

A Usage-Centric Take on Intent Understanding in E-Commerce

Anonymous ACL submission

Abstract

Identifying and understanding user intents is a pivotal task for E-Commerce. Despite its popularity, intent understanding has not been consistently defined or accurately benchmarked. In this paper, we focus on predicative user intents as “how a customer use a product”, and pose intent understanding as a natural language reasoning task, independent of product ontologies. Through topology analysis, we highlight two weaknesses in FolkScope, the SOTA E-Commerce Intent Knowledge Graph, precluding it from effectively reasoning about user intents and recommending diverse useful products. Following these observations, we propose a product recovery benchmark to isolate intent understanding abilities from confounders. We verify the identified weaknesses, and discuss future directions for intent understanding.¹

1 Introduction

User intents are a crucial source of information for E-Commerce (Zhang et al., 2016; Hao et al., 2022). Intents reveal users’ motivation in E-Commerce interactions: suppose a user plans to go for **outdoor barbecue**, their intent may not refer only to barbecue smoker grills but also to other items that can be potentially useful for outdoor barbecue, such as disposable cutlery or plates. In these cases, traditional product recommendation approaches would fail to handle these queries or to remind customers of the products they may need but have forgotten.

Intent Understanding offers the crucial ability to recommend distinct products based on common user intents they fulfil. It involves identifying user intents and connecting them with products. To identify intents, each method summarizes a profile of user intents for each product listing, from user interactions (e.g. co-buy records, reviews). Then, at intent-product relation mining, each method learns to predict useful products based on user intents.

One significant challenge towards effective intent understanding is the poor definition of user intents, which precludes effective intent identification and can easily result in contaminated intent-product associations. In prior work (Yu et al., 2023; Luo et al., 2021), user intents are often blended with “product properties” or “similar products”. While these are also related to user intents, we argue that they are, in nature, shortcuts which benefit existing product recommendation benchmarks, but do not necessarily align with the objective for intent understanding, namely, to retrieve superficially distinct kinds of products serving common intents.

Therefore, we propose a usage-centric paradigm of intent understanding. In this paradigm, user intents are focused on natural language predicative phrases, describing how customers **use** a product; also, instead of individual product listings, we aim to predict *kinds of products* **useful** for an intent. Specifically, user intents are activities to accomplish (e.g. **outdoor barbecue**) or situations to resolve (**lower-back pain**); kinds of products are clusters of product listings of the same category (e.g. **scrub brush**) with a common property (e.g. **stiff bristle**). The task then is a natural language reasoning task, closely related to commonsense reasoning (Sap et al., 2019; Bosselut et al., 2019), in the form of: “The user has intent I ” entails “The kind of product P is useful for the user.”

Under the usage-centric paradigm, we present an analysis of a SOTA E-Commerce intent KG, FolkScope (Yu et al., 2023), which reported positive results on an intrinsic co-buy prediction task. Refactoring their KG to model associations between kinds of products and their usage intents, we find two unsatisfactory characteristics in their KG topology: 1) *property-ambiguity*: generated user intents are poorly aligned with relevant product properties, such that the KG often map user intents to kinds of product with the relevant category but fairly random properties; 2) *category-rigidity*: each

¹We will release our code and datasets.

081 intent is strongly associated with a single category
082 of product, such that the KG is unable to recom-
083 mend distinct products that serve common intents.

084 In light of these findings, we introduce a chal-
085 lenge dataset for usage-centric intent understand-
086 ing, the Product Recovery Benchmark. This bench-
087 mark isolates intent understanding abilities from
088 confounders, where each method aims to recover
089 the kinds of useful products for customers having a
090 particular intent profile. Here, we further validate
091 the impact of the weaknesses in current SOTA.

092 To summarize, in this paper: 1) we propose a
093 usage-centric paradigm of intent understanding, as
094 natural language reasoning; 2) we present the prod-
095 uct recovery benchmark for usage-centric intent
096 understanding, and report results with SOTA base-
097 lines; 3) we identify crucial weaknesses in existing
098 SOTA as *category-rigidity* and *property-ambiguity*,
099 where we propose intent mining from genuine user
100 reviews as a promising future direction.

101 2 Usage-Centric Intent Understanding

102 We propose a usage-centric paradigm of intent un-
103 derstanding, focusing on usage user intents and
104 the kinds of useful products, where each method
105 aims to ground usage user intents in kinds of useful
106 products. Differently from the “informal queries”
107 in Luo et al. (2021), and similarly to Ding et al.
108 (2015), our usage user intents are generic eventual-
109 ities/situations, independent of product ontologies.

110 We introduce *kinds of products* as the target gran-
111 ularity level, as it abstracts away the nuanced dif-
112 ferences among individual listings, and yields a
113 purely natural language setup, independent of prod-
114 uct ontologies. It contains just enough information
115 (category + property) to represent the product list-
116 ings inside for intent understanding.

117 User intents rarely require combinations of prop-
118 erties in a product category. Therefore, to avoid
119 generating factorial numbers of kinds of product,
120 we impose a mild constraint that only one property
121 is specified for each kind of product.

122 We demonstrate the specificity trade-off with
123 an example below: for **outdoor barbecues**, a **stiff-**
124 **bristle scrub brush** is useful for cleaning the grease
125 on the grill. To that end, there are many list-
126 ings of hard-bristle scrubs but the exact choice
127 among them is irrelevant to the user intent and
128 could be identified by downstream recommenda-
129 tion systems using other factors (e.g. customer
130 habit, geo-location, etc.). However, the **stiff bristle**

property is essential for a candidate to be suitable
for **outdoor barbecues**. In short, grouping based on
kinds of products strikes a balance between spar-
sity that comes with specificity, and ambiguity that
comes with generality.

131 3 FolkScope Analysis 136

137 3.1 KG Refactoring 137

138 We refactor FolkScope to our usage-centric intent
139 understanding paradigm. FolkScope KG connects
140 product listings with their user intents, which are
141 generated with OPT-30B (Zhang et al., 2022) when
142 given pairs of co-bought products sourced from
143 Amazon Reviews Dataset (ARD) (Ni et al., 2019),
144 along with commonsense relations.

145 Among their 18 commonsense relations, we fil-
146 ter out all “item” relations as well as 3 “function”
147 relations (*SymbolOf*, *MannerOf*, and *DefinedAs*),
148 since they are nominal in nature, and are irrelevant
149 to product usage. We keep the remaining 5 predica-
150 tive relations, *UsedFor*, *CapableOf*, *Result*, *Cause*,
151 *CauseDesire*, as legitimate user intents.

152 To group the product listings into kinds of prod-
153 ucts, we take the fine-grained product categories
154 from ARD (e.g. *Kids’ Backpacks*), and borrow
155 the attributes under the relation *PropertyOf* in the
156 original FolkScope KG as properties.²

157 We compute the association strengths from le-
158 gitimate user intents to common kinds of products
159 by aggregation. Let $e(I_i, P_j)$ be the connection of
160 intent I_i with product listing P_j , P_j belongs to a
161 kind of products K_k . The association strength for
162 edges in the refactored KG are then computed as:
163 $e'(I_i, K_k) = \sum_{P_j \in K_k} pm_i(P_j, K_k) * e(I_i, P_j)$.³

164 3.2 Statistical Analysis 164

165 Through statistical analysis of the KG, we identify
166 two major weaknesses in the FolkScope KG: it is
167 over-specific about categories of useful products,
168 but under-specific about the required properties in
169 these categories. Intents in FolkScope KG tend to
170 be associated with products with vague properties
171 from a single category, rather than specific kinds
172 of products from diverse categories.

173 **Property-Ambiguity** For each user intent, we
174 look into the distribution of its edge weights among

²These attributes do not fit the criteria for usage user in-
tents, but they are acquired through generic LLM prompted
summarization, and thus are borrowed as product properties.

³The *pmi* term penalizes product listings with multiple
kinds of products (e.g. multiple properties in one listing).

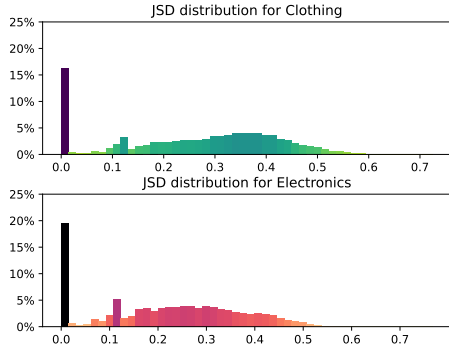


Figure 1: Histograms of Jensen-Shannon Divergence for each user intent. Values are packed around 0: distributions of edge weights conditioned on intents are close to unconditioned frequency priors.

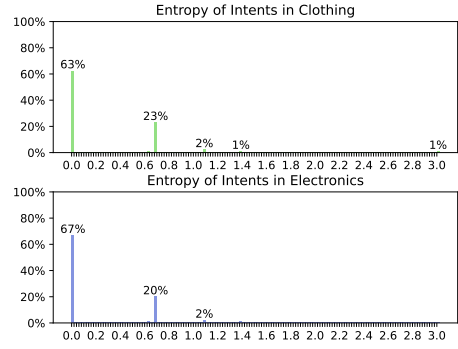


Figure 2: Histograms of category-entropy for each user intent. Values are concentrated at 0.0 and 0.7, meaning the intent is associated with only 1 / 2 categories.

different kinds of products. We compare these posterior edge-weight distributions, conditioned on the intent, against the prior frequency distributions. We calculate Jensen-Shannon Divergence (JSD) between these conditional and prior distributions (see Figure 1): for up to 20% of cases, JSD is < 0.1 , where only 2% of cases have $JSD > 0.5$.

This shows, the KG’s edge weights among different kinds of products are strongly predicted by their prior distribution, and are insensitive to the specific usages depicted by user intents. We credit this to the mismatch between property and intent mining: each product listing may have multiple properties and may serve multiple intents, but the mappings between these properties and intents are underspecified.

Category-Rigidity For each user intent, we measure the category-diversity of its edge weights in the refactored KG: we compute the entropy of its edge weights grouped by product categories.⁴

Figure 2 shows the entropy meta-distributions: entropy values are concentrated in 2 narrow ranges, $[0, 0.02)$ and $[0.68, 0.70)$. We notice that an entropy in $[0, 0.02)$ indicates that the associations about this intent are focused on only one product category; $[0.68, 0.70)$ indicates that the associations are focused on two product categories. Therefore, from Figure 2 we can conclude that over 80% of the intents are associated with only one or two categories. This category-rigidity in FolkScope hampers its ability to recommend diverse kinds of products, as we will discuss in §4.2.

⁴Please see Appendix B for an example.

4 The Product Recovery Benchmark

We introduce a product recovery challenge to measure success in usage-centric intent understanding.

4.1 Dataset Elicitation

We develop the product recovery benchmark based on the Amazon Reviews Dataset (Ni et al., 2019), a pool of product listings, enriched with English product descriptions, category information, anonymized user purchase records and reviews.⁵

For each data entry, we take a product listing from the Amazon Reviews Dataset, and acquire the kinds of products as in §3.1 for prediction targets.

At inference time, given a product listing, each method predicts a set of user intents (using product description, user reviews, etc.). Then, using **solely** the predicted intent profile as input, the method recovers useful kinds of products based on its knowledge of E-Commerce demands (either in symbolic KGs, or in LLMs). The predicted kinds of products are compared against: 1) *bought-product-recovery*: kinds of product to which the current product belongs; 2) *co-bought-product-recovery*: kinds of products co-bought with the current product.

We take *bought-product-recovery* as our main evaluation setup, since it focuses on intent-to-kinds-of-product associations. Compared to the product recommendation evaluation in Yu et al. (2023), it marginalizes over confounder factors inciting co-buy behaviour (e.g. brand loyalty, geolocation, etc.). We also include the co-bought-product-recovery setup⁶ to evaluate cross-category recom-

⁵Note that our elicitation procedure is corpus-agnostic, we empirically select ARD because it is the largest available, and for consistency with evaluation in Yu et al. (2023).

⁶To re-cap, in this setup we also predict co-buy behavior, as in product recommendation, but here we only predict kinds of products in other categories than the bought product.

Models	Clothing	Electronics
FolkScope	0.192	0.263
FolkScope – properties	0.116	0.166
FolkScope + GPT	0.187	0.257

Table 1: MRR_{\max} for *bought-product-recovery* task.

mendation performance between distinct products.

Evaluation metric Following prior work (Chen and Wang, 2013), we measure success by Mean Reciprocal Rank (MRR) of gold kinds of products in the predicted distributions. We focus our evaluation on being able to successfully predict the gold kind of products associated with user intent. In case multiple gold kinds of products are assigned for a product listing, the highest-ranking hit is taken to calculate the MRR_{\max} (see Appendix Eq. 2).

4.2 Experiments and Results

We evaluate the FolkScope KG (refactored in §3.1) with the Product Recovery benchmark. We offer the baseline results in Table 1, and highlight below the impact of weaknesses discussed in §3.2.

Property-Ambiguity To understand how property ambiguity affects FolkScope performance, we compare it with another prior property baseline derived from it: for each evaluation entry, we corrupt the FolkScope predictions by replacing the property in the predicted kinds of products based on the property popularity. (see Appendix A.2 for details)

From Table 1, we observe that *FolkScope – properties* reached respectable performance with only moderate regression from FolkScope predictions. This limited MRR gap shows the impact of property-ambiguity, where performance gains could be expected with better property alignment.

Category-Rigidity To validate the category-rigidity observation in §3.2, we also evaluate the FolkScope KG in the co-bought-product-recovery setup, where we specifically use it to predict kinds of co-bought products in **other categories**.

In this setup, we observe low MRR_{\max} of 0.077 and 0.033 for *Clothing* and *Electronics* domains, respectively: the FolkScope KG cannot effectively recommend superficially distinct kinds of products connected by common user intents.

Notably, between the two domains, FolkScope reaches a slightly higher MRR_{\max} in *Clothing*. This is consistent with our findings in Figure 2, where category-entropy values are slightly more spreaded than in *Electronics*.

LLM Rerank We also evaluate LLM performance in our product recovery benchmark with GPT-3.5-turbo (Brown et al., 2020). Ideally, we would like the LLM to predict useful kinds of products end-to-end. However, due to the difficulty of reliably matching LLM predictions with gold kinds of products⁷, we instead adopt a re-ranking paradigm, where we prompt the LLM to re-rank the top-10 kinds of products predicted by FolkScope (see Appendix A.4 for details).

As Table 1 shows, we observe no substantial improvement with LLM-reranking. We investigate this failure by looking into where hits are met in the predictions: the MRR_{\max} of 0.192 and 0.263 actually consist of 16% and 22% of hits-at-1 ($RR_{\max} = 1.0$), 73% and 63% of no-hits-in-top10 ($RR_{\max} < 0.1$), and only 11% and 15% of hits between 2 to 10 ($0.1 < RR_{\max} < 1.0$). These polarized distributions leave little room for re-ranking to take effect.

We raise the warning that dataset artefacts from the common source corpus (AWD) could be behind this abnormally high hit-at-1 rate (compared with the MRR_{\max} value), where the reported MRR_{\max} values may have been inflated. Due to the lack of another large E-Commerce Reviews corpus, we leave further investigations for future work.

5 Discussions and Conclusion

In this paper, we revisit intent understanding from a usage-centric perspective as a natural language reasoning task, aiming to detect superficially distinct kinds of products useful for common usage intents. We have introduced a novel Product Recovery benchmark, and have investigated two weaknesses of the SOTA FolkScope KG in supporting usage-centric intent understanding: *Property Ambiguity* and *Category-Rigidity*.

We advocate for adopting the usage-centric intent understanding paradigm, and we believe that the above weaknesses can be alleviated by considering user reviews, in addition to co-buy records. The former weakness can be addressed by the fact that relevant product properties and usage intents are likely to co-occur in product reviews; whereas the latter, can be ameliorated by the increased likelihood of the same usage intent to appear consistently in reviews across different categories.

⁷In Appendix C, we also include an LLM-only experiment using GPT-4 as the matching metric: we find no evidence of the LLM-only method outperforming the FolkScope baseline, and find GPT-4 matching metric is over permissive.

328 Limitations

329 In this paper, we have proposed to study E-
330 Commerce intent understanding from a usage-
331 centric perspective. Due to the lack of consis-
332 tent task definition, we are only able to anal-
333 yse one SOTA intent understanding KG (namely
334 FolkScope) and one SOTA LLM. We encourage
335 more research attention on the usage-centric E-
336 commerce intent understanding task for a more
337 diverse landscape.

338 We have established that weaknesses of Prop-
339 erty Ambiguity and Category Rigidity exist in the
340 SOTA KG, and we have offered a principled hy-
341 pothesis that utilizing genuine user reviews could
342 help with these weaknesses. However, due to lim-
343 its to the scope of this paper, we do not provide
344 empirical evidence for this hypothesis and leave it
345 as a promising direction of future work.

346 We note that as this paper is related to recommen-
347 dation, there exists risks that methods developed
348 on the Product Recovery Benchmark may be used
349 to bias customer decisions; on the other hand, we
350 also note that our task definition is purely natural
351 language and does not involve any individual prod-
352 uct listings, therefore it would not bias customer
353 choices among directly competing listings of the
354 same kinds of products.

355 References

356 Antoine Bosselut, Hannah Rashkin, Maarten Sap, Chai-
357 tanya Malaviya, Asli Celikyilmaz, and Yejin Choi.
358 2019. [COMET: Commonsense Transformers for](#)
359 [Automatic Knowledge Graph Construction](#). In *Pro-*
360 *ceedings of the 57th Annual Meeting of the Asso-*
361 *ciation for Computational Linguistics*, pages 4762–
362 4779, Florence, Italy. Association for Computational
363 Linguistics.

364 Tom B. Brown, Benjamin Mann, Nick Ryder, Melanie
365 Subbiah, Jared Kaplan, Prafulla Dhariwal, Arvind
366 Neelakantan, Pranav Shyam, Girish Sastry, Amanda
367 Aspell, Sandhini Agarwal, Ariel Herbert-Voss,
368 Gretchen Krueger, Tom Henighan, Rewon Child,
369 Aditya Ramesh, Daniel M. Ziegler, Jeffrey Wu,
370 Clemens Winter, Christopher Hesse, Mark Chen, Eric
371 Sigler, Mateusz Litwin, Scott Gray, Benjamin Chess,
372 Jack Clark, Christopher Berner, Sam McCandlish,
373 Alec Radford, Ilya Sutskever, and Dario Amodei.
374 2020. [Language Models are Few-Shot Learners](#).
375 [ArXiv:2005.14165 \[cs\]](#).

376 Li Chen and Feng Wang. 2013. [Preference-based clus-](#)
377 [tering reviews for augmenting e-commerce recom-](#)
378 [mendation](#). *Knowledge-Based Systems*, 50:44–59.

Xiao Ding, Ting Liu, Junwen Duan, and Jian-Yun Nie.
2015. [Mining User Consumption Intention from So-](#)
[cial Media Using Domain Adaptive Convolutional](#)
[Neural Network](#). *Proceedings of the AAAI Confer-*
ence on Artificial Intelligence, 29(1). Number: 1.

Zhenyun Hao, Jianing Hao, Zhaohui Peng, Senzhang
Wang, Philip S. Yu, Xue Wang, and Jian Wang. 2022.
[Dy-hien: Dynamic evolution based deep hierarchi-](#)
[cal intention network for membership prediction](#). In
Proceedings of the Fifteenth ACM International Con-
ference on Web Search and Data Mining, WSDM '22,
page 363–371, New York, NY, USA. Association for
Computing Machinery.

Xusheng Luo, Le Bo, Jinhang Wu, Lin Li, Zhiy Luo,
Yonghua Yang, and Keping Yang. 2021. [AliCoCo2:](#)
[Commonsense Knowledge Extraction, Representa-](#)
[tion and Application in E-commerce](#). In *Proceedings*
of the 27th ACM SIGKDD Conference on Knowledge
Discovery & Data Mining, pages 3385–3393, Virtual
Event Singapore. ACM.

Jianmo Ni, Jiacheng Li, and Julian McAuley. 2019. [Jus-](#)
[tifying Recommendations using Distantly-Labeled](#)
[Reviews and Fine-Grained Aspects](#). In *Proceedings*
of the 2019 Conference on Empirical Methods in
Natural Language Processing and the 9th Interna-
tional Joint Conference on Natural Language Pro-
cessing (EMNLP-IJCNLP), pages 188–197, Hong
Kong, China. Association for Computational Lin-
guistics.

Maarten Sap, Ronan Le Bras, Emily Allaway, Chan-
dra Bhagavatula, Nicholas Lourie, Hannah Rashkin,
Brendan Roof, Noah A. Smith, and Yejin Choi. 2019.
[ATOMIC: An Atlas of Machine Commonsense for](#)
[If-Then Reasoning](#). *Proceedings of the AAAI Confer-*
ence on Artificial Intelligence, 33:3027–3035.

Changlong Yu, Weiqi Wang, Xin Liu, Jiabin Bai,
Yangqiu Song, Zheng Li, Yifan Gao, Tianyu Cao, and
Bing Yin. 2023. [FolkScope: Intention Knowledge](#)
[Graph Construction for E-commerce Commonsense](#)
[Discovery](#). [ArXiv:2211.08316 \[cs\]](#).

Chenwei Zhang, Wei Fan, Nan Du, and Philip S. Yu.
2016. [Mining user intentions from medical queries:](#)
[A neural network based heterogeneous jointly mod-](#)
[eling approach](#). In *Proceedings of the 25th Interna-*
tional Conference on World Wide Web, WWW '16,
page 1373–1384, Republic and Canton of Geneva,
CHE. International World Wide Web Conferences
Steering Committee.

Susan Zhang, Stephen Roller, Naman Goyal, Mikel
Artetxe, Moya Chen, Shuohui Chen, Christopher De-
wan, Mona Diab, Xian Li, Xi Victoria Lin, Todor Mi-
haylov, Myle Ott, Sam Shleifer, Kurt Shuster, Daniel
Simig, Punit Singh Koura, Anjali Sridhar, Tianlu
Wang, and Luke Zettlemoyer. 2022. [Opt: Open pre-](#)
[trained transformer language models](#).

A Implementation Details

A.1 Benchmark data split

We follow Yu et al. (2023), and we split product instance in FolkScope KG into training, validation and test splits with respective portions of 80%, 10% and 10%. Please refer to Table 2 for detailed statistics. Note that *Clothing* stands for the “Clothing, Shoes and Jewelry” domain in the Amazon Reviews Dataset, and *Electronics* simply stands for the “Electronics” domain in the Amazon Reviews Dataset.

Categories	Train	Validation	Test
Clothing	30296	2027	2088
Electronics	85086	7853	7900

Table 2: Number of product listings in the training, validation and test set. Please note that we drop product listings that lack related kinds of products, so the ratio of the number of instances across the splits are not exactly equal to 8:1:1.

A.2 Prior Property Baseline

For each kind of product in the prediction list, we corrupt its property part with respect to its prior popularity within its fine-grain category in the Amazon Reviews Dataset. Popularity is defined as the frequency of a property appearing with the product listing having this corresponding fine-grained category. To avoid repeated kinds of products in the predictions, when multiple predicted kinds of products from the same category are predicted, we draw properties top-down w.r.t. popularity for each prediction.

A.3 Evaluation Metric

Our evaluation metric MRR_{\max} can be formally defined as follows:

$$RR_{\max}(l) = \max_{c \in C_{\text{gold}}(l)} (\text{rank}(c)^{-1}) \quad (1)$$

$$MRR_{\max} = \frac{\sum_{l \in L} RR_{\max}(l)}{|L|} \quad (2)$$

where RR represents the Reciprocal Rank, $C_{\text{gold}}(l)$ are the gold clusters for the listing l and L is the set of all listings in the benchmark.

A.4 GPT-3.5-turbo Re-ranking

For each product listing l , when there is no predicted kind of products given a set of related user

intents, we mark the $RR_{\max}(l)$ as 0 both before and after re-ranking.

A.4.1 Re-ranking Prompt

A product is suitable for the following purposes:

$\{Intents\}$

Please rank the following categories in order of likelihood that the product belongs to them (most likely to least likely):

$\{kinds\ of\ products\ list\} \dots$

Answer:

1.

We fill *Intents* with a set of mined user intents and *kinds of products list* with the top 10 predictions for kinds of products.

Note that in this setting and in § C.1.1, we still use the term “category” in LLM prompts to refer to kinds of products, because during preliminary experiments we found that LLMs do not respond well to the term “kind of product”.

B Category Rigidity

In a E-Commerce user intent KG, a non-negligible amount of usage user intents should entail the demand for diverse products from different categories.

For the example of **outdoor barbecues**, for outdoor barbecues one may need not only **scrub brush**, but also other categories of products, such as **picnic blankets**, **grill gloves**, etc.

Therefore, we take the category-entropy of the edges for each user intent, to measure how diverse the KG edges are w.r.t. categories. We add up the edge weights grouped by product categories (e.g. edge weights to **stiff bristle scrub brush** and **scrub brush with wooden handle** are added together), and compute the entropy of the converted category distribution. As discussed in §3.2, we found severe category rigidity in the FolkScope KG, where very few user intents have diverse category distributions, the majority of user intents are associated with only one category, followed by those associated with two.

C GPT End-to-End Evaluation

We perform an additional experiment to directly predict kinds of products in an end-to-end setup,

with an LLM, for our proposed product recovery task. Again, we use GPT-3.5-turbo as the LLM and design the zero-shot prompt as in §C.1.1. However, due to the absence of the complete ontology of the Amazon Reviews Dataset, it is challenging for GPT-3.5-turbo to predict the exact ground truth kinds of products. To sidestep the difficulty of evaluating whether the predicted strings are semantically identical to the ground truth labels, we use GPT-4 to judge whether there is a match between predicted and ground truth labels. The relevant prompt is specified in §C.1.2. The detailed evaluation results is presented in Table 3.

	Clothing	Electronics
GPT-3.5-turbo	0.511	0.543
FolkScope	0.527	0.671

Table 3: MRR_{max} score when evaluating using GPT-4 as the judge for matching. Values for GPT-3.5-turbo and our baseline refactored FolkScope KG are both higher in absolute values due to the more benign matching criterion; the LLM baseline with GPT-3.5-turbo does not outperform the KG baseline.

From Table 3, we can observe that GPT-3.5-turbo does not outperform the FolkScope KG baseline on the product recovery benchmark. Compared to the strict string matching results in Table 1, GPT-4 evaluation has a significantly more permissive criterion on matching, yielding much higher MRR_{max} values. We find many of these “matched” verdicts by GPT-4 to be spurious (see Table 4), and conclude that GPT-4 cannot easily achieve reliable matching for the product recovery benchmark, and more robust criteria are needed before replacing the exact match criterion.

C.1 Prompt Examples

C.1.1 Kinds of Products Prediction

Intents:

{*intents*}

Given the intents, please predict the top 10 kinds of products that will be useful for these intents.

A kind of product is the concatenation of a fine-grained category from the Amazon Review Dataset and a useful property. For example: Clothing, Shoes & Jewelry|Men|Watches|Wrist Watches ### leather.

Kinds of products:

1.

C.1.2 Prediction Evaluation

Here is a list of predicted categories:

{*prediction*}

Validate each prediction based on the ground truth categories[T/F].

Each prediction can be considered true when it is similar to one of the ground truth categories.

Ground truth categories:

{*ground truth*}

D Computational Budget

D.1 Main Experiments

All the benchmark construction and evaluation has been performed using 2 x Intel(R) Xeon(R) Gold 6254 CPUs @ 3.10GHz.

FolkScope KG Refactoring We converted all the intents generated by FolkScope without applying any of its proposed filters based on the graph evaluation results on the validation set. The whole graph generation for both domains takes around 24 hours in total.

FolkScope Intents Evaluation We need around 71 and 6 hours for evaluating the intents for the test set of the Clothing and Electronics domain respectively.

D.2 LLM Experiments

We mainly use GPT-3.5-turbo and GPT-4 for our LLM-related experiments. Please refer to Table 5 for details about the relevant costs. For both models, we keep the default parameters from OpenAI, and set the temperature to 0 to facilitate reproducibility.

E Artifact Licenses

Amazon Reviews Dataset: Limited license for academic research purposes and for non-commercial use (subject to [Amazon.com Conditions of Use](#))

FolkScope: MIT license

Ground truth kinds of products
1. Clothing, Shoes & Jewelry Costumes & Accessories Men Accessories ### Wandering Gunman 2. Clothing, Shoes & Jewelry Costumes & Accessories Men Accessories ### Holster 3. Clothing, Shoes & Jewelry Costumes & Accessories Men Accessories ### Western
GPT-3.5-turbo prediction
1. Clothing, Shoes & Jewelry Men Costumes Western ### authentic ...
Ground truth kinds of products
1. Clothing, Shoes & Jewelry Women Jewelry Earrings Stud ### Jewelry 2. Clothing, Shoes & Jewelry Women Jewelry Earrings Stud ### Gemstone 3. Clothing, Shoes & Jewelry Women Jewelry Earrings Stud ### Sterling Silver
GPT-3.5-turbo prediction
1. Clothing, Shoes & Jewelry Women Earrings Stud Earrings ### elegant and beautiful ...

Table 4: Here we list two examples that GPT-4 validate with $RR_{\max} = 1$. In the first example, it validates the first prediction as true by matching the “property” part of the ground truth 3 with the main category of prediction 1. In the second example, the “property” part of prediction 1 is too general compared to all the ground truth kinds of products, but it still validates it as true.

Experiment	Clothing	Electronics
LLM Rerank	3.86 \$	1.38 \$
LLM End-to-End	15.57 \$	14.56\$

Table 5: API costs of our LLM-related experiments. For the LLM Rerank experiment, we re-rank all the data samples in the test set while for the End-to-End evaluation, we only sample 1000 data samples in the test set.