:	loafer, target consumers: men, size:	1313
1264	Category: Kitchen Appliance	3.5, price: \$200, structure: cushioned	1314
1265	Attributes: brand: KitchenAid model:	insole).	1315
1266	Artisan Series capacity: 5-Quart	Output:	1316
1267	type: Stand Mixer color: Empire Red	New Intention: Invest in premium quality	1317
1268	...	footwear for long-lasting style and	1318
1269	Input:	comfort.	1319
1270	<INPUT MESSAGE>	Attribute: design: luxury material and	1320
1271	Output:	craftsmanship.	1321
1272		Rationale: durability: The LV Glove	1322
1273	A.2 Intention Tree Construction	Loafer is crafted from high-quality	1323
1274	To construct the intention tree by filling in the nec-	materials, offering durability and	1324
1275	essary intention entries, pivotal attributes and un-	style that ensures it will last longer	1325
1276	derlying comparison, we use the following 5-shots	than ordinary shoes.	1326
1277	template:	Comparison: collectability: compared to	1327
1278	Act as a customer who is browsing a	the Nike Free Metcon 5, which focuses	1328
1279	series of products.	on performance, the LV offers a blend	1329
1280	For each input, you are required to	of luxury and longevity, making it a	1330
1281	generate several intentions as output,	worthy investment.	1331
1282	and each intention should only contain	New Intention: Own a versatile pair	1332
1283	the following lines of information:	of shoes suitable for both casual and	1333
1284	New Intention: new intention you may	formal settings.	1334
1285	have after interacting with the new	Attribute: varieties: loafer.	1335
1286	product	Rationale: usages: The loafer style of	1336
1287	Attribute: attribute(s) of the new	the LV Glove Loafer makes it versatile	1337
1288	product that caused the change in	enough to be worn in both casual	1338
1289	intention. You can infer new ones from	and formal settings, unlike the more	1339
1290	the image. Generate in the format of	specialized athletic design of the Nike	1340
1291	attribute: value	Free Metcon 5.	1341
1292	Rationale: a short rationale explaining	Comparison: versatility: While the Nike	1342
1293	why the attribute of the new product	Free Metcon 5 is primarily designed for	1343
1294	reflects the new intention. Generate in	workouts, the LV Glove Loafer's loafer	1344
1295	the format of facets: reasoning	style offers versatility for various	1345
1296	Comparison: a comparison between the	occasions.	1346
1297	new product and the previous product to	New Intention: Enhance your wardrobe	1347
1298	justify why the new product caused the	with a statement piece that reflects	1348
1299	change in intention. Generate in the	personal style.	1349
1300	format of aspects: comparison	Attribute: design: unique and luxurious.	1350
1301	Here is one example with five	Rationale: aesthetics: The unique and	1351
1302	intentions:	luxurious design of the LV Glove Loafer	1352
1303	...	makes it a statement piece that can	1353
1304	Input:	elevate any outfit, reflecting personal	1354
1305	Previous Intention: Looking for stylish	style.	1355
1306	and modern footwear that complements	Comparison: uniqueness: Unlike the more	1356
1307	their athletic look.	common athletic design of the Nike Free	1357
1308	Previous Product: Nike Free Metcon	Metcon 5, the LV Glove Loafer stands	1358
1309	5 Women's Workout Shoes (varieties:	out as a unique and stylish addition to	1359
1310	runner, target consumers: women, size:	the wardrobe.	1360
1311	3.5, price: \$100).	New Intention: Prioritize comfort	1361
1312	New Product: LV Glove Loafer (varieties:	without compromising on style.	1362
		Attribute: comfort: cushioned insole.	1363

1364 Rationale: comfort: The cushioned
1365 insole of the LV Glove Loafer ensures
1366 comfort for long periods, making it a
1367 practical choice without compromising
1368 on style.
1369 Comparison: comfort: While the Nike
1370 Free Metcon 5 is designed for athletic
1371 performance, the LV Glove Loafer offers
1372 a balance of comfort and style for
1373 everyday wear.
1374 New Intention: Choose a high-end brand
1375 to reflect social status.
1376 Attribute: brand: Louis Vuitton.
1377 Rationale: status: Owning a product
1378 from a high-end brand like Louis Vuitton
1379 reflects social status and prestige.
1380 Comparison: brand prestige: Compared to
1381 Nike, which is known for athletic wear,
1382 Louis Vuitton is a luxury brand that
1383 signifies higher social status.
1384 ...
1385 Input:
1386 Previous Intention: <Previous
1387 Intention>
1388 Previous Product: <PREVIOUS PRODUCTS>
1389 New Product: <THE LAST PRODUCT>
1390 Output:
1391

1392 For smaller branching sizes, we just need to
1393 delete some of the examples provided. And for
1394 extending to more branches every step, we add
1395 additional examples as needed.

1396 A.3 Intention Generator Model Selection

1397 We first tried out free open LVLM models like
1398 Mantis and LLaVA families. However, the mod-
1399 els fail to achieve the desired outcome since most
1400 of them cannot output in pre-assigned formatting.
1401 This is possible because models are not able to han-
1402 dle the potentially long, complex product textual
1403 descriptions and attributes provided. For example,
1404 repeatedly generating a single word or outputting a
1405 large number of special symbols like "####." After
1406 careful examination, it is not caused by prompts
1407 and product information included. We switched to
1408 GPT-4o-mini later and found the generated inten-
1409 tion and metadata result is in the desired format,
1410 and it demonstrates comparable results with GPT-
1411 4o. Therefore, we opt for GPT-4o-mini as the major
1412 generating force for the intention tree construction
1413 part.

A.4 Fintuning Model Selection 1414

1415 We selected models for fine-tuning based on the fol-
1416 lowing criteria to ensure diversity and compatibility
1417 with our experimental setup: (1) *Organizational Di-*
1418 *versity*: Models were chosen from different organi-
1419 zations to ensure a broad representation of architec-
1420 tures and training methodologies. (2) *Model Size*:
1421 Models were required to have fewer than 11 billion
1422 parameters to be trainable on our hardware. (3)
1423 *Performance*: We selected top-performing models
1424 from their respective organizations. For example,
1425 we chose Llama-3.1-8B from Meta, as it is among
1426 the best 7B/8B models released. Later versions,
1427 such as Llama-3.2, focus on smaller models (e.g.,
1428 Llama-3.2-1B and Llama-3.2-3B), while Llama-
1429 3.3 targets larger models (e.g., Llama-3.3-70B),
1430 making Llama-3.1-8B the optimal choice for our
1431 study.

A.5 Training-Test Splits 1432

1433 The detailed process is outlined as follows: (1)
1434 *Indexing*: We created an index for all annotated
1435 questions. Each questionnaire included four ques-
1436 tions corresponding to Tasks 1–4, which ensured
1437 an equal number of samples per task. Therefore,
1438 the proportion of each question’s indices are equal.
1439 (2) *Index Set Creation*: A unified set of indices
1440 was established, with each index uniquely corre-
1441 sponding to a specific Task and Session number for
1442 traceability. (3) *Splitting*: We adopted a 4:1 train-
1443 test split ratio. Indices were randomly sampled to
1444 create training and test sets. While the number of
1445 samples for Tasks 1–4 may vary slightly due to
1446 random sampling, the counts remain largely bal-
1447 anced. (4) The resulting training and test sets were
1448 used across different models and training schemes
1449 (e.g., zero-shot, fine-tuning with SIB, or sequential
1450 fine-tuning with MIND followed by SIB).

A.6 Few-shot Example Curation 1451

1452 When selecting examples for prompt templates, our
1453 primary criterion is clarity, ensuring the examples
1454 are easily understandable by large language models
1455 (LLMs). Additionally, examples must be concise,
1456 free of conflicts or ambiguities, and possess a rela-
1457 tively clear answer from a human perspective.

1458 Given that small, finite discrete point masses can-
1459 not fully approximate the continuous product and
1460 intention distribution space, the chosen examples
1461 are not intended to be representative of the *entire*
1462 product distribution.

1463	We initially generated a large set of examples	1512
1464	randomly using GPT-4o (8 examples across 5 at-	1513
1465	tempts of different categories and features, total-	1514
1466	ing 40 examples). Our researchers then manually	1515
1467	selected, refined, and annotated these examples,	1516
1468	followed by a thorough validation process.	1517
1469	Impact of Prompt Variations: To identify the	1518
1470	optimal prompt, we experimented with multiple	1519
1471	templates by (1) varying example categories (e.g.,	1520
1472	switching from clothing to electronic devices) and	1521
1473	(2) increasing the number of examples. Testing on	
1474	the GPT-4o API showed that these modifications	
1475	resulted in minimal accuracy fluctuations (within	
1476	1%).	
1477	A.7 Model Evaluation	
1478	We evaluate the model on the tasks using differ-	1523
1479	ent prompt techniques including zero-shot prompts	1524
1480	(see Table 8), 5-shots prompts (see Table 9, 10, 11,	1525
1481	12) and Chain-of-Thought prompts (see Table 13).	1526
1482	The detailed prompts we used can be found in the	1527
1483	corresponding tables.	1528
1484	Note that for consistency, the “ <i>Few-shots</i> ” evalu-	1529
1485	ation mentioned in the Main Table 2 uses the same	
1486	prompt as the 5-shot evaluation prompt.	
1487	A.8 Finetuning Methods	
1488	We employed Supervised Fine-Tuning (SFT) with	1530
1489	Low-Rank Adaptation (LoRA) for all fine-tuning	1531
1490	experiments. To ensure reproducibility and ease	1532
1491	of implementation, we utilized the open-source	1533
1492	LLaMA-Factory framework (Zheng et al., 2024).	1534
1493	Upon acceptance, we will release the complete	1535
1494	implementation code and configuration files for	1536
1495	the data generation pipeline, evaluation, and fine-	1537
1496	tuning stages. Additionally, we will revise the rele-	1538
1497	vant section with clear descriptions for improved	1539
1498	transparency.	1540
1499	A.9 Role of GPT-4o-mini	1541
1500	We utilized GPT-4o-mini in both the dataset gener-	1542
1501	ation pipeline and closed-source model evaluation,	
1502	as detailed in Appendix D.1. During dataset con-	
1503	struction, GPT-4o-mini outperformed many open-	
1504	source models due to its ability to process multiple	
1505	images while maintaining consistency with textual	
1506	instructions, a critical requirement given the large	
1507	scale of information processed. To ensure a robust	
1508	evaluation, we also included GPT-4o as a base-	
1509	line and compared various prompting techniques	
1510	(zero-shot, 5-shot, and Chain-of-Thought) across	
1511	models.	
	Notably, despite GPT-4o-mini’s role in dataset	1512
	generation, it did not exhibit significantly superior	1513
	performance on our dataset’s question-answering	1514
	tasks, highlighting the dataset’s effectiveness in	1515
	challenging current models. However, GPT-4o-	1516
	mini did outperform GPT-4o on several metrics,	1517
	validating its inclusion in the evaluation. To avoid	1518
	bias in error analysis (Section 5.5), we used GPT-	1519
	4o’s Chain-of-Thought responses instead of GPT-	1520
	4o-mini’s.	1521
	A.10 Expert Selection for Error Analysis	1522
	The experts were three Computer Science PhD stu-	1523
	dents from our institution’s research community,	1524
	each with at least five publications in natural lan-	1525
	guage processing (NLP) or related fields, ensur-	1526
	ing their expertise in the domain. Selecting three	1527
	experts helped minimize bias and enhance consis-	1528
	tency in the error analysis.	1529
	For the error analysis, we randomly sampled	1530
	200 questions from Tasks 1–4 where the model	1531
	(GPT-4o with Chain-of-Thought) provided incor-	1532
	rect answers. These sampled questions were an-	1533
	alyzed collectively to identify patterns in model	1534
	errors across all tasks, rather than focusing on a sin-	1535
	gle task. Each expert independently reviewed the	1536
	sampled questions and assigned error labels based	1537
	on predefined criteria, such as misinterpretation of	1538
	user intent, failure to utilize session context, or log-	1539
	ical inconsistencies in reasoning. The final error	1540
	labels were determined through a consensus pro-	1541
	cess among the three experts to ensure reliability.	1542
	B Theoretical Framework	1543
	B.1 Intention Tree	1544
	The Intention Tree \mathbf{T} is defined inductively. At	1545
	each discrete time step $t \in \mathbb{N}^+$, we model the	1546
	branching process of the tree by extending the ex-	1547
	isting $\mathbf{T}_{1,2,3,\dots,t-1}$ to $\mathbf{T}_{1,2,3,\dots,t}$. We denote the in-	1548
	formation stored in the intention nodes up to time	1549
	step t as \mathcal{H}_t , which represents the session interac-	1550
	tion history observed up to time t .	1551
	Note that traditional models typically infer in-	1552
	tentions based on direct transitions of products or	1553
	intentions, such as $\mathbb{P}(P_{t+1} P_t)$ or $\mathbb{P}(I_{t+1} I_t)$.	1554
	Our model instead adopts a two-step prediction	1555

process:

$$\begin{aligned} \mathbb{P}(P_{t+1}|\mathcal{H}_t) &= \mathbb{P}(\mathcal{M}_t, P_{t+1}|\mathcal{H}_t) \\ &= \mathbb{P}(\mathcal{M}_t|\mathcal{H}_t) \cdot \mathbb{P}(P_{t+1}|\mathcal{H}_t, \mathcal{M}_t) \\ &\approx \mathbb{P}(\mathcal{M}_t^\phi|\mathcal{H}_t) \cdot \mathbb{P}(P_{t+1}|\mathcal{H}_t, \mathcal{M}_t^\phi) \end{aligned} \quad (1)$$

Here, \mathcal{M}_t^ϕ is the model’s approximation of the session-level information \mathcal{M}_t at time t .

Rather than modeling the rich information in I_t as a monolithic whole, we explicitly separate out the key components. This gives rise to the inferred elements:

- I_t^ζ : the extracted intention,
- C_t^ζ : relevant comparisons, and
- A_t^ζ : associated attributes,

which together serve as an approximation of the full information space at time step t :

$$\mathcal{M}_t^\phi = (A_t^\zeta, I_t^\zeta, C_t^\zeta)$$

From the model’s perspective, the inference process is thus simplified to:

$$\mathbb{P}(P_{t+1}|\mathcal{H}_t, I_t^\zeta, C_t^\zeta, A_t^\zeta)$$

Here, the components $(A_t^\zeta, I_t^\zeta, C_t^\zeta)$ are inferred from (P_{t+1}, \mathcal{H}_t) , with the superscript ζ indicating a branched approximation of consumer intentions, generated by GPT-4o-mini.

As an example of how this formulation enables targeted modeling, consider the task of Valued Attribute Regularization. We first sample a subset of attribute profiles \mathcal{A}^π from the full attribute space \mathcal{A} . The model is tasked with distinguishing which attributes are of crucial influence based on the provided profiles. This task is represented as a cascade of inference steps:

$$P_{t+1}, \mathcal{H}_t | \mathcal{A}^\pi \Rightarrow I_t^\zeta | P_{t+1}, \mathcal{H}_t, \mathcal{A}^\pi \Rightarrow A_t^\zeta | \mathcal{A}^\pi, I_t^\zeta$$

constraint to $A_t^\zeta \in \mathcal{A}^\pi$. with the constraint that $A_t^\zeta \in \mathcal{A}^\pi$.

B.2 Intuition

Consider a session with two products and a 5-branching scheme. One intention is generated for the first product, while five intentions are generated for the second product. These five intentions for the second product serve as branches stemming from the intention of the first product.

When a third product is added to the session, each intention associated with the second product becomes a new initial intention, branching out into five additional intentions linked to the third product. This process continues iteratively for subsequent products.

Each intention is uniquely associated with a set of metadata (e.g., comparison, justification, attributes) in a one-to-one correspondence.

C Task and Evaluation Design

C.1 Design Criteria for Choice Options

The scores ranging from 0 to 3 in the task definition serve as symbolic representations of specific answer formulations. Concrete examples illustrating the implementation of these choices can be found in the tables on the following pages.

For example, for the first task, the score represents the likelihood estimation $\mathcal{S}_1(P_t, I_{t-1}) |_{\mathcal{H}_{t-1}} \in \{0, 1, 2, 3\}$, which indicates the probability of a customer interacting with P_t . A score of 3 signifies the highest likelihood, while a score of 0 represents the lowest. Specifically, as shown in Table 8 of our paper:

- Score = 3 corresponds to *A. Yes: The product is a logical and reasonable outcome of the purchasing intention.*
- Score = 2 corresponds to *B. Maybe Yes: I may consider this, but it’s not a strong impulse.*
- Score = 1 corresponds to *C. Maybe No: The product is not directly related to my intention.*
- Score = 0 corresponds to *D. No: I would never purchase it if I were the customer with the given intention.*

We hope this explanation resolves any doubts regarding the formulation of the scores. For the ground truth label, we combine “Maybe Yes” (B) with “Yes” (A) and “Maybe No” (C) with “No” (D) to create a clearer decision boundary. This approach aligns with common practices in survey analysis and simplifies the classification problem into two categories. The design of the ground truth label *accounts for neutral response bias to better reflect the true distribution of opinions*. Otherwise, the overrepresentation of neutral responses can lead to unexpected model bias. By grouping the responses, the problem is reduced to a binary

1641
1642
1643
1644
1645
1646
1647
1648
1649
1650
1651
1652
1653
1654
1655

1656
1657

1658
1659
1660
1661
1662
1663
1664
1665
1666
1667
1668
1669
1670
1671
1672
1673
1674
1675
1676
1677
1678
1679
1680
1681
1682
1683
1684
1685
1686
1687
1688
1689
1690
1691

classification setting (e.g., approve/reject or positive/negative). This not only simplifies analysis but also makes the results more actionable for the model training and fine-tuning.

For Task 4, we designate a score of 0 as the true positive label, as it represents “firm preferences for products within the same category and feature set.” The other two labels are grouped together, as they reflect the customer’s inclination to explore products from different features or categories. This results in a binary label design for Task 4, distinct from the label structures used in Tasks 1–3.

For further details, please kindly refer to Table 8, where we outline the three-choice design for Task 4 in contrast to Tasks 1–3.

C.2 Why Not Use Open-Ended Question Answering

We considered using open-ended question answering but identified several practical and theoretical challenges that led us to adopt our current framework, which we believe is more efficient, scalable, and cost-effective. These challenges include: (1) *Standardization and Scalability*: Allowing annotators to propose intentions would result in highly diverse responses, making it difficult to standardize and cluster them for consistent model evaluation. Our framework prioritizes a unified, scalable pipeline that facilitates easy adoption and production use. Open-ended intent proposals would complicate dataset standardization and increase the time required for annotation, hindering scalability. (2) *Cost Efficiency*: Proposing intentions requires annotators to generate detailed, meaningful content, which demands higher expertise and increases labor costs. Given the large number of questions in our dataset, relying on a small group of highly qualified annotators (e.g., PhD or postdoc students) is impractical. Our current approach, using structured multiple-choice questions, reduces the cognitive and financial burden of annotation. (3) *Annotator Reliability*: Generating reliable intent proposals requires a high level of responsibility and task understanding. Our Worker Selection Protocol (see Appendix D.1) and annotation filtering process (Appendix E.1) revealed that, even with structured tasks, many of the 300 initial candidates provided arbitrary labels, with only 11 passing our rigorous selection criteria. Allowing open-ended intent proposals would likely exacerbate reliability issues, raising concerns about the trustworthiness of the annotations.

Given these considerations, we believe our current methodology strikes an effective balance for practical implementation.

D Annotation Process

D.1 Worker Selection Protocol

We carry out strict quality control to ensure high quality human annotation result. To start with, we send qualification round invitations only to workers who satisfy the following constraints: (i) pass over 2,000 HITs, (ii) score over 90% on historical approval rate. We curate a qualification test with sampled sessions, whose gold label are provided by the authors. We first retain those with an accuracy of over 75% on the qualification test and complete over 20 questions at the same time.

Following that, we disqualify those spammers or underperforming workers. More specifically, we filter out those workers who are simply picking one side of the choices for the majority of the time. After conducting another round of testing, results show that least 7 people picking one side of the answer 80 percent of the time, so we eliminate out those people and proceed to main round annotations.

11 workers left after the selection process out of 300 initial candidates. This gives a worker selection rate of 3.67%.

D.2 Annotation Instructions

We give instruction to workers in layman’s terms, with both detailed question definitions and specific explanations for important information included. We tried to include more information and less distractions. The question definitions are closely aligned to what we defined earlier in Section 3. For each of the first three questions, workers are asked to annotate using a four-point 0 to 3 likelihood scale, where 0 stands for the least probable or the least plausible, and 3 means the most likely or the most plausible. For the forth question, annotation results are constraint to a three point scale, from 1 to 3. Larger the number, larger the likelihood of exploring more diverse products.

E Annotation Result Analysis

E.1 Raw Label Result

The distribution of labels is reflected in the *Majority* score presented in the Table 2. We have also summarized the raw label distribution below. Instead of reporting individual choices (e.g., A or

1692
1693
1694

1695

1696

1697
1698
1699
1700
1701
1702
1703
1704
1705
1706
1707
1708
1709
1710
1711
1712
1713
1714
1715
1716
1717
1718

1719

1720
1721
1722
1723
1724
1725
1726
1727
1728
1729
1730
1731
1732
1733

1734

1735
1736
1737
1738
1739

B), we grouped them to mitigate individual annotator biases observed during the annotation process, where some annotators consistently favored extreme answers while others preferred intermediate ones.

Task_Ind	Label	Count	Percentage
1	A_B	5844	62.30%
1	C_D	3536	37.70%
2	A_B	4282	45.65%
2	C_D	5098	54.35%
3	A_B	6735	71.80%
3	C_D	2645	28.20%
4	A	3456	36.84%
4	B_C	5924	63.16%

Table 4: Summarized counts and percentages of question labels. *Task_Ind* denotes the task index, ranging from 1 to 4. *A_B* indicates whether option *A* or *B* was selected, and similarly, *C_D* represents whether option *C* or *D* was chosen.

E.2 Consistency

We employ a majority vote rule to establish ground truth. In Table 5, the notations “2:1” and “3:0” indicate the label distribution for each question. A “3:0” label means all three annotators selected the same answer, while “2:1” indicates that one annotator’s choice differed from the other two (e.g., annotators 1 and 2 selected *A* or *B*, while annotator 3 selected *C* or *D*). This distribution shows that for over 50% of the questions, human annotators consistently agreed on the answer. Thus, these results further suggest that current models lag behind human performance in understanding these tasks.

Task_Ind	2:1	3:0
1	6041	7959
2	9170	4830
3	5390	8610
4	3934	10066

Table 5: Consistency analysis of binary answer label distribution. *Task_Ind* denotes the task index, ranging from 1 to 4.

E.3 Annotation Quality Filter

In addition to dataset-level analysis, we conducted a detailed examination of individual annotator statistics. For instance, during the annotator filtering process, we compiled Table 6 to summarize individual label distribution.

During the annotator filtering process, we identified individuals who consistently favored one set of choices (e.g., always selecting *A/B* or *C/D*). Such behavior may indicate a lack of engagement, where annotators select options indiscriminately to complete the task and receive payment. To ensure quality, we excluded these annotators in the official annotation round by administering a separate qualification test and filtering out those exhibiting this pattern.

Annotator_ID	A	B	C	D
A1***1A	3201	719	893	31
A2***EZ	3402	1208	437	1
A2***2M	106	5540	1950	48
A1***SU	633	186	135	18
A3***TX	2113	173	610	28
A2***BO	287	32	49	0
A2***YO	919	221	196	24
A2***E0	466	129	140	5
AF***9P	60	23	14	3

Table 6: Exclude annotators who consistently select the same option the majority of the time.

E.4 Benchmark and Data Quality Validation

As noted previously, incorrect labels may arise from objective factors, such as internal conflicts within session products and metadata in the intention tree or the inclusion of complex metadata in the labeling instructions. **These issues stem from customer behaviors**, which can exhibit self-conflicting patterns or random product-to-product jumps, introducing inherent randomness that is challenging to eliminate. However, one can expect that preprocessing to identify these patterns, though effort-intensive, can mitigate such issues and enhance benchmark quality.

E.5 Clarifications on the Majority Vote Score

The majority voting we defined in the Table 2 means that, for each task, after assigning ground-truth labels to all questions in the dataset, we calculate the most frequently selected option (i.e., with highest label frequencies) across all questions for that task and determine its corresponding score. For instance, in Task 3, options *A/B* indicate that the provided context offers valid justifications for user behavior, while options *C/D* indicate insufficient justifications. The correct option is different for each specific question, even the questions are from the same Task 3, as it is not logical to assume a

universal “correct” response (e.g., always choosing A/B or C/D) that consistently outperforms others on that Task 3. *What we can only say is that in our dataset, the correct answer is skewed towards the A/B side.*

Thus, the majority vote serves as a baseline for comparison with model performance, **not as a measure of human performance or data quality.**

E.6 Missing F1 Score for Tasks 2 and 4

The NaN F1 scores for Tasks 2 and 4 arise from the F1 score calculation process, which is defined as follows:

1. Precision = $\frac{TP}{TP + FP}$
2. Recall = $\frac{TP}{TP + FN}$
3. F1 = $\frac{2 \cdot \text{Precision} \cdot \text{Recall}}{\text{Precision} + \text{Recall}}$

Here, TP, FP, TN, and FN represent true positives, false positives, true negatives, and false negatives, respectively. For Tasks 1–3, we define options A/B as positive (with options C/D as negative), and for Task 4, option A is positive (with options B/C as negative). In Tasks 1 and 3, the majority vote favors A/B, resulting in $TP > 0$ and $FP > 0$, yielding valid Precision and Recall values. However, in Tasks 2 and 4, the majority vote favors negative options, leading to $TP = 0$ and $FP = 0$. This results in $TP + FP = 0$, rendering Precision undefined and the F1 score NaN.

F Model Performance Insights

F.1 Evaluation Task Performances Metrics

We display our summarized task performance metrics for each of the tasks in Table 7. Statistics in both counting and percentage format are included.

F.2 Finetuning

Our analysis suggests that models fine-tuned solely on the Session Intention Benchmark (SIB) may underperform due to the diverse and widely dispersed distribution of SIB data. This diversity creates a significant gap between the training and test sets, impacting model generalization (Wang et al., 2023b,a, 2024b; Wang and Song, 2025).

For instance, some models exhibit degraded performance when fine-tuned only on SIB for certain tasks. However, incorporating external intention injection (e.g., fine-tuning on SIB combined with

Metric	Task No.				
	TASK 1	TASK 2	TASK 3	TASK 4	
Count	# TP	781	415	527	57
	# FN	354	434	775	585
	# TN	234	515	309	949
	# FP	458	445	215	131
Percentage	TP (%)	42.75%	22.94%	28.86%	3.31%
	FN (%)	19.38%	24.00%	42.44%	33.97%
	TN (%)	12.81%	28.47%	16.92%	55.11%
	FP (%)	25.07%	24.60%	11.77%	7.61%

Table 7: Task performance metrics for error analyses of GPT-4o with Chain-of-Thought answering SESSIONINTENTBENCH. Where *TP*, *FN*, *TN*, *FP* stand for true positive, false negative, true negative, and false positive answers respectively.

MIND, an external intention knowledge base) improves outcomes. As an example, Llama-3.1-8B achieved 56.78% accuracy on Task 1 in a zero-shot setting, which dropped to 52.82% after fine-tuning on SIB alone but increased to 60.10% when fine-tuned on both SIB and MIND. This improvement likely results from the general intention knowledge provided by the external source, which mitigates distribution dispersion and enhances the model’s ability to generalize across the training-test gap.

F.3 Imbalanced Task Performance Gain with Intention Injection

As discussed in Section 5.4, a key factor affecting model performance is the ability to understand user intent, particularly for Task 1, which requires models to leverage intent information to infer future product preferences. Direct fine-tuning on SESSIONINTENTBENCH (SIB) led to a significant performance drop in Task 1, as models struggled to interpret intent effectively. However, sequential fine-tuning i.e., first on the MIND dataset followed by SIB, resulted in substantial performance improvements, supporting our hypothesis that pre-training on intent understanding enhances performance on intent-driven tasks like Task 1.

F.4 Solely Fine-Tuning with MIND

we have evaluated model performance after fine-tuning solely on the MIND dataset (i.e., without fine-tuning on SIB). When tested on SIB, *models fine-tuned only on MIND struggled to adhere to SIB’s question-answering format* (i.e., selecting specific options). Instead, models produced verbose reasoning without clear answers or embedded answers within the reasoning, resulting in uniformly low performance across models (<10%).

We attribute this to two factors: 1. Model Capacity: The models, limited to fewer than 11 billion parameters, struggled to adapt quickly to SIB’s distinct question-answering format during testing. 2. Dataset Differences: The MIND dataset contains open-ended questions designed to elicit detailed intent descriptions, which differ significantly from SIB’s structured, option-based format. This mismatch hindered the models’ ability to generalize to SIB’s requirements. *But that is also why we chose MIND.* We intentionally selected MIND for its focus on detailed intent descriptions to ensure models learn robust intent understanding, as exemplified in the MIND dataset samples. Here we appended some examples from the MIND dataset we used for finetuning, for your reference.

```
[ "instruction": "Act as a customer who
bought these two product: ", "input":
"Product A: Lincoln Stain Wax Shoe
Polish 3 Fl Oz (Selection of Colors);
Product B: Angelus Shoe Wax Polish
3fl Oz (Color Variety); What is your
possible co-buy intention for these two
products?", "output": "The potential
co-buy intention could be that the
person wants to purchase both shoe
polish products to have a variety of
colors to choose from when polishing
their shoes" , "instruction": "Act as a
customer who bought these two product: ",
"input": "Product A: BMC Mens 6 pc Mixed
Design Self Tie Bowtie Pocket Square
Suit Accessories; Product B: Tenby
Living 2-Pack Black Tie Rack, Organizer,
Hanger, Holder - Affordable Ti.; What
is your possible co-buy intention
for these two products?", "output":
"The potential co-buy intention for
people purchasing these two products
simultaneously could be to enhance their
wardrobe and maintain an organized and
stylish appearance" , ...]
```

F.5 The BERT-based Models

We evaluated pretrained BERT-based models like RoBERTa-large-355M and DeBERTa-v3-large, but found them unsuitable for testing in our context.

For RoBERTa-large-355M, the raw output is exemplified as follows:

```
["task_counter": 25248,
"session_counter": 6311,
```

```
"question_idx": 3, "response":
**A**, "task_counter": 27563,
"session_counter": 6890,
"question_idx": 2, "response":
**Yes**, "task_counter": 2654,
"session_counter": 663, "question_idx":
1, "response": **A**, "task_counter":
16969, "session_counter": 4242,
"question_idx": 0, "response":
**A**, "task_counter": 33507,
"session_counter": 8376,
"question_idx": 2, "response":
**Yes**, ...]
```

Upon analyzing the output, we observed that RoBERTa-large-355M consistently produces responses such as “A” or “Yes” regardless of the question, failing to align with the query’s requirements. This behavior is not observed in larger models. We hypothesize that smaller models like RoBERTa-large-355M struggle with complex tasks, often defaulting to predicting the most frequent answer in the masked space based on the nearest question description, without effectively processing the provided session product information, let alone metadata restrictions.

For DeBERTa-v3-large, the raw output is entirely consisted of random strings such as "IBILITY" and "Measurement" instead of providing answers to the questions:

```
["task_counter": 25248,
"session_counter": 6311,
"question_idx": 3, "response":
**IBILITY**, "task_counter":
27563, "session_counter": 6890,
"question_idx": 2, "response":
**IBILITY**, "task_counter":
2654, "session_counter": 663,
"question_idx": 1, "response":
**Measurement**, "task_counter":
16969, "session_counter": 4242,
"question_idx": 0, "response":
**IBILITY**, "task_counter":
33507, "session_counter": 8376,
"question_idx": 2, "response":
**IBILITY**, ...]
```

Despite experimenting with various prompting techniques, the model consistently failed to process the queries effectively. This suggests that DeBERTa-v3-large, due to its limited capacity,

1978 struggles with the complexity of these tasks. Sim-
1979 ilar issues were encountered when using Mantis
1980 family models as large vision-language models
1981 (LVLMs) for intention generation, where outputs
1982 often exhibited broken structures or degenerated
1983 into pure noise.

1984 Based on the observed limitations, we believe
1985 that models such as RoBERTa-large-355M and
1986 DeBERTa-v3-large should not be included in the
1987 current evaluation phase due to their inability to
1988 effectively process the complex tasks in our frame-
1989 work.

Survey Instructions (Click to Collapse)

How Intentions Evolve with Changing Attributes?

Welcome to our Main Round HITs. Congratulations on passing the qualification test and thanks for participating in our HITs!

In this survey, you will be provided a session of products and asked to evaluate alterations in purchasing intentions as the product attributes changes.

Before the questions: You will be provided with a list of Session Products that will be used throughout the questions.

Answer each question: Select the option that best describes your evaluation of the model's output based on the criteria provided.

Question Formalization

Q1: Changing Intentions

After reviewing the listed products (including their titles, attributes, images, etc.), and assuming you have the provided purchasing intention, we want to understand how likely you are to purchase a specific product based on this intention. Your task is to decide whether you would consider purchasing the product given your current intentions.

You'll be provided with four rating options: Yes, Maybe yes, Maybe no and No.

Q2: Attribute that Matters

After reviewing the listed products, and assuming you highly value a specific attribute of the listed products, we want to understand how likely you are to purchase another product based on this valued attribute. Your task is to decide whether you would consider purchasing the product given your focus on the specific characteristic.

You'll be provided with four rating options: Yes, Maybe yes, Maybe no and No.

Q3: Comparisons

After reviewing the listed products, and assuming you have the provided purchasing intentions, we want to understand if the comparison between the products provides a detailed and reasonable justification for your purchasing impulse. Your task is to decide whether the comparison is thorough enough to justify your change in intention.

You'll be provided with four rating options: Yes, Maybe yes, Maybe no and No.

Q4: Changing Desire

After reviewing the listed products, and assuming you have the provided purchasing intention, we want to understand if you still wish to explore similar products. Your task is to decide whether you want to continue exploring products within the same category or look for products in different categories.

You'll be provided with three rating options: Yes, Maybe yes, and No.

Session Products List

Session Products List is a list of products that you browsed (possibly consider purchasing) in a short period of time on Amazon.

The list of products will contain the following information:

- (1) **Product title:** The name of the product you viewed.
- (2) **New intention:** You should imagine yourself as a customer who has the mentioned intention/impulse when browsing the products. The word "New" means it's the intention you hypothetically have after seeing the last product in the current list.
- (3) **Attributes:** The features, functions, or characteristics of the product that you may consider when making a purchase decision. They are complementary information for the title/image to facilitate your decision process.

Each Session Products List is in one-to-one correspondence with the question following it.

Additional Hints

- Read the Session Products List carefully: Understand the previous intention, previous product, and new product details.
- Submit your response: Once you have answered all questions, click the Submit button to complete the HIT.

Figure 5: The annotation instruction we shown to workers, with detailed question definitions in layman's terms and specific explanations for important information (e.g., a preview of information contained in *Session Product List*).

Task	Zero-shot Prompt
TASK 1	<p>Act as a customer who is browsing a series of products given as follows. <session product information></p> <p>After seeing <previous products>, and assuming you are a customer who has the intention of <second last intention>. How likely are you to purchase <last product> based on the assumed intention?</p> <p>A. Yes: The product is a logical and reasonable outcome of the purchasing intention. B. Maybe yes: I may consider this, but it's not a strong impulse. C. Maybe no: The product is not directly related to my intention. D. No: I would never purchase it if I were the customer with the given intention.</p> <p>Your Answer (Answer A or B or C or D only):</p>
TASK 2	<p>Act as a customer who is browsing a series of products given as follows. <session product information></p> <p>After seeing <previous products>, and assuming you are a customer who highly value the feature <second last intention attribute> of <second last product>. How likely are you to purchase <last product>?</p> <p>A. Yes: The product logically and reasonably matches the characteristics I value. B. Maybe yes: I might consider this product, but it doesn't strongly appeal to me. C. Maybe no: The product does not directly relate to the characteristic I value. D. No: I would not purchase this product if I were focused on the given characteristic.</p> <p>Your Answer (Answer A or B or C or D only):</p>
TASK 3	<p>Act as a customer who is browsing a series of products given as follows. <session product information></p> <p>Comparing between <last two products>, and assuming you have the intention of <last two intention>, Does this comparison <last intention comparison> provide an in-depth justification of your impulse?</p> <p>A. Yes: the comparison is reasonable and detailed enough to justify the change. B. Maybe yes: The comparison could be more detailed and thorough but can be ignored. C. Maybe no: The comparison is not entirely reasonable or lacks sufficient in-depth detail. D. No: The comparison does not provide any underlying reasons or insights.</p> <p>Your Answer (Answer A or B or C or D only):</p>
TASK 4	<p>Act as a customer who is browsing a series of products given as follows. <session product information></p> <p>After seeing <previous products>, and assuming you have the intention of <previous intention>, do you still want to explore similar products?</p> <p>A. Yes: I want to explore products under the same category. B. Maybe yes: I want to explore products under the same category but with different features. C. No: I want to explore products under other categories.</p> <p>Your Answer (Answer A or B or C only):</p>

Table 8: Zero-shot prompts for model evaluation. TASK 1 stands for *Intent-Based Purchasing Likelihood Estimation*, TASK 2 stands for *Purchasing Likelihood Inference via Valued Attributes Regularization*, TASK 3 stands for *Intention Justification via Comparison*, TASK 4 stands for *Intention Evolution Modeling*.

Task	5-shots Prompt
TASK 1	<p>Act as a customer who is browsing a series of products given as follows. <session product information> You hold an assumed intention, which will be provided later. After seeing the products, you will be asked to determine the likelihood of purchasing the last product \ based on the assumed intention. You will be given four options to choose from: Yes, Maybe yes, Maybe no, No. Please select the most appropriate option based on the given context. A. Yes: The product is a logical and reasonable outcome of the purchasing intention. B. Maybe yes: I may consider this, but it's not a strong impulse. C. Maybe no: The product is not directly related to my intention. D. No: I would never purchase it if I were the customer with the given intention.</p> <p>Here are a few examples: Q: After seeing Eco-friendly laundry detergent, bamboo dish brush, reusable kitchen cloths, and assuming you are a customer who have the intention of Reducing household chemical usage. How likely are you to purchase A biodegradable dish soap based on the assumed intention? A: A. Yes</p> <p>Q: After seeing Instant Pot, KitchenAid Stand Mixer, Ninja Air Fryer, and assuming you are a customer who have the intention of Upgrading kitchen equipment for home cooking. How likely are you to purchase A set of gourmet spices based on the assumed intention? A: C. Maybe no</p> <p>Q: After seeing Columbia hiking boots, North Face backpack, Garmin GPS watch, and assuming you are a customer who have the intention of Planning for outdoor adventures. How likely are you to purchase A formal suit for weddings based on the assumed intention? A: D. No</p> <p>Q: After seeing "1984" by George Orwell, "To Kill a Mockingbird" by Harper Lee, \ "The Catcher in the Rye" by J.D. Salinger, and assuming you are a customer who have the intention of Finding new reading material for leisure. How likely are you to purchase "The Da Vinci Code" by Dan Brown based on the assumed intention? A: B. Maybe yes</p> <p>Q: After seeing Rolex Submariner, Omega Seamaster, Tag Heuer Monaco, and assuming you are a customer who have the intention of Finding a timeless gift for a special occasion. How likely are you to purchase A limited edition Patek Philippe watch based on the assumed intention? A: A. Yes</p> <p>Q: After seeing <previous products>, and assuming you are a customer who have the intention of <second last intention>. How likely are you to purchase <last product> based on the assumed intention? A:</p>

Table 9: 5-shots prompts for model evaluation. TASK 1 stands for *Intent-Based Purchasing Likelihood Estimation*

Task	5-shots Prompt
TASK 2	<p>Act as a customer who is browsing a series of products given as follows. <session product information> You have a valued feature/attribute, which will be provided later. After seeing the products, you will be asked to determine the likelihood of purchasing the last product \ based on the valued attribute. You will be given four options to choose from: Yes, Maybe yes, Maybe no, No. Please select the most appropriate option based on the given context. A. Yes: The product logically and reasonably matches the characteristics I value. B. Maybe yes: I might consider this product, but it doesn't strongly appeal to me. C. Maybe no: The product does not directly relate to the characteristic I value. D. No: I would not purchase this product if I were focused on the given characteristic.</p> <p>Here are a few examples: Q: After seeing Noise-canceling headphones, wireless earbuds, Bluetooth speaker, and assuming you are a customer who highly value the feature High audio quality of Bluetooth speaker. How likely are you to purchase A premium soundbar? A: A. Yes</p> <p>Q: After seeing adjustable standing desk, monitor with blue light filter, Ergonomic office chair, and assuming you are a customer who highly value the feature Ergonomics of Ergonomic office chair. How likely are you to purchase A desk lamp with a USB port? A: C. Maybe no</p> <p>Q: After seeing Organic facial cleanser, natural moisturizer, chemical-free sunscreen, and assuming you are a customer who highly value the feature Natural ingredients of chemical-free sunscreen. How likely are you to purchase A synthetic fragrance? A: D. No</p> <p>Q: After seeing DSLR camera, camera tripod, external flash, and assuming you are a customer who highly value the feature Professional photography of external flash. How likely are you to purchase A photo editing software? A: A. Yes</p> <p>Q: After seeing High SPF sunscreen, UV-blocking sunglasses, wide-brimmed hat, and assuming you are a customer who highly value the feature Sun protection of wide-brimmed hat. How likely are you to purchase An aloe vera gel? A: B. Maybe yes</p> <p>Q: After seeing <previous products>, and assuming you are a customer who highly value the feature <second last intention attribute> \ of <second last product>. How likely are you to purchase <last product>? A:</p>

Table 10: 5-shots prompts for model evaluation. TASK 2 stands for *Purchasing Likelihood Inference via Valued Attributes Regularization*.

Task	5-shots Prompt
TASK 3	<p>Act as a customer who is browsing a series of products given as follows. <session product information> You have an assumed intention, which will be provided later. You will be asked to evaluate the provided comparison between the last two products \ based on the assumed intention. You will be given four options to choose from: Yes, Maybe yes, Maybe no, No. Please select the most appropriate option based on the given context. A. Yes: the comparison is reasonable and detailed enough to justify the change. B. Maybe yes: The comparison could be more detailed and thorough but can be ignored. C. Maybe no: The comparison is not entirely reasonable or lacks sufficient in-depth detail. D. No: The comparison does not provide any underlying reasons or insights.</p> <p>Here are a few examples: Q: Comparing between a budget smartphone with a long battery life and A high-end smartphone with \ superior low-light performance, and assuming you have the intention of Finding a device with the best camera quality, Does this comparison The high-end smartphone boasts advanced camera technology \ provide in-depth justification of your impulse? A: A. Yes</p> <p>Q: Comparing between A compact car and a mid-size SUV, and assuming you have the intention of Prioritizing fuel efficiency, Does this comparison the mid-size SUV, although spacious, consumes more fuel due to its larger engine \ and heavier body provide in-depth justification of your impulse? A: B. Maybe yes</p> <p>Q: Comparing between A luxury wristwatch and a fitness tracker, and assuming you have the intention of Tracking health metrics, Does this comparison Finding a more affordable watch provide in-depth justification of your impulse? A: D. No</p> <p>Q: Comparing between A leather office chair with plush cushioning and \ a mesh office chair with lumbar support and assuming you have the intention of Seeking maximum comfort during long working hours, Does this comparison The mesh office chair offers better breathability and ergonomic support \ provide in-depth justification of your impulse? A: A. Yes</p> <p>Q: Comparing between A hardcover book and an e-reader, and assuming you have the intention of Enhancing the reading experience, Does this comparison The hardcover book provides a tactile, while the e-reader offers portability, \ adjustable text size provide in-depth justification of your impulse? A: C. Maybe no</p> <p>Q: Comparing between <last two products>, and assuming you have the intention of <last two intention>, Does this comparison <last intention comparison> provide in-depth justification of your impulse? A:</p>

Table 11: 5-shots prompts for model evaluation. TASK 3 stands for *Intention Justification via Comparison*.

Task	5-shots Prompt
TASK 4	<p>Act as a customer who is browsing a series of products given as follows. <session product information> You will be provided with a sequence of intention. You will be asked to determine whether you still want to explore similar products \ based on the sequence of intention. You will be given three options to choose from: Yes, Maybe yes, No. Please select the most appropriate option based on the given context. A. Yes: I want to explore products under the same category. B. Maybe yes: I want to explore products under the same category but with different features. C. No: I want to explore products under other categories.</p> <p>Here are a few examples:</p> <p>Q: After seeing Stainless steel kitchen knives, non-stick frying pans, silicone spatulas, and assuming you have the intention of Upgrading kitchen tools for home cooking, do you still want to explore similar products? A: B. Maybe yes</p> <p>Q: After seeing Fitness tracker, yoga mat, resistance bands, and assuming you have the intention of Tracking fitness progress, do you still want to explore similar products? A: A. Yes</p> <p>Q: After seeing Stainless steel refrigerator, smart oven, induction cooktop, and assuming you have the intention of Making the kitchen more energy efficient, do you still want to explore similar products? A: C. No</p> <p>Q: After seeing Smart thermostat, LED light bulbs, energy-efficient washing machine, and assuming you have the intention of Saving on utility bills, do you still want to explore similar products? A: B. Maybe yes</p> <p>Q: After seeing Indoor plants, plant stands, watering can, and assuming you have the intention of Creating a greener living space, do you still want to explore similar products? A: A. Yes</p> <p>Q: After seeing <previous products>, and assuming you have the intention of <previous new intention>, do you still want to explore similar products? A:</p>

Table 12: 5-shots prompts for model evaluation. TASK 4 stands for *Intention Evolution Modeling*.

Task	Chain-of-Thought Prompt
TASK 1	<p>Act as a customer who is browsing a series of products given as follows. <session product information></p> <p>After seeing <previous products>, and assuming you are a customer who have the intention of <second last intention>. How likely are you to purchase <last product> based on the assumed intention?</p> <p>A. Yes: The product is a logical and reasonable outcome of the purchasing intention. B. Maybe yes: I may consider this, but it's not a strong impulse. C. Maybe no: The product is not directly related to my intention. D. No: I would never purchase it if I were the customer with the given intention.</p> <p>Answer with a brief rationale then make your final choice \ by answering the option alphabet A/B/C/D only in the last line of your response. Your Answer:</p>
TASK 2	<p>Act as a customer who is browsing a series of products given as follows. <session product information></p> <p>After seeing <previous products>, and assuming you are a customer who highly value the feature <second last intention attribute> \ of <second last product>. How likely are you to purchase <last product>?</p> <p>A. Yes: The product logically and reasonably matches the characteristic I value. B. Maybe yes: I might consider this product, but it doesn't strongly appeal to me. C. Maybe no: The product does not directly relate to the characteristic I value. D. No: I would not purchase this product if I were focused on the given characteristic.</p> <p>Answer with a brief rationale then make your final choice \ by answering the option alphabet A/B/C/D only in the last line of your response. Your Answer:</p>
TASK 3	<p>Act as a customer who is browsing a series of products given as follows. <session product information></p> <p>Comparing between <last two products>, and assuming you have the intention of <last two new intention>, Does this comparison <last intention comparison> provide in-depth justification of your impulse?</p> <p>A. Yes: the comparison is reasonable and detailed enough to justify the change. B. Maybe yes: The comparison could be more detailed and thorough but can be ignored. C. Maybe no: The comparison is not entirely reasonable or lacks sufficient in-depth detail. D. No: The comparison does not provide any underlying reasons or insights.</p> <p>Answer with a brief rationale then make your final choice \ by answering the option alphabet A/B/C/D only in the last line of your response. Your Answer:</p>
TASK 4	<p>Act as a customer who is browsing a series of products given as follows. <session product information></p> <p>After seeing <previous products>, and assuming you have the intention of <previous new intention>, do you still want to explore similar products?</p> <p>A. Yes: I want to explore products under the same category. B. Maybe yes: I want to explore products under the same category but with different features. C. No: I want to explore products under other categories.</p> <p>Answer with a brief rationale, then make your final choice \ by answering the option alphabet A/B/C only in the last line of your response. Your Answer:</p>

Table 13: Chain-of-Thought prompts for model evaluation. TASK 1 stands for *Intent-Based Purchasing Likelihood Estimation*, TASK 2 stands for *Purchasing Likelihood Inference via Valued Attributes Regularization*, TASK 3 stands for *Intention Justification via Comparison*, TASK 4 stands for *Intention Evolution Modeling*.