

A Joint Optimization Method for Multi-Turn Dialogue Intent Prediction and Adaptive Human-Machine Transfer in Intelligent Customer Service Scenarios

Yanyan Zhang^{1,*}

Carnegie Mellon University, Pittsburgh, PA

* Corresponding author: Yanyan Zhang
zhangyanyanusa@gmail.com

Abstract: In multi-turn customer service dialogue, intent prediction and human handoff are often treated as separate tasks, leading to inconsistent decisions when model confidence is low. In this paper, we propose an uncertainty-aware joint modeling framework that connects these two processes. A shared dialogue representation is learned, followed by dual modules for intent prediction and transfer decision. We introduce an entropy-based uncertainty measure to capture prediction confidence and formulate handoff as a conditional decision dependent on both intent and uncertainty. A lightweight consistency constraint is further applied to align confidence with system behavior. The method we put forward offers a straightforward and flexible alternative to conventional rule-based schemes.

Keywords: Multi-turn dialogue; Intent prediction; Human-machine transfer; Uncertainty estimation; Intelligent customer service; Joint modeling;

1. Introduction

With the rapid development of intelligent customer service systems, multi-turn dialogue has become a mainstream interaction method in real-world scenarios such as e-commerce, financial services, and technical support. Unlike single-round dialogue systems, multi-turn conversation models can fully utilize contextual information to generate more flexible and coherent replies. Nevertheless, this flexibility also introduces new difficulties. As the dialogue proceeds, user intention may become ambiguous, shift between rounds, or be closely tied to previous content.

This often leads to uncertainty in intention recognition during multi-turn conversations, particularly in complex or vaguely expressed situations^[1].

At the same time, deciding the proper timing for transferring a conversation to human agents remains a critical issue in practical deployment. Too early a transfer raises operational costs, while an overly delayed handover can damage user experience and satisfaction. Most current approaches adopt rule-based logic or static confidence thresholds, which fail to effectively reflect the dynamic characteristics of real-world dialogues.

Despite their close relationship, intent prediction and human handoff are usually treated as independent tasks. But in real-world use, these two processes are closely linked. When the system isn't sure about a user's intent, it usually makes more sense to hand the conversation over to a human agent. Conversely, when the model is confident, automated responses are usually sufficient^[2].

Ignoring this interaction may lead to inconsistent or suboptimal decisions. In this paper, we propose an uncertainty-aware joint modeling framework that explicitly connects intent prediction with human handoff decisions. Instead of relying on predefined rules, the model learns to use uncertainty as a signal to guide decision-making^[3]. By introducing an entropy-based uncertainty measure and a conditional decision mechanism, we provide a unified and interpretable approach to coordinating automation and human intervention.

2.Related Work

2.1 Intent Prediction in Multi-Turn Dialogue

Intent prediction serves as a core task in dialogue systems, which is usually modeled as a classification problem. Early approaches usually analyze each user message on its own, which makes it hard to capture contextual relationships. To fix this problem, newer models use recurrent neural networks or Transformer structures to incorporate full dialogue history, allowing them to better understand multi-turn conversations^[4]. In spite of such advances, most existing methods concentrate mainly on prediction accuracy and take the reliability of model outputs for granted. In practical scenarios, user intent may be ambiguous or change progressively, resulting in predictions with inconsistent confidence levels. The absence of explicit uncertainty modeling weakens the robustness of such systems in real-world applications.

2.2 Human-Machine Transfer Strategies

Deciding the appropriate timing for transferring a conversation to a human agent represents a core module in contemporary intelligent customer service systems^[5]. Conventional strategies typically depend on rule-based mechanisms—for instance, predefined keywords or explicit user demands—or fixed confidence thresholds, under which handover is initiated once the model’s predictive confidence drops under a set value.

Although such methods are straightforward to implement and launch, they are insufficient in terms of adaptive capacity. Fixed confidence thresholds often fail to generalize across different domains or dialogue contexts, and rule-based systems are prone to breakdown when faced with diverse user expressions that fall outside pre-set rules. More significantly, such approaches usually make escalation decisions separately from intent prediction, without accounting for how predictive uncertainty should shape the system’s overall behavior^[6].

2.3 Joint Modeling

Multi-task^[7] learning is widely used in dialogue systems to boost performance by sharing features across related tasks such as intent detection and dialogue state tracking^[8]. This allows models to learn general patterns and generalize better to unfamiliar situations.

However, many existing joint-learning frameworks do not explicitly model how these tasks interact with one another^[9]. In particular, the relationship between intent prediction and human handoff decisions has not been fully explored. Simply sharing latent representations does not guarantee consistent or reasonable system behavior, especially when predictive uncertainty is a key factor^[10].

In contrast, this work explicitly models the interaction between the two tasks by introducing an uncertainty-aware mechanism^[11]. By linking prediction confidence with transfer decisions, the proposed approach provides a more coherent and adaptive framework for real-world dialogue systems.

3. Methodology

3.1 Problem Formulation

We represent a multi-turn dialogue as a sequence of utterances:

$$D = \{u_1, s_1, u_2, s_2, \dots, u_t\}$$

where u_i and s_i denote the user and system utterances at turn i , respectively.

Given a dialogue context, the system is expected to handle two closely related tasks. The first is intent prediction, which aims to identify the underlying user intent:

$$P(y_{\text{intent}} \mid D)$$

The second is the human handoff decision, which determines whether the dialogue should be transferred to a human agent:

$$P(y_{\text{transfer}} \mid D)$$

In many current approaches, these two tasks are constructed and optimized independently. Although this design simplifies system development, it fails to fully capture the real-world decision-making logic. Intuitively, when the system has low confidence in identifying user intent, transferring the conversation to a human agent tends to be a more dependable solution. In contrast, high-confidence predictions generally allow the system to handle the dialogue autonomously. To capture this interaction, we consider a joint formulation:

$$P(y_{\text{intent}}, y_{\text{transfer}} \mid D)$$

This perspective makes it possible to model how intent prediction influences transfer decisions. In addition, we introduce an uncertainty variable u , and formulate the transfer decision as a conditional process: $P(y_{\text{transfer}} \mid h, \hat{y}_{\text{intent}}, u)$ where \hat{y}_{intent} is the predicted intent distribution and u reflects the model's confidence.

3.2 Model Overview

The overall model follows a unified architecture with a shared encoder and task-specific components.

Given the dialogue D , we first compute a contextual representation:

$$h = \text{Encoder}(D)$$

where h encodes the semantic information of the entire dialogue context. In this work, the encoder can be implemented using a Transformer-based architecture^[12].

On top of this shared representation, two prediction modules are constructed. One is responsible for intent prediction, and the other for transfer decision. The key difference from standard multi-task setups is that the transfer decision is not made independently, but instead conditioned on both the predicted intent and its associated uncertainty.

3.3 Intent Prediction

The intent prediction module is implemented as a standard classification layer:

$$P(y_{\text{intent}} \mid h) = \text{softmax}(W_i h + b_i)$$

where W_i and b_i are learnable parameters.

This module produces a probability distribution over all possible intent categories, which is later used not only for prediction but also for estimating uncertainty.

3.4 Uncertainty Estimation

To measure how confident the model is in its predictions, we use entropy^[13] as a simple and direct way to represent uncertainty:

$$u = - \sum_{k=1}^K P(y_k) \log P(y_k)$$



When u is higher, the model's predicted probabilities are more evenly distributed, meaning it is less certain about the user's intent. A lower u , on the other hand, means the model is more confident in its prediction.

3.5 Uncertainty-Aware Transfer Decision

Instead of making transfer decisions solely based on the dialogue representation h , we incorporate both the predicted intent and

its uncertainty:
$$P(y_{\text{transfer}} \mid h, \hat{y}_{\text{intent}}, u) = \sigma(W_t[h; \hat{y}_{\text{intent}}; u] + b_t)$$

where $[\cdot][\cdot]$ denotes vector concatenation and σ is the sigmoid function.

This formulation allows the model to adjust its behavior dynamically. When uncertainty is high, the model can assign a higher probability to transferring the dialogue. When confidence is strong, it is more likely to continue automated interaction.

3.6 Training Objective

The overall training objective combines both tasks:

$$L = L_{\text{intent}} + \lambda L_{\text{transfer}} + \gamma L_{\text{consistency}}$$

where:

L_{intent} is the cross-entropy loss for intent prediction

L_{transfer} is the binary classification loss for transfer decision

$L_{\text{consistency}}$ is a regularization term that links uncertainty with transfer behavior

3.7 Consistency Constraint

To explicitly align uncertainty with decision behavior, we introduce a simple consistency

constraint:
$$L_{\text{consistency}} = (u - P(y_{\text{transfer}}))^2$$

This term guides the model to produce a higher transfer probability when uncertainty is elevated, and a lower probability when its predictions are confident. In this manner, the entire decision-making process aligns more naturally with intuitive practical expectations.

3.8 Discussion

Compared with common practice, this framework changes how the transfer decision is made. Rather than depending on fixed rules or manually set thresholds, the decision is learned directly from data and driven by the model's internal confidence signal. This makes the system more adaptable across diverse dialogue scenarios. At the same time, the use of uncertainty provides a degree of interpretability. When a handoff is triggered, it can usually be explained by low confidence in intent prediction, which matches how human agents typically make similar judgments in practice^{[20][21]}.

Overall, the proposed method provides a simple way to link intent prediction with transfer decisions, while maintaining a relatively simple and easy-to-implement model structure.

4. Discussion

4.1 Comparison with Threshold-Based Methods

Many real-world systems still rely on simple rules to trigger human handoff, like fixed confidence thresholds. While these are easy to put into practice, they usually need careful manual adjustment and do not work well across different dialogue situations. A threshold that works fine in one setting may not work at all in another^{[17][19]}.

By contrast, the method proposed in this paper incorporates the transfer decision into the learning framework. Instead of relying on manually set rules, the model adjusts its decisions based on both the predicted intent and its own confidence. This approach makes the whole decision-making process more adaptive and better suited to the unique context of each individual conversation.

4.2 Uncertainty in Decision-Making

A central idea in this work is to use uncertainty as an explicit signal when making decisions. In real-world interactions, uncertainty often serves as an indicator of prediction reliability. When a system lacks confidence in its predictions, deferring to a human agent is generally the safer course of action.

By integrating uncertainty into the model, the system can naturally learn this intuitive behavior—eliminating the need for manually designed heuristics. This also makes the decision-making process easier to understand. For instance, a higher chance of transferring to a human agent usually comes from greater uncertainty in predictions, which matches how human support staff would judge similar situations^[18].

4.3 Practical Implications

From a practical standpoint, this framework provides a simple way to strike a balance between automated processing and human intervention. In large-scale customer service systems, excessive unnecessary transfers raise operational costs, while inadequate human support can damage the user experience. A mechanism that adjusts based on varying confidence levels helps manage this trade-off more efficiently^{[15][16]}.

Another benefit lies in the model's relative simplicity. Since the model only uses standard parts like shared encoders and classification layers, it can be added directly into existing systems without major changes to the overall structure.

4.4 Limitations and Future Directions

There are several limitations to this work. First, the uncertainty here is based only on the predicted intent, without accounting for other factors that might affect handoff decisions. In practice, user sentiment, dialogue complexity, and other details also matter for determining whether to transfer to a human agent.

Second, this framework works best for shorter dialogues. For longer and more complex dialogues, the current model may not fully capture the subtle details needed to make reliable transfer decisions.

In future work, we can include more useful signals such as user sentiment and conversation complexity, and also explore how reinforcement learning can help the system adapt better across different scenarios over time^[14].

5. Conclusion

In this study, we examined how intent prediction and human handoff decisions work together in multi-turn dialogue systems. Rather than treating the two tasks as separate independent modules, this paper proposes a simple method to connect them via shared feature representations and uncertainty indicators.

The core idea is intuitive: when the model lacks confidence in its predictive results, it should adopt a more cautious strategy and show a stronger tendency to transfer the dialogue to human agents. By employing entropy to measure uncertainty and integrating this indicator into the decision-making process, the model can naturally exhibit such prudent decision-making behavior. Adding an extra consistency constraint also helps reinforce this logic, making the overall decision process much easier to understand.

Although the framework is relatively simple, it shows that even lightweight adjustments to the way these tasks are connected can result in more consistent and reasonable system performance. In practical scenarios, this method helps achieve a more reasonable

balance between automated responses and human intervention. Future research can be extended in several meaningful directions. For instance, integrating additional indicators like user sentiment and dialogue complexity can provide more comprehensive contextual clues, supporting more robust decision-making. Besides, exploring how reinforcement learning can be applied to iteratively optimize the handover mechanism during continuous conversational interactions is also a worthwhile direction.

Data Availability Statement

Data will be made available on request.

Funding

This work was supported without any funding.

Conflicts of Interest

The author(s) declare no conflicts of interest.

Ethical Approval and Consent to Participate

Not applicable.

References

- [1] Devi, M. (2022). *Context-Aware Dialogue Systems for Human-Computer Interaction Applications*. *International Journal of Science, Research and Technology*, 5(3), 7769-7776.
- [2] Yi, Z., Ouyang, J., Xu, Z., Liu, Y., Liao, T., Luo, H., & Shen, Y. (2025). *A survey on recent advances in llm-based multi-turn dialogue systems*. *ACM Computing Surveys*, 58(6), 1-38.
- [3] Li, K., Chen, X., Song, T., Zhou, C., Liu, Z., Zhang, Z., ... & Shan, Q. (2025). *Solving situation puzzles with large language model and external reformulation*. *arXiv preprint arXiv:2503.18394*.
- [4] Luo, M., Du, B., Zhang, W., Song, T., Li, K., Zhu, H., ... & Wen, H. (2023). *Fleet rebalancing for expanding shared e-Mobility systems: A multi-agent deep reinforcement learning approach*. *IEEE Transactions on Intelligent Transportation Systems*, 24(4), 3868-3881.
- [5] Zhu, H., Luo, Y., Liu, Q., Fan, H., Song, T., Yu, C. W., & Du, B. (2019). *Multistep flow prediction on car-sharing systems: A multi-graph convolutional neural network with attention mechanism*. *International Journal of Software Engineering and Knowledge Engineering*, 29(11n12), 1727 - 1740.
- [6] Frayne, E., & Elowen, C. (2025). *Cognitive Interaction Models: Deep Learning Approaches for Intelligent and Context-Aware Human-Machine Systems*. *Artificial Intelligence and Computing Innovations*, 5(11).
- [7] Hao, Z. (2026). *Low-Overhead Scheduling for Real-Time AI Workloads on Multi-Core Edge Chips*. *International Journal of Advance in Applied Science Research*, 5(3), 15-25.
- [8] Wang, C. (2026). *A Study on Data-Driven Budget Optimization for US Enterprises' Cross-Border Marketing*. *Frontiers in Management Science*, 5(1), 41-46.
- [9] Soudani, H., Petcu, R., Kanoulas, E., & Hasibi, F. (2024). *A survey on recent advances in conversational data generation*. *ACM Computing Surveys*.
- [10] Li, K., Chen, X., Song, T., Zhang, H., Zhang, W., & Shan, Q. (2024). *GPTDrawer: Enhancing Visual Synthesis through ChatGPT*. *arXiv preprint arXiv:2412.10429*.
- [11] Hao, Z. (2025). *Fault-Tolerant Real-Time Scheduling for Edge AI in US Critical Infrastructure*. *Engineering Frontiers*, 1(4).
- [12] Zhang, Z., Li, S., Zhang, Z., Liu, X., Jiang, H., Tang, X., ... & Jiang, M. (2025). *IHEval: Evaluating language models on following the instruction hierarchy*. *arXiv preprint arXiv:2502.08745*.
- [13] Wu, Y. (2026). *Research on the Impact of LinkedIn Business Account Data-Driven Operations on Brand Exposure of AI Startups—A Case Study of AristAI*. *International Academic Journal of Social Science*, 2, 27-37.

- [14] Wang, H., Li, Q., & Liu, Y. (2023). Adaptive supervised learning on data streams in reproducing kernel Hilbert spaces with data sparsity constraint. *Stat*, 12(1), e514.
- [15] Natarajan, G. (2025). Dialogue Management Systems: How Conversational Agents Decide What to Say. *Journal Of Engineering And Computer Sciences*, 4(7), 230-239.
- [16] Luo, M., Zhang, W., Song, T., Li, K., Zhu, H., Du, B., & Wen, H. (2021, January). Rebalancing expanding EV sharing systems with deep reinforcement learning. In *Proceedings of the Twenty-Ninth International Conference on International Joint Conferences on Artificial Intelligence* (pp. 1338-1344).
- [17] Hao, Z. (2026). Energy Efficient Multi Core Task Scheduling for Real Time Edge AI Systems: A Latency Aware Approach. *International Journal of Advance in Applied Science Research*, 5(3), 1-14.
- [18] Chaudhuri, D. (2022). *Enriching Text-Based Human-Machine Interactions with Additional World Knowledge* (Doctoral dissertation, Universitäts-und Landesbibliothek Bonn).
- [19] Hao, Z. (2025). Task Affinity-Aware Scheduling for Multi-Core Edge Devices in Autonomous Vehicles. *Engineering Frontiers*, 1(2).
- [20] Zhang, Z., Li, S., Zhang, Z., Liu, X., Jiang, H., Tang, X., ... & Jiang, M. (2025). IHEval: Evaluating language models on following the instruction hierarchy. *arXiv preprint arXiv:2502.08745*.
- [21] Mo, F., Mao, K., Zhao, Z., Qian, H., Chen, H., Cheng, Y., ... & Nie, J. Y. (2025). A survey of conversational search. *ACM Transactions on Information Systems*, 43(6), 1-50.