

Building Multi-turn Intent Classification with LLM-based Labeling

Anonymous ACL submission

Abstract

001 Intent classification is essential for customer
002 service routing, connecting customers to the
003 appropriate agents and reducing handling time
004 and operational cost. Developing a real-world
005 multi-turn intent classification system is chal-
006 lenging due to complex intent taxonomies, dy-
007 namic intent switching within conversations,
008 and limited labeled training data. To ad-
009 dress these challenges, we propose a scalable
010 multi-turn intent classification framework for e-
011 commerce customer service that models intent
012 along multiple dimensions. We introduce LLM-
013 based labeling strategies to annotate real cus-
014 tomer transcripts at scale and augment training
015 with LLM-simulated multi-turn dialogues that
016 expand coverage of topic and intent switches,
017 which are rare in existing transcripts. Through
018 extensive experiments, we find that explanation-
019 guided labeling with a self-critique step pro-
020 duces the most accurate training labels. Fine-
021 tuned models built on a RoBERTa backbone
022 outperform zero-shot LLM prompting while
023 achieving substantially lower inference latency.
024 Finally, we show that a hybrid approach that
025 combines the fine-tuned classifier with LLM
026 prompting further improves accuracy over ei-
027 ther component alone. Overall, our results pro-
028 vide practical guidance for building and de-
029 ploying high-accuracy, low-latency, large-scale
030 multi-turn intent classification systems.

031 1 Introduction

032 Agentic customer service has become increasingly
033 important for e-commerce platforms (Cui et al.,
034 2017; Zhou et al., 2023). Different agents, either
035 LLMs (Large Language Models) or workflows, are
036 designed or trained to handle specific customer
037 issues. Intent detection is therefore crucial for effi-
038 cient routing. Intent classification failures may lead
039 to irrelevant responses or unnecessarily escalate to
040 human agents, increasing operational cost and de-
041 grading customer experience (Qi et al., 2021).

In recent years, LLMs have demonstrated strong
042 potential for improving intent detection due to their
043 few-shot generalization and broad world knowl-
044 edge (Zhao et al., 2023; Arora et al., 2024). How-
045 ever, deploying intent models that achieve both
046 high accuracy and low latency at scale remains
047 challenging in customer service. 048

049 First, **scalability** becomes a bottleneck as busi-
050 ness domains or product lines expand. In real-world
051 e-commerce systems, a high-level intent such as
052 “Cancel” or “Refund” may be associated with multi-
053 ple products. Modeling intents at a product-specific
054 level leads to label explosion and increasing model
055 complexity. Moreover, maintaining a consistent
056 intent taxonomy and obtaining sufficient labeled
057 training data becomes difficult at scale (Qi et al.,
058 2021; Liu et al., 2024a). Additional **scalability**
059 challenges arise from compound intents. User ut-
060 terances may express multiple, non-mutually ex-
061 clusive intents, such as requesting an order can-
062 cellation while reporting an unrecognized charge.
063 Traditional flat intent classifiers are ill-suited to
064 this. 064

065 Second, **context carryover and intent switch-**
066 **ing** increase the difficulty of intent detection in
067 multi-turn conversations. The domain of user in-
068 tent might be inferred from prior context, while
069 user goals can shift in mid-dialogue. For example,
070 a customer may initially seek help troubleshooting
071 a service issue and later shift to requesting service
072 cancellation. Without modeling conversational his-
073 tory, a system may misattribute the current intent
074 or fail to link it to the relevant context. Prior work
075 has incorporated contextual signals for intent pre-
076 diction (Wu et al., 2021; Nandi et al., 2024), but
077 typically assumes access to comprehensive labeled
078 multi-turn data or context features. Addressing
079 topic shifts and intent switches in real customer-
080 service applications remains underexplored. 080

081 Finally, **low-latency requirements** constrain prac-
082 tical deployment in production. Although LLMs

083 have strong potential for intent classification, high
084 latency makes them unsuitable for real-time infer-
085 ence. Production deployments must therefore bal-
086 ance classification accuracy with computational
087 efficiency (Liu et al., 2024a).

088 To address these challenges, we (i) decompose in-
089 tent understanding into three dimensions—*Domain*,
090 *Intent*, and *Issue*—so that adding new domains or
091 intents expands only the relevant dimension. This
092 reduces the effective label space and supports com-
093 pound intent modeling by separating actionable
094 intent types (e.g., “Cancel”, “Informational In-
095 quiry”) from issue attributes (e.g., “Unrecognized
096 Charge”, “Sync/Download”). We then (ii) fine-
097 tune lightweight models leveraging various LLMs-
098 based methods to generate labels from customer
099 agent transcripts at scale. The model processes
100 the concatenation of the dialogue history and the
101 current turn, enabling context-aware, turn-level in-
102 tent detection. Finally, (iii) we augment training
103 data with LLM-simulated multi-turn dialogues that
104 inject topic and intent switches—patterns that are
105 rare in transcripts but critical for robustness.

106 2 Related work

107 2.1 Multi-turn intent classification

108 Multi-turn intent classification incorporates dia-
109 logue context to improve intent prediction. Prior
110 work uses contextual encoders for customer ser-
111 vice intent detection and models cross-turn depen-
112 dencies via hierarchical or graph-based structures
113 (Wang et al., 2021; Senese et al., 2020; Liu and
114 Chen, 2019; Qin et al., 2021). We follow this line
115 but focus on data-driven robustness to context car-
116 rryover and intent switching via controllable simu-
117 lation.

118 LLMs have also been explored for intent detec-
119 tion through prompting and hybrid routing (Arora
120 et al., 2024), as well as retrieval-augmented or
121 demonstration-based pipelines for few-shot intent
122 prediction (Yu et al., 2021; Zhang et al., 2025; Liu
123 et al., 2024b). In contrast, we primarily use LLMs
124 for automatic labeling and data generation to train
125 lightweight models, with optional LLM fallback at
126 inference time.

127 2.2 User simulation

128 User simulation is widely used for synthetic dia-
129 logue generation and system evaluation, including
130 agenda-based and neural approaches (Schatzmann
131 et al., 2007; Lin et al., 2021; Sun et al., 2022). Re-

cent work shows LLMs can act as user simulators
(Balog and Zhai, 2025; Balog et al., 2025). Unlike
simulators aimed at general task completion, our
simulator targets real-world multi-turn behaviors—
especially intent switching—using control to main-
tain coherence and constrain generation.

132 2.3 Chain-of-Thought and self-critique 133 reasoning 134 135 136 137

Chain-of-Thought and self-refinement improve
LLM outputs via intermediate reasoning and it-
erative revision (Wei et al., 2022; Kojima et al.,
2022; Madaan et al., 2023; Shinn et al., 2023).
We apply these ideas to transcript labeling (rather
than inference-time in-context learning) to gener-
ate higher-quality training data for low-latency intent
models. Compact transformer students further sup-
port efficient deployment (Sanh et al., 2019; Jiao
et al., 2020; Wang et al., 2020).

138 3 Preliminary 139 140 141 142 143 144 145 146 147 148 149

We formulate multi-turn intent classification as a
multi-class text classification problem. Given a
conversation history \mathcal{C} and the customer’s current
utterance q , the objective is to predict an intent label
 t from a predefined set $T = \{t_1, \dots, t_k\}$ using a
model \mathcal{M} . The predicted intent \hat{t} is obtained by
maximizing the posterior probability:

$$\hat{t} = \arg \max_{t \in T} P(t | q, \mathcal{C}; \theta), \quad (1)$$

where $P(t | q, \mathcal{C}; \theta)$ denotes the probability of in-
tent t conditioned on (q, \mathcal{C}) , parameterized by θ of
 \mathcal{M} . As mentioned in Section 1, incorporating both
 q and \mathcal{C} in our model is essential because of con-
text carryover and intent switching in multi-turn
settings.

The number of intents k grow rapidly as product
lines expand. To address the scalability issue, we
propose an ontology that captures complementary
aspects of intent understanding (Figure 1): **Do-**
main, representing product categories; **Intent**, dis-
tinguishing conversational intents from actionable
intents; and **Issue**, representing slot-level attributes
associated with each intent. This decomposition
reduces classification complexity as well as sup-
ports **compound intent modeling** naturally by de-
coupling intent identification from issue detection,
allowing the system to jointly predict an intent and
its associated issue.

We train separate models to classify *Domain*, *In-*
tent, and *Issue*. Their outputs are then com-

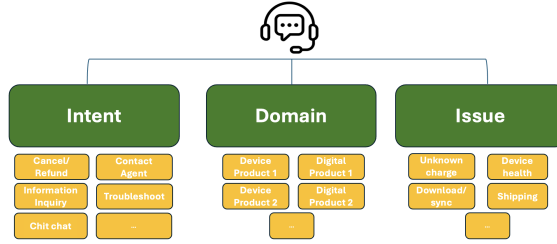


Figure 1: Intent classification ontology.

posed and mapped to downstream actions, including workflows, LLM-based agents, or human-agent hand-offs (Figure 2). Considering the latency requirement, we leverage a foundation model, Claude 3.7 Sonnet¹, to generate high-quality labeled training data. In the following sections, we introduce and evaluate multiple LLM-based labeling and data augmentation strategies that leverage existing transcript data to finetune our models under this ontology.

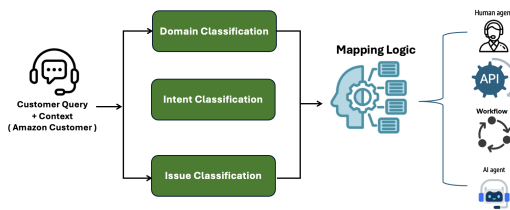


Figure 2: Intent detection system architecture

4 Method

4.1 Labeling strategy for intent classification

In real conversations, either the agent or the customer may consecutively input several utterances. However, the dialogue between a customer and a chatbot is typically conducted in an alternating manner. To construct data that aligns with the chatbot format, we merge consecutive utterances from the same role into a single unit, resulting in a dialogue $d = [q_1, a_1, \dots, q_n, a_n]$ where q_i represents the user’s current query and a_i represents the agent’s response. For a query q_i , we define its context or conversation history as $C_i = [q_1, a_1, \dots, q_{i-1}, a_{i-1}]$. In the following, we propose different label generation strategies for building our intent classification models.

Single-stage reasoning-guided labeling

¹<https://www.anthropic.com/news/claude-3-7-sonnet>

We prompt the LLM to infer the customer’s intent given the current user query q_i and its associated conversation history C_i . The prompt provides the taxonomy along with descriptions of each label, and instructs the LLM to output a single label accompanied by a free-form explanation. This explanation cites the key evidence in both the current query q_i and the context C_i that supports the model’s decision. We hypothesize that asking the LLM to articulate its reasoning and supporting evidence leads to more accurate intent classification by itself.

Two-stage labeling with self-critique

We adopt a two-stage labeling strategy that decomposes intent annotation into an *initial prediction* followed by an explicit *self-critique and revision*. **Stage 1** applies the single-stage reasoning-guided labeling procedure described in the previous section, producing an intent label with a brief explanation. In **Stage 2**, the LLM is given the original input together with the Stage-1 prediction and rationale, and is prompted to act as a critic. It assesses consistency with the dialogue evidence and intent taxonomy in Stage-1, flags failures (e.g., reliance on spurious keywords, missed contextual signals in C_i , confusion between closely related intents, or hallucinated assumptions), and then either confirms the Stage-1 label or revises it with a short justification. Stage 2 must explicitly indicate whether the label is *kept* or *revised*; when revised, it must cite minimal supporting evidence (e.g., a specific utterance in C_i or phrase in q_i) motivating the correction. We hypothesize that this two-stage self-critique improves accuracy by correcting errors introduced during initial reasoning.

To evaluate our hypothesis, we manually reviewed more than 5,500 randomly sampled test instances and report the results in Table 1. The single-stage and two-stage strategies produced identical annotations for approximately 87% of samples, and these were largely accurate. In the remaining 13% of

cases where the strategies disagreed, the two-stage approach was correct in the majority of instances, suggesting that the second-stage reassessment reduces hallucinations and overall annotation errors.

4.2 Data augmentation strategy for intent classification

Although our e-commerce platform provides access to millions of real customer service transcripts between agent and customers that can be leveraged to train models, these transcripts are typically single-topic and follow a largely linear progression. In contrast, customer–bot interactions are potentially dynamic: users may switch intents or change goals in a single session. Consequently, models trained solely on real transcripts often fail to generalize to these complex patterns.

Data generation using dialogue simulator: To bridge this gap, we develop a multi-turn dialogue generation framework that simulates realistic customer–bot interactions. Our framework uses an LLM-based user query simulator that interacts with an LLM-based response simulator, while a *Simulator Controller* orchestrates the conversation by initializing dialogues, managing turn-by-turn flow, and introducing topic/intent shifts when appropriate. The controller combines intent-level planning with controlled randomness to elicit underrepresented behaviors in real transcripts (e.g., chit-chat, intent changes, follow-up and clarification questions). At each turn, it selects the next user intent and prompts the user simulator to produce the corresponding utterance, after which the response simulator generates the bot reply. We also inject alternative seed queries sampled from a multi-domain top-query database to facilitate topic and intent transitions. The overall interaction flow is shown in Figure 3 and simulation prompt template can be found in Appendix C.3. Finally, we use the two stage labeling strategy in Section 4.1 to automatically label simulated dialogues.

4.3 Hybrid approach of LLM and fine-tuned model

To demonstrate how a fine-tuned lightweight model can be combined with an LLM to balance accuracy and latency, we introduce a *hybrid intent detection strategy* that couples fast local classification with selective LLM escalation. As shown in Appendix C.5 and Figure 4, the lightweight model first outputs an intent probability distribution along with a confidence score. If the confidence score exceeds

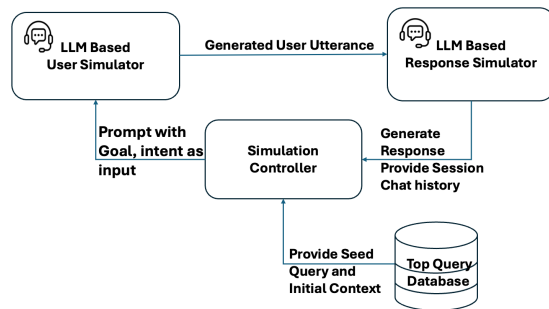


Figure 3: Multi-turn conversation simulator

a threshold τ , we accept the lightweight model’s prediction. Otherwise, we extract the model’s top K intent candidates and invoke an LLM to perform *constrained* disambiguation over this candidate set. This hybrid design routes most intent detection requests through the low-latency lightweight model, while reserving LLM calls for a small subset of ambiguous cases. The choice of τ and K is motivated by the accuracy vs coverage analyses in Appendix C.6 (Figures 5 and 6).

5 Experimental setup

5.1 Datasets

We apply the anonymization procedure described in Appendix A. Each sample consists of conversation history, the current customer utterance, and its corresponding annotations for intent, issue, and product category. To build a balanced dataset, we use stratified sampling across months and domains. Since some categories are significantly overrepresented, we further downsample high-frequency categories when forming the final training and evaluation splits. We provide additional dataset construction details are provided in Appendix B.

5.2 Evaluation benchmarks

We consider the following baselines to evaluate and benchmark our LLM-based labeling and data augmentation methods:

- **RoBERTa:** We fine-tune RoBERTa-base (Liu et al., 2019) with multi-head classification components (*Intent*, *Domain*, and *Issue*), enabling simultaneous prediction across multiple label spaces. To compare prompting and data-generation choices for fine-tuning (Section 4), we train variants using: (i) single-stage reasoning-guided labeling, (ii) two-stage labeling with self-critique, and (iii) two-stage labeling with

Outcome	Count
Stage 2 is correct	566
Stage 1 is correct	84
Both Stage 1 & Stage 2 are incorrect	70
Both Stage 1 & Stage 2 are correct	4802
Total	5522

Outcome	Accuracy
Stage 2 is correct	97%
Stage 1 is correct	88%

Table 1: Stage 1 and Stage 2 performance analysis

self-critique plus dialogue-simulator augmentation. We provide implementation details in appendix C.4.

- **Claude Sonnet 3.7 zero-shot:** We evaluate Claude Sonnet 3.7 in a zero-shot setting using a prompt that enumerates the intent taxonomy and provides brief definitions for each label.
- **Nova Pro zero-shot:** We mirror the **Claude Sonnet 3.7 Zero-shot** setup but replace Claude with Nova Pro, a smaller LLM that offers lower inference latency.
- **Hybrid approach (fine-tuned RoBERTa + Claude Sonnet 3.7):** Following Section 4.3, we first use the best-performing fine-tuned RoBERTa model to retrieve the top-3 candidate labels for each classifier (*Intent, Domain, Issue*). We then prompt Claude Sonnet 3.7 to select the final label conditioned on these candidates. This hybrid baseline tests whether zero-shot prompting can further improve over RoBERTa fine-tuning.

6 Results

6.1 Automatic evaluation on fine-tuned RoBERTa models

As described in Sections 4.1 and 4.2, we employ multiple data labeling and augmentation strategies to generate training data for fine-tuning RoBERTa models. Specifically, we consider three approaches for comparison: (1) a single-stage prompting strategy, (2) a two-stage prompting strategy, and (3) a two-stage prompting strategy + simulated dialogue data augmentation. We generate ground-truth labels using an LLM according to each strategy. We report automatic evaluation results in Table 2. Models trained under all three strategies achieve comparable performance across classification tasks, indicating that the fine-tuned models are able to effectively learn from their corresponding LLM teacher models. Notably, models fine-tuned using the two-stage prompting strategy exhibit substan-

Model Type	Single-Stage Labeling	Two-Stage Labeling	Two-Stage Labeling + Data Augmentation
Intent Model	80.70%	80.30%	79.50%
Product Model	80.90%	80.80%	81.70%
Issue Model	76.70%	79.10%	78.20%

Table 2: Automated evaluation accuracy results

tial performance gains on the issue classification task. We hypothesize that issue classification involves greater ambiguity and is therefore more challenging. The two-stage prompting strategy, which encourages additional self-verification and refinement, produces higher-quality labels and improves fine-tuning performance.

6.2 Human evaluation

To comprehensively evaluate the effectiveness of our labeling and data augmentation strategies, as well as the benefits of the hybrid approach, we compare across several baseline approaches (see Section 5.2) using a 1,206 human-annotated dataset. The dataset is constructed by random sampling, with 50% from the evaluation dataset described in Section 5.1 and the remaining 50% from production traffic. The results, presented in Table 3, illustrate performance across all 3 types of classification models and an overall aggregate assessment based on Precision, Recall, and F1-score.

Overall, our results show that models fine-tuned with the two-stage labeling strategy and augmented with simulated dialogue data outperform alternative approaches. While training with data generated via single-stage prompting already yields strong performance, adding self-critique stage further improves labeling accuracy and model performance. Moreover, incorporating simulated multi-turn dialogues helps the model better handle intent switches, contributing to additional gains (F1 score increases by 2%-5% compared to single-stage prompting). Finally, we observe that intent classification perfor-

	Intent Model			Product Model			Issue Model		
	Precision	Recall	F1 Score	Precision	Recall	F1 Score	Precision	Recall	F1 Score
RoBERTa(Single-stage labeling)	73%	71%	72%	73%	65%	69%	67%	61%	64%
RoBERTa(Two-stage labeling)	74%	73%	73%	72%	69%	70%	72%	64%	68%
RoBERTa(Two-stage labeling+data augmentation)	75%	74%	74%	75%	72%	73%	73%	66%	69%
Hybrid approach	77%	75%	76%	78%	74%	76%	74%	68%	71%
Nova Pro zero-Shot	48%	42%	45%	61%	52%	56%	42%	35%	38%
Claude Sonnet 3.7 zero-shot	75%	67%	71%	61%	69%	65%	54%	51%	52%

Table 3: Performance comparison across different baselines for intent, product, and issue classification

mance is further improved by a hybrid approach that combines a fine-tuned RoBERTa model with LLM-based prompting.

Purely prompting-based methods (e.g., Claude zero-shot and Nova zero-shot) exhibit inferior performance. This degradation likely stems from the absence of domain-specific customer service training data, which increases the likelihood of hallucinations and intent misclassification. Although the hybrid approach also employs prompting when model confidence is low, it substantially outperforms standalone prompting. This gain stems from using the fine-tuned RoBERTa model to first retrieve a limited set of relevant intent candidates, thereby constraining the LLM’s decision space and enabling more accurate final predictions.

We also evaluate P50 and P90 latency across different methods, with results summarized in Table 4. The findings indicate that RoBERTa achieves the

Model	P50 Latency	P90 Latency
RoBERTa	0.08s	0.1s
Nova Pro Zero-shot	1.98s	3.97s
Claude Sonnet 3.7 Zero-Shot	5.93s	7.52s
Hybrid (RoBERTa + Zero-shot)	1.99s	3.1s

Table 4: Model latency performance comparison

lowest latency, approximately 20–50× faster than LLM zero-shot. While the hybrid approach delivers improved accuracy, it incurs higher latency than RoBERTa, suggesting a trade-off between performance gains and inference efficiency.

6.3 Online deployment performance

We deployed the intent classification model (fine-tuned RoBERTa-base) in an e-commerce customer service production systems. Compared to the previously deployed single-turn intent detection system, which could not support topic shifts or context carryover in an ongoing sessions, our model enables seamless and dynamic intent routing in multi-turn interactions. In a one-month online A/B test, our

model increased the bot automation rate by 4.91% and improved the positive customer response rate by 7.89%, demonstrating benefits for both customer experience and operational efficiency while achieving low end-to-end production latency (P50: 0.12 s; P90: 0.16 s; P99: 0.20 s).

7 Conclusion

In this paper, we address multi-turn intent detection for customer service applications. To handle the scalability and heterogeneity of intent taxonomies, we propose an ontology that captures complementary facets of user intent. To mitigate the scarcity of annotated data, we introduce LLM-based labeling methods that generate high-quality supervision from existing customer transcripts, and augment training with LLM-simulated multi-turn dialogues that explicitly model topic shifts and intent switches.

Our experiments show that two-stage labeling with self-critique, combined with simulated dialogue augmentation, consistently outperforms alternative labeling strategies. The resulting fine-tuned RoBERTa models outperform pre-trained LLMs in zero-shot settings while achieving substantially lower latency. A hybrid routing strategy that combines fine-tuned RoBERTa with an LLM further improves performance on ambiguous cases.

Our findings provide actionable guidance for practitioners building production multi-turn intent detection systems by effectively combining real transcript data with LLM-generated dialogues.

8 Limitations

Our evaluation is based on customer service conversations from a specific set of products, locales, and workflow designs, so results may not fully generalize to other domains with different intent taxonomies or dialogue patterns. In addition, model quality depends on the consistency of upstream labels (human or LLM-assisted); any ambiguity in intent definitions or noise in automatic labeling can

479	propagate to training and inflate offline estimates.	Yinhan Liu, Myle Ott, Naman Goyal, Jingfei Du, Man-	532
480	Finally, offline metrics may not perfectly translate	dar Joshi, Danqi Chen, Omer Levy, Mike Lewis, Luke	533
481	to end-to-end customer impact because real de-	Zettlemoyer, and Veselin Stoyanov. 2019. Roberta: A	534
482	ploysments involve additional constraints (policy,	robustly optimized bert pretraining approach. <i>arXiv</i>	535
483	UI, latency, and fallback behavior) and are sub-	<i>preprint arXiv:1907.11692</i> .	536
484	ject to distribution shift over time; broader cross-	Zhengyuan Liu and Nancy Chen. 2019. Reading turn	537
485	domain/locale testing and controlled online stud-	by turn: Hierarchical attention architecture for spoken	538
486	ies are needed to validate robustness and user out-	dialogue comprehension. In <i>Proceedings of the 57th</i>	539
487	comes.	<i>Annual Meeting of the Association for Computational</i>	540
		<i>Linguistics</i> .	541
488	References	Zihan Liu, Yiming Chen, Hao Zhang, et al. 2024b. Lara:	542
489	Gaurav Arora, Shreya Jain, and Srujana Merugu. 2024.	Linguistic-adaptive retrieval-augmentation for multi-	543
490	Intent detection in the age of LLMs. In <i>Proceedings of</i>	turn intent classification. In <i>Proceedings of the 2024</i>	544
491	<i>the 2024 Conference on Empirical Methods in Natural</i>	<i>Conference on Empirical Methods in Natural Language</i>	545
492	<i>Language Processing: Industry Track</i> .	<i>Processing</i> .	546
493	Krisztian Balog, Nolwenn Bernard, Saber Zerhoubi, and	Aman Madaan, Niket Tandon, Prakhar Gupta, Skyler	547
494	ChengXiang Zhai. 2025. Theory and toolkits for user	Hallinan, Luyu Gao, Sarah Wiegrefe, Nanyun Raja,	548
495	simulation in the era of generative AI: User modeling,	Shivang Gulati, Shubham Tan, et al. 2023. Self-refine:	549
496	synthetic data generation, and system evaluation.	Iterative refinement with self-feedback. <i>arXiv preprint</i>	550
497	In <i>Proceedings of the 48th International ACM SIGIR Con-</i>	<i>arXiv:2303.17651</i> .	551
498	<i>ference on Research and Development in Information</i>	Subhadip Nandi, Neeraj Agrawal, Anshika Singh, and	552
499	<i>Retrieval (SIGIR '25)</i> , pages 4138–4141, Padua, Italy.	Priyanka Bhatt. 2024. Enhancing customer service chat-	553
500	Association for Computing Machinery.	bots with context-aware nlu through selective attention	554
501	Krisztian Balog and ChengXiang Zhai. 2025. User	and multi-task learning. In <i>Proceedings of the 8th Inter-</i>	555
502	simulation in the era of generative AI: User modeling,	<i>national Conference on Data Science and Management</i>	556
503	synthetic data generation, and system evaluation. <i>arXiv</i>	<i>of Data (CODS-COMAD 2024)</i> , pages 220–228. Asso-	557
504	<i>preprint arXiv:2501.04410</i> .	ciation for Computing Machinery.	558
505	Lei Cui, Shaohan Huang, Furu Wei, Chuanqi Tan, Chao-	Haode Qi, Lin Pan, Atin Sood, Abhishek Shah, Ladislav	559
506	qun Duan, and Ming Zhou. 2017. Superagent: A cus-	Kunc, Mo Yu, and Saloni Potdar. 2021. Benchmark-	560
507	tommer service chatbot for e-commerce websites. In <i>Pro-</i>	ing commercial intent detection services with practice-	561
508	<i>ceedings of ACL 2017, system demonstrations</i> , pages	driven evaluations. In <i>Proceedings of the 2021 Con-</i>	562
509	97–102.	<i>ference of the North American Chapter of the Associa-</i>	563
510	Xiaoqi Jiao, Yichun Yin, Lifeng Shang, Xin Jiang, Xiao	<i>tion for Computational Linguistics: Human Language</i>	564
511	Chen, Linlin Li, Fang Wang, and Qun Liu. 2020. Tiny-	<i>Technologies: Industry Papers</i> , pages 304–310, Online.	565
512	BERT: Distilling BERT for natural language understand-	Association for Computational Linguistics.	566
513	ing. In <i>Findings of the Association for Computational</i>	Libo Qin, Zhou Chen, Wanxiang Che, Hang Li, and	567
514	<i>Linguistics: EMNLP 2020</i> , pages 4163–4174, Online.	Ting Liu. 2021. Knowing where to leverage: Context-	568
515	Association for Computational Linguistics.	aware graph convolutional network with an adaptive	569
516	Takeshi Kojima, Shixiang Shane Gu, Machel Reid, Yu-	fusion layer for contextual spoken language understand-	570
517	taka Matsuo, and Yusuke Iwasawa. 2022. Large lan-	ing.	571
518	guage models are zero-shot reasoners. <i>arXiv preprint</i>	Victor Sanh, Lysandre Debut, Julien Chaumond, and	572
519	<i>arXiv:2205.11916</i> .	Thomas Wolf. 2019. DistilBERT, a distilled version	573
520	Hsien-chin Lin, Nurul Lubis, Songbo Hu, Carel van	of BERT: smaller, faster, cheaper and lighter. <i>arXiv</i>	574
521	Niekerk, Christian Geishauer, Michael Heck, Shutong	<i>preprint arXiv:1910.01108</i> .	575
522	Feng, and Milica Gasic. 2021. Domain-independent	Jost Schatzmann, Blaise Thomson, Karl Weilhammer,	576
523	user simulation with transformers for task-oriented di-	Hui Ye, and Steve Young. 2007. Agenda-based user	577
524	alogue systems. In <i>Proceedings of the 22nd Annual</i>	simulation for bootstrapping a POMDP dialogue sys-	578
525	<i>Meeting of the Special Interest Group on Discourse and</i>	tem. In <i>Human Language Technologies 2007: The</i>	579
526	<i>Dialogue</i> , pages 445–456, Singapore and Online. Asso-	<i>Conference of the North American Chapter of the As-</i>	580
527	ciation for Computational Linguistics.	<i>sociation for Computational Linguistics; Companion</i>	581
528	Junhua Liu, Yong Keat Tan, Bin Fu, and Kwan Hui Lim.	<i>Volume, Short Papers</i> , pages 149–152, Rochester, New	582
529	2024a. Balancing accuracy and efficiency in multi-turn	York. Association for Computational Linguistics.	583
530	intent classification for LLM-powered dialog systems	Matteo Antonio Senese, Giuseppe Rizzo, Mauro Drag-	584
531	in production.	oni, and Maurizio Morisio. 2020. MTSI-BERT: A	585
		session-aware knowledge-based conversational agent.	586
		In <i>Proceedings of the Twelfth Language Resources</i>	587

588 *and Evaluation Conference*, pages 717–725, Marseille,
589 France. European Language Resources Association.

590 Noah Shinn, Federico Cassano, Edward Berman, Ash-
591 win Gopinath, Karthik Narasimhan, and Shunyu Yao.
592 2023. Reflexion: Language agents with verbal rein-
593 forcement learning. *arXiv preprint arXiv:2303.11366*.

594 Weiwei Sun, Shuyu Guo, Shuo Zhang, Pengjie Ren,
595 Zhumin Chen, Maarten de Rijke, and Zhaochun Ren.
596 2022. Metaphorical user simulators for evaluat-
597 ing task-oriented dialogue systems. *arXiv preprint*
598 *arXiv:2204.00763*.

599 Peiyao Wang, Joyce Fang, and Julia Reinspach. 2021.
600 [CS-BERT: a pretrained model for customer service dia-](#)
601 [logues](#). In *Proceedings of the 3rd Workshop on Natural*
602 *Language Processing for Conversational AI*, pages 130–
603 142, Online. Association for Computational Linguistics.

604 Wenhui Wang, Furu Wei, Li Dong, Hangbo Bao, Nan
605 Yang, and Ming Zhou. 2020. MiniLM: Deep self-
606 attention distillation for task-agnostic compression of
607 pre-trained transformers. In *Advances in Neural Infor-*
608 *mation Processing Systems*.

609 Jason Wei, Xuezhi Wang, Dale Schuurmans, Maarten
610 Bosma, Brian Ichter, Fei Xia, Ed H. Chi, Quoc V. Le,
611 and Denny Zhou. 2022. Chain-of-thought prompting
612 elicits reasoning in large language models. In *Advances*
613 *in Neural Information Processing Systems*.

614 Ting-Wei Wu, Ruolin Su, and Biing-Hwang Juang.
615 2021. [A context-aware hierarchical BERT fusion net-](#)
616 [work for multi-turn dialog act detection](#). In *Proceedings*
617 *of Interspeech 2021*, pages 1239–1243.

618 Dian Yu, Luheng He, Yuan Zhang, Xinya Du, Panupong
619 Pasupat, and Qi Li. 2021. [Few-shot intent classification](#)
620 [and slot filling with retrieved examples](#). In *Proceedings*
621 *of the 2021 Conference of the North American Chapter*
622 *of the Association for Computational Linguistics: Hu-*
623 *man Language Technologies*, pages 734–749, Online.
624 Association for Computational Linguistics.

625 Ziji Zhang, Michael Yang, Zhiyu Chen, Yingying
626 Zhuang, Shu-Ting Pi, Qun Liu, Rajashekar Maragoud,
627 Vy Nguyen, and Anurag Beniwal. 2025. [REIC: RAG-](#)
628 [enhanced intent classification at scale](#). In *Proceedings*
629 *of the 2025 Conference on Empirical Methods in Natu-*
630 *ral Language Processing: Industry Track*, pages 1072–
631 1080, Suzhou (China). Association for Computational
632 Linguistics.

633 Wayne Xin Zhao, Kun Zhou, Junyi Li, Tianyi Tang,
634 Xiaolei Wang, Yupeng Hou, Yingqian Min, Beichen
635 Zhang, Junjie Zhang, Zican Dong, Yifan Du, Chen Yang,
636 Yushuo Chen, Zhipeng Chen, Jinhao Jiang, Ruiyang
637 Ren, Yifan Li, Xinyu Tang, Zikang Liu, Peiyu Liu, Jian-
638 Yun Nie, and Ji-Rong Wen. 2023. [A survey of large](#)
639 [language models](#).

640 Yunhua Zhou, Jiawei Hong, and Xipeng Qiu. 2023. [To-](#)
641 [wards open environment intent prediction](#). In *Findings*
642 *of the Association for Computational Linguistics: ACL*
643 *2023*, pages 2226–2240, Toronto, Canada. Association
644 for Computational Linguistics.


```

701 19. NON_AMAZON_TOPIC - Completely unrelated to Amazon products or services
702
703 ## Issue Options (Choose Exactly One)
704
705 1. SETUP/REGISTRATION/ACTIVATION/INSTALLATION
706 2. CANCEL/REFUND/RETURN/REPLACEMENT/RETAIN/SUBSCRIBE/TRADE-IN/TRANSACTIONAL
707 3. CONNECT/PAIR
708 4. WIFI/NETWORK
709 5. SYNC/DOWNLOAD
710 6. PLAY/STREAMING/DISPLAY
711 7. AUDIO/SOUND
712 8. DAMAGE/REPAIR
713 9. WARRANTY
714 10. DEVICE_HEALTH
715 11. DEFECT
716 12. BATTERY
717 13. SOFTWARE/OTA_UPDATE
718 14. LOST/STOLEN
719 15. PROMOTIONS/CREDITS/LOOT
720 16. RESTRICTIONS/PARENTAL_CONTROL
721 17. HOUSEHOLD
722 18. PAYMENT_ISSUE/PURCHASE/COINS
723 19. PRICING
724 20. UNRECOGNIZED_CHARGES/UNKNOWN_CHARGES/FRAUD_CHARGES
725 21. CONTENT_AVAILABILITY
726 22. SHIPPING/DELIVERY
727 23. OTHER_ISSUE
728 24. NONE
729
730 ## Product Options (Choose Exactly One)
731
732 1. Music
733 2. Video
734 3. SmartTV
735
736 [... additional product categories ...]
737
738 ## Output Format (JSON)
739 {
740     "predicted_conversational_intent": "One of the 19 intent labels",
741     "predicted_issue": "One of the 24 issue labels",
742     "predicted_product": "One of the 31 product labels",
743     "reason": "Brief explanation justifying the choice"
744 }
745

```

Classification prompt for stage 2:

```

746 ## Taxonomy Intents
747
748 1. CANCEL - Request to cancel subscription/service/product or turn off
749    auto renewal. Excludes refunds.
750 2. REFUND - Request for refund related to service/subscription.
751 3. CANCEL&REFUND - Explicitly cancel and request refund.
752 4. RETURN - Request for return related to service/product.
753 5. REPLACEMENT - Request for replacement related to service/product.
754 6. RETAIN/SUBSCRIBE - Request to retain/subscribe, NOT cancel.
755 7. TRADE-IN - Customer requesting to trade-in devices.
756 8. ADS-FREE/ADS-REMOVAL - Request to remove ads or subscribe ads free.
757 9. TROUBLESHOOT - Customer describes issue requiring troubleshooting.
758 10. INFORMATIONAL_INQUIRY - General "how-to" or "what-is" questions
759     answerable from help pages. Excludes vague utterances like "I need
760     help".
761 11. ISSUE_DEPENDENT_INQUIRY - Inquiries tied to specific issue/service/
762     product requiring customer context. Exclude transactional intents.
763 12. Yes - Customer responded affirmatively. Don't confuse with
764     END_CONVERSATION.
765 13. No - Customer responds "No" or "Nope". Don't confuse with
766     END_CONVERSATION.
767 14. CHIT_CHAT/CONVERSATION_FILLER - Polite fillers (hi, thanks),
768

```

navigation phrases. Exclude frustration/complaint and Yes/No.	769
15. FRUSTRATION/COMPLAINT - General dissatisfaction without clear request.	770
16. REQUEST_HUMAN_AGENT - Customer explicitly asks for human agent.	771
17. END_CONVERSATION - Customer states issue resolved. Don't confuse with polite chit-chat.	772
18. NON_RELATED_TOPIC - Completely unrelated to the e-commerce. Must be completely off-topic.	773
	774
	775
	776
## Taxonomy Issues	777
	778
1. SETUP/REGISTRATION/ACTIVATION/INSTALLATION - Issues registering, setting up, activating, or installing device/subscription/service/app.	779
2. CANCEL/REFUND/RETURN/REPLACEMENT/RETAIN/SUBSCRIBE/TRADE-IN/ TRANSACTIONAL - Issue with transaction already occurred or attempted.	780
3. CONNECT/PAIR - Issues connecting or bluetooth pairing.	781
4. WIFI/NETWORK - Issues with networking or Wifi.	782
5. SYNC/DOWNLOAD - Trouble syncing/downloading (often eBook/App related).	783
6. PLAY/STREAMING/DISPLAY - Trouble playing content or visualizing.	784
7. AUDIO/SOUND - Audio issues (no audio, low audio, etc.).	785
8. DAMAGE/REPAIR - Device damage, inquiring about repair.	786
9. WARRANTY - Inquiring about or checking warranty.	787
10. DEVICE_HEALTH - Device performance/stability issues (crashes, reboots, bricked, frozen, responsiveness).	788
11. DEFECT - General mention of device defect.	789
12. BATTERY - Battery issues.	790
13. SOFTWARE/OTA_UPDATE - Software update issues (stuck on update screen, no update available).	791
14. LOST/STOLEN - Reporting lost/stolen device.	792
15. PROMOTIONS/CREDITS/LOOT - Issues with promotions, bundles, credits, loot.	793
16. RESTRICTIONS/PARENTAL_CONTROL - Trouble with Parental Controls, Pins, child purchases.	794
17. HOUSEHOLD - Issues with household (sharing content, adding members). If PINS/Child Controls, select Parental Controls instead.	795
18. PAYMENT_ISSUE/PURCHASE/COINS - Issue/inquiry related to payment, gift cards, or coins.	796
19. PRICING - Price related issues (price match, price adjustment).	797
20. UNRECOGNIZED_CHARGES/UNKNOWN_CHARGES/FRAUD_CHARGES - Generic unrecognized charges customer not aware of or deems fraud.	798
21. CONTENT_AVAILABILITY - Issues with content availability.	799
22. SHIPPING/DELIVERY - Shipping or delivery related issues.	800
23. OTHER_ISSUE - Other action-oriented/transactional inquiries not mentioned above.	801
24. NONE - No specific issue, purely chit-chat, complaint or frustration.	802
	813
## Taxonomy Products	814
	815
1. Music.	816
2. Video	817
3. Video Channels	818
4. SmartTV Cube	819
[... additional product categories ...]	820
	821
<hr/>	
Classification prompt template for LLM zero-shot:	822
## Dialogue Context	823
{dialogue_context}	824
	825
## Turn T - Current Customer Utterance	826
{turn_T_utterance}	827
	828
## Conversational/Actional Intent Options (Choose Exactly One)	829
	830
1. CANCEL - Request to cancel a subscription or service or product	831
2. REFUND - Request for a refund related to a service or subscription	832
3. CANCEL&REFUND - Explicitly state desire to cancel and request a refund	833
4. RETURN - Request for a return related to a service or product	834
5. REPLACEMENT - Request for a replacement related to a service or product	835
6. RETAIN/SUBSCRIBE - Request to retain/subscribe subscription and NOT to cancel	836

837 7. TRADE-IN - Customer is requesting to trade-in devices
838 8. ADS-FREE/ADS-REMOVAL - Customer is requesting to remove ads or subscribe ads free
839 9. TROUBLESHOOT - Customer describes an issue requiring troubleshooting
840 10. INFORMATIONAL_INQUIRY - General "how-to" or "what-is" questions
841 11. ISSUE_DEPENDENT_INQUIRY - Inquiries tied to a specific issue or service/product
842 12. Yes - The customer responded affirmatively
843 13. No - Customer generally respond with "No" or "Nope"
844 14. TRANSACTIONAL_OTHER - Other action-oriented or transactional inquiries
845 15. CHIT_CHAT/CONVERSATION_FILLER - Polite fillers, navigation phrases
846 16. FRUSTRATION/COMPLAINT - General dissatisfaction expressed
847 17. REQUEST_HUMAN_AGENT - Customer explicitly asks for a human agent
848 18. END_CONVERSATION - Customer clearly states issue is resolved
849 19. NON_AMAZON_TOPIC - Completely unrelated to Amazon products or services

850

851 ## Issue Options (Choose Exactly One)

852

853 1. SETUP/REGISTRATION/ACTIVATION/INSTALLATION
854 2. CANCEL/REFUND/RETURN/REPLACEMENT/RETAIN/SUBSCRIBE/TRADE-IN/TRANSACTIONAL
855 3. CONNECT/PAIR
856 4. WIFI/NETWORK
857 5. SYNC/DOWNLOAD
858 6. PLAY/STREAMING/DISPLAY
859 7. AUDIO/SOUND
860 8. DAMAGE/REPAIR
861 9. WARRANTY
862 10. DEVICE_HEALTH
863 11. DEFECT
864 12. BATTERY
865 13. SOFTWARE/OTA_UPDATE
866 14. LOST/STOLEN
867 15. PROMOTIONS/CREDITS/LOOT
868 16. RESTRICTIONS/PARENTAL_CONTROL
869 17. HOUSEHOLD
870 18. PAYMENT_ISSUE/PURCHASE/COINS
871 19. PRICING
872 20. UNRECOGNIZED_CHARGES/UNKNOWN_CHARGES/FRAUD_CHARGES
873 21. CONTENT_AVAILABILITY
874 22. SHIPPING/DELIVERY
875 23. OTHER_ISSUE
876 24. NONE

877

878 ## Product Options (Choose Exactly One)

879

880 1. Music
881 2. Video
882 3. SmartTV

883

884 [... additional product categories ...]

885

886 ## Output Format (JSON)

887

```
888 {
889     "predicted_conversational_intent": "One of the 19 intent labels",
890     "predicted_issue": "One of the 24 issue labels",
891     "predicted_product": "One of the 31 product labels",
892 }
```

892

C.2 System prompt

893

894 We define three separate prompt templates corresponding to Stage 1, Stage 2 annotation and
895 LLM zero-shot. Note LLM zero-shot and Stage 1 annotation shares the same system prompt.

896

Stage 1 or LLM zero-shot system prompt:

897

898 You are an expert annotation assistant specializing in analyzing conversations
899 between customers and bots/agents. Your task is to classify each customer message
900 (Turn T) into its primary intent and the most relevant Amazon product or
901 service discussed.

902

903 Use the provided dialogue history for context, and ensure that classifications
adhere strictly to the predefined categories. Always select exactly one intent

and one product for each message, even if the product is inferred from the context. 904

If the product is ambiguous but likely Amazon-related, choose 'Other'. If the message is unrelated to Amazon, select 'NON_AMAZON_TOPIC' as the intent. Provide clear reasoning for your classifications, referencing specific dialogue cues and your decision-making process. 905
906
907
908
909
910

Stage 2 system prompt: 911

You are a strict annotation reviewer. Your job is to AUDIT a prior classification (Stage 1) for a customer's Turn T, using the SAME taxonomy as Stage 1. 912
913

Goals: 914

- Verify that the predicted intent, issue, and product each match their definitions. 915

- Challenge the original reasoning (try to find contradictions or missing evidence). 916

- Correct any mistakes; otherwise confirm the original labels. 917

- Be conservative with ambiguous cases: only use NON_AMAZON_TOPIC when clearly 918

unrelated to Amazon; do not confuse CHIT_CHAT with END_CONVERSATION; do not confuse 919

YES/NO with functional intents. 920

- Always choose EXACTLY ONE intent, ONE issue, and ONE product from the Stage-1 921

taxonomy (no new labels). 922

- Keep reasoning concise and reference concrete spans from the dialogue 923

(short quotes). IMPORTANT: Do NOT repeat the Stage-1 JSON. 924

Produce the reviewer JSON ONLY using the final_* keys. 925
926

C.3 Dynamic conversation simulation prompt templates 927

The following prompt templates were used in this study to serve a two-stage approach for generating realistic customer service conversations. The system_prompt_dialogue_helper and user_intent_helper 928

work together to analyze existing conversation history and identify all possible customer intents (such as ChitChat, Frustration, IntentChange, or FollowUpQuestion) that could naturally occur next in the 929

dialogue flow. Once potential intents are identified, the system randomly selects one and employs the system_prompt_talk_to_bot and user_prompt_turn_helper templates to generate authentic customer 930

responses that align with the chosen intent. This dual-phase prompting strategy ensures that 931

simulated conversations maintain conversational coherence while introducing realistic variability in 932

customer behavior, enabling comprehensive testing of chatbot performance across diverse interaction 933

scenarios. 934

935

936

937

938

system_prompt_dialogue_helper: 939

ROLE: 940

You are the dialogue helper for a user simulator helping find the intent for the 941

user given a conversation history. 942

943

944

TASK: 945

Select out all POSSIBLE intent from CANDIDATE LIST to carry on the conversation 946

given the previous Conversation history. 947

948

GUIDELINES: 949

1. Read through the whole conversation and identify the subset of 950

POSSIBLE INTENTS 951

952

CANDIDATE LIST: 953

1.ChitChat: Small talk loosely related to previous chat history 954

2.Frustration: Expression of frustration in the middle of conversation 955

3.intentChage: The user changes their request midway through a conversation, the 956

request can be related to previous Conversation history 957

4.FollowUpQuestion: The user asks follow-up questions that are related to previous 958

Conversation history 959

5.Clarification: The user asks for clarification for certain points in previous 960

Conversation history 961

6.Rambling: Speaker(s) ramble and repeat themselves. They may paraphrase themselves 962

7.ContactRealAgent: Request to speak to a real agent 963

8.EndConversation: End Conversation naturally when the issue or problem is resolved 964

965

966 Here is the input format
967 <Conversation history>
968 [provide chat_history here]
969 </Conversation history>
970
971 Here is the output format
972 <Possible Intents>
973 [provide possible intents here]
974 </Possible Intents>
975

976 **user_intent_helper:**

977 Here is your input:

978
979 <Conversation history>
980 {chat_history}
981 </Conversation history>
982

983 Now respond with what the customer would say next:
984

985 **system_prompt_talk_to_bot:**

986 ROLE:

987 You are a user engaging in a natural conversation with a customer service bot or
988 agent. Your goal is to generate the next user turn based on the conversation
989 history and the intent provided below.
990

991 OBJECTIVES

992 Conversation-Level Goal:

993 Seek a resolution (e.g., HOW-TO answer) to the seed query provided below.

994 Current Turn Goal:

995 Generate a user response that aligns with the current intent described below.
996

997 INTENT DEFINITIONS

998 ChitChat - Casual or light-hearted comments loosely related to the conversation.

999 Frustration - Expressions of annoyance or dissatisfaction.

1000 IntentChange - The user changes their goal mid-conversation, potentially related
1001 to prior turns.

1002 FollowUpQuestion - User asks a question that builds directly on prior discussion.

1003 Clarification - User requests clarification about something mentioned previously.

1004 Rambling - User paraphrases, repeats, or meanders while staying within the topic.

1005 ContactRealAgent: Request to speak to a real agent

1006 EndConversation: End Conversation naturally when the issue or problem is resolved
1007

1008 GUIDELINES:

1009 -If no chat history exists, begin with the seed query.

1010 -Respond naturally: ask relevant questions, express preferences, or make decisions
1011 as needed.

1012 -If the bot successfully resolves the task and provides a reference number, reply
1013 only with: "I'm all set" (no additional text).

1014 -If the bot is repetitive or unhelpful across multiple turns, escalate by using
1015 Contact Real Agent intent and say "talk to a real agent".

1016 -Do not impersonate a bot or break character.

1017 -Be concise and speak like a real customer in real life. Each response should be
1018 less than 25 words.

1019 -If the intent is either ContactRealAgent or EndConversation. Be concise and the
1020 response should be straight forward.

1021 Here is the input format

1022 <seed_query>

1023 [provide seed_query here]

1024 </seed_query>

1025 <Provided Intent>

1026 [provide intent here]

1027 </Provided Intent>

1028 <Conversation history>

1029 [provide chat_history here]

1030 </Conversation history>

1031 Here is the output format

1032 <Current Turn>

```
[provide the current turn]
</Current Turn>
"""
```

1033
1034
1035
1036

user_prompt_turn_helper:

Here is your input:

```
<Conversation history>
{chat_history}
```

1037
1038
1039
1040
1041
1042
1043
1044

```
<Provided Intent>
{provided_intent}
</Provided Intent>
```

1045
1046
1047
1048
1049

Now respond with what the customer would say next:

C.4 RoBERTa model fine-tuning implementation details

1050

We use RoBERTa-base for intent detection model, optimizing with cross-entropy and the Adam optimizer. Models are trained for 10 epochs on 8 NVIDIA A10 GPUs, with learning rates of $1e-5$, batch sizes of 32 and early stopping(patience = 5) to prevent overfitting.

1051
1052
1053

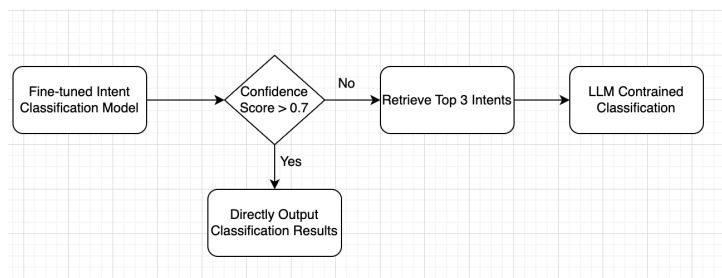


Figure 4: Hybrid deployment approach

C.5 Hybrid approach work flow implementation Details

1054

Figure 4 illustrates the workflow of our hybrid approach. In our implementation, we set the confidence threshold to $\tau = 0.7$ and use $K = 3$ candidate intents for LLM disambiguation (see Appendix C.6). We use a fine-tuned RoBERTa model as the lightweight classifier, and invoke Claude Sonnet 3.7 for intent classification when the confidence score falls below the threshold.

1055
1056
1057
1058

C.6 Analysis of τ and K in hybrid approach

1059

We select the confidence threshold by balancing accuracy and coverage. As shown in Figure 5, setting the threshold to 0.7 allows our fine-tuned model to cover roughly 80% of intent-detection requests while maintaining about 85% accuracy. We therefore choose 0.7 as the operating point because it offers a practical tradeoff: the low-latency model can handle the majority of traffic with sufficiently high accuracy for direct deployment. For the remaining low-confidence cases (confidence ≤ 0.7), we defer intent detection to an LLM. Importantly, for these instances, the ground-truth intent appears in our fine-tuned model’s top-3 predictions nearly 90% of the time (Figure 6). This suggests that low confidence typically reflects ambiguity among a small set of plausible intents rather than a complete failure. Accordingly, we ask the LLM to select among the top-3 candidate intents produced by the fine-tuned model.

1060
1061
1062
1063
1064
1065
1066
1067
1068

This yields an effective hybrid intent detection strategy: high-confidence requests (confidence ≥ 0.7) are handled directly by the lightweight model, while only a small fraction of ambiguous cases trigger a secondary LLM call over a constrained top-3 label set. This design improves accuracy on difficult utterances while keeping overall latency under control.

1069
1070
1071
1072

This is a section in the appendix.

1073

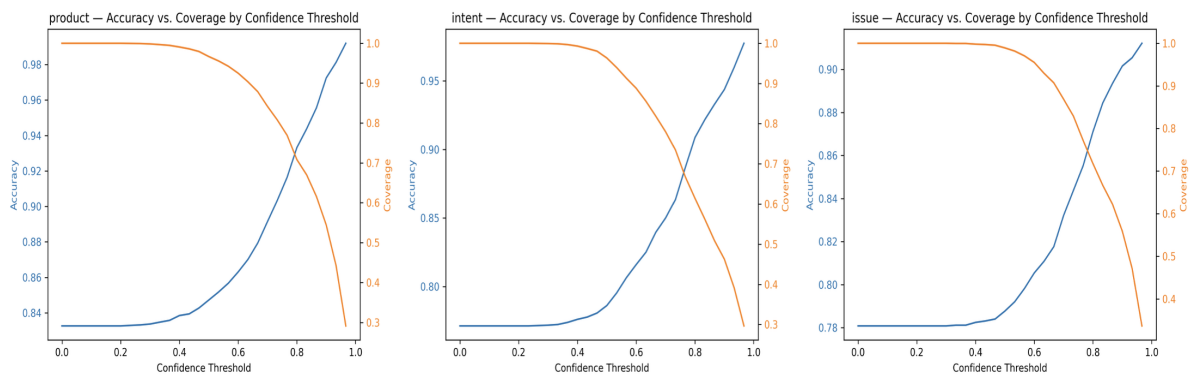


Figure 5: Accuracy–coverage tradeoff under different confidence thresholds

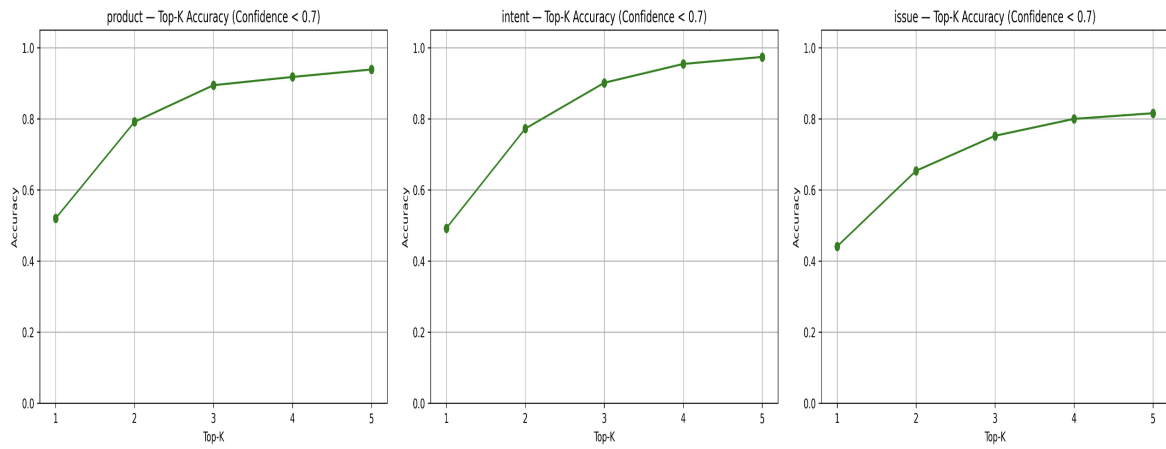


Figure 6: Top- K accuracy for low-confidence predictions