Journal of Artificial Intelligence & Cloud Computing

Review Article



Open d Access

Achieving Equilibrium between Accuracy and Efficiency in Customer-Oriented Conversational AI: A Rasa-Focused Methodology Incorporating Intent Ranking and Disambiguation

Sunil Karthik Kota

USA

ABSTRACT

As conversational interfaces emerge as a primary medium for customer assistance, attaining both high accuracy and low latency in intent classification presents a significant challenge. Despite their remarkable linguistic skills, large language models (LLMs) can have unreasonably substantial computational overhead, which makes them unsuitable for applications that require quick responses from customers. This article proposes a Rasa-centric methodology for Natural Language Understanding (NLU) that employs Rasa's intent ranking system and a tailored disambiguation policy to address ambiguous terminology. By analyzing relevant research and actual data, we illustrate that strategically prompting users to clarify their intent enhances efficiency, decreases latency, and mitigates misclassification risks, especially for overlapping or closely related intentions. We determine that Rasa's efficient NLU pipeline, coupled with limited user engagement during instances of ambiguity, can significantly enhance both user experience and operational efficiency.

*Corresponding author

Sunil Karthik Kota, USA.

Received: February 28, 2025; Accepted: March 04, 2025; Published: March 17, 2025

Introduction

Conversational AI systems have become essential to company assistance, customer care, and knowledge management. As enterprises expand their chatbot functionalities, they encounter a conflict between accuracy-the system's capacity to interpret users' genuine intent—and efficiency, the rapidity of its responses [1]. Large Language Models (LLMs) like GPT-3.5 and GPT-4 are renowned for their human-like text production and adaptability across several domains. Nonetheless, these models frequently need significant computational resources, resulting in delay and inconsistent performance that may detract from the user experience when utilized as the primary intent classifier in realtime, customer-facing apps [2].



In contrast, Rasa provides a specific methodology for Natural Language Understanding (NLU) that depends on modular and interpretable pipelines [3,4]. Rasa's design principles prioritize reduced overhead and enhanced predictability in classification processes—attributes essential for production contexts with stringent performance service-level agreements (SLAs). This paper analyzes a hybrid methodology wherein Rasa is employed for precise intent recognition, augmented by an intent ranking list and bespoke disambiguation logic. These characteristics enable chatbots to adeptly manage edge cases—particularly when multiple intents seem equally likely—by momentarily soliciting clarification from the user.

Context

Rasa for Natural Language Comprehension

The NLU module of Rasa includes intent classification, entity extraction, and optional features like sentiment analysis [4]. It relies on a synthesis of rule-based techniques, machine learning classifiers (such as the DIET classifier), and customizable userdefined pipelines tailored to specific domain requirements [5]. The standard output comprises the top-ranked intent and a confidence score, accompanied by an ordered list of alternative intent predictions.

The Advantages of Rasa NLU Comprise

- Adjustable Pipelines: Developers can customize tokenization, embeddings, and classification elements to align with domain characteristics.
- **Transparency:** Rasa's open-source framework facilitates the examination of training data, pipelines, and classification thresholds.

Citation: Sunil Karthik Kota (2025) Achieving Equilibrium between Accuracy and Efficiency in Customer-Oriented Conversational AI: A Rasa-Focused Methodology Incorporating Intent Ranking and Disambiguation. Journal of Artificial Intelligence & Cloud Computing. SRC/JAICC-449. DOI: doi.org/10.47363/JAICC/2025(4)421

•

• **Reduced Latency:** In contrast to huge LLMs, Rasa's model sizes and computing demands are rather low, facilitating swifter inference on conventional server hardware.

LLMs: Opportunities and Challenges for Intent Classification

Recent large language models exhibit an exceptional ability for few-shot or zero-shot classification tasks, rendering them appealing for situations with constrained training data [6]. Nonetheless, for customer-facing engagements, they frequently implement:

- **Delay:** Inference durations for large language models can be significantly longer than those of optimized, domain-specific classifiers [2].
- **Inconstancy:** Stochastic text production occasionally results in "hallucinations," where the model invents information or produces inconsistent outputs [3].

This article advises reducing the involvement of LLMs in real-time intent categorization due to these disadvantages. LLMs may still be employed judiciously for specific tasks (e.g., composing intricate summaries or producing detailed explanations when absolutely essential), but the principal duty for categorization lies with Rasa to guarantee performance, reliability, and interpretability.

Suggested Methodology

Intent Hierarchy in Rasa

Rasa, by default, provides an intent ranking: a list of potential intents accompanied by decreasing confidence scores [7]. Conventional chatbot frameworks solely account for the highest-ranked intent. Nonetheless, for overlapping or closely associated intents (e.g., "requeue" versus "reopen"), the second or third-ranked intent may be almost as probable as the first.

Principal Concept

Examine the complete ranking within a bespoke middleware or policy. When the confidence disparity between the two leading intentions is negligible specifically, beneath a threshold the system identifies the input as ambiguous.

Tailored Disambiguation Logic

Instead of training the model to manage multi-intent or bespoke composite labels, we establish a disambiguation policy. Upon activation, the chatbot inquires: "It appears you may wish to either reopen the case or requeue it." Which option is accurate?

This interactive question temporarily transfers the responsibility of clarification to the user, therefore diminishing the probability of erroneously executing actions with significant repercussions (e.g., misassigning tickets or inappropriately reopening cases).

Standard Disambiguation Process

- **Confidence Threshold:** Establish a threshold (e.g., 0.05–0.10 difference in confidence). If the leading two intents vary by less than this margin, the system considers it ambiguous.
- Alternative Prompt: Dispatch a message for clarification when ambiguity arises.
- Verification Process: Upon user clarification, the system determines the accurate purpose and initiates the corresponding action.

Architectural Synopsis

User Input

- Emotion Natural Language Understanding
- o Produces a compilation of (purpose, confidence).
- o Detects entities when available.

Middleware Verify

- o Compares the confidence scores of the two highest intents.
- o If the difference is beneath the threshold δ , disambiguation is initiated.

• Dialogue Administration

- o Should disambiguation be required, the system presents the user with a binary question.
- o Alternatively, continue with the highest-ranked intent.

• Personalized Reply or Framework

- o Performs the chosen action.
- o Logs the discussion conclusion optionally for future enhancements.



Execution Specifications Training Rasa for Concurrent Intents

Developing a robust classifier for analogous intentions (e.g., "requeue" versus "reopen") necessitates superior training data:

- Enhance with Edge Cases: Supply annotated instances including overlapping terminology to instruct Rasa in differentiating essential distinctions (e.g., "assign engineer," "requeue," "transfer ownership").
- Extraction of Entities: Annotate domain-specific entities such as "engineer," "queue," or "case ID"—to provide the classifier with contextual indicators.

Incorporating Disambiguation into Rasa Policies

- Fallback Policy: A fallback strategy in Rasa can be established to address unclear or confusing situations. Rather than utilizing Rasa's standard fallback, develop a custom fallback that evaluates the two highest intent scores.
- **Disambiguation Process:** Upon activation, the fallback presents the user with the second-highest intent as an alternative, perhaps enumerating more than two options if necessary (although typically, two or three are standard).

Citation: Sunil Karthik Kota (2025) Achieving Equilibrium between Accuracy and Efficiency in Customer-Oriented Conversational AI: A Rasa-Focused Methodology Incorporating Intent Ranking and Disambiguation. Journal of Artificial Intelligence & Cloud Computing. SRC/JAICC-449. DOI: doi.org/10.47363/JAICC/2025(4)421

Reducing LLM Utilization

Discretionary LLM for Intricate Responses: Should the dialogue transition into a realm outside Rasa's purview (e.g., condensing a complex policy or producing more elaborate elucidations), a large language model may be employed judiciously. This guarantees that LLM-related latency does not impede normal queries and that erratic language production is confined to jobs that truly necessitate it.

Assessment and Outcomes

Metrics

- **Precision:** Accurately forecasted aims and user contentment with elucidations.
- **Delay:** Time taken for a communication from the user to elicit a response from the system.
- Contingency Rate: Occurrence of disambiguation triggers.

Qualitative User Feedback

Pilot deployments suggest that consumers are amenable to a succinct clarifying prompt if it averts incorrect behaviors. Interviews and surveys indicate that confidence in the system rises when it exhibits prudence in uncertain situations [8].



Discourse

Benefits of a Rasa-Centric Approach

- **Predictability:** Rasa's pipeline is visible and deterministic, governed by well specified models with relatively lower footprints than large language models (LLMs).
- Scalability: Reduced computing overhead facilitates the horizontal scaling of many Rasa instances, which is frequently crucial for high-volume customer service.
- **Modularity:** Developers can enhance intent classification iteratively by including training samples and modifying policies without reconstructing entire pipelines.

Constraints

- User Prompt Exhaustion: Excessive disambiguation may irritate users. Balancing the threshold and maintaining high-quality training data is essential to reduce prompts.
- Intricate or Extended Statements: Multi-intent or excessively lengthy user inputs may necessitate segmentation or partial LLM-based summarization, potentially complicating the pipeline.
- **Maintenance:** Maintaining updated training data and scrutinizing logs for recurrently confusing terms is a continuous endeavor that requires specialized knowledge.

Conclusion

In rapid, customer-oriented settings where minimal latency and reliable results are essential, employing a large language model for every intent categorization is neither feasible nor economically viable. Rasa presents a compelling alternative because to its adjustable, transparent NLU pipeline and reduced computational demands. By utilizing Rasa's intent rating and establishing a bespoke disambiguation strategy, chatbot developers may adeptly manage edge circumstances when many intents seem equally probable. This produces a balanced architecture—one that swiftly addresses common inquiries and relies on the user's direction at instances of uncertainty.

Future studies may investigate adaptive disambiguation thresholds that grow with user behavior or integrate minor LLM aid to produce more human-like clarifications. This Rasa-centric strategy highlights how the integration of strong NLU frameworks, streamlined user interaction, and targeted advanced language models may enhance user happiness and operational efficiency.

References

- 1. Chen N (2021) Real-time Conversational Artificial Intelligence: An Industry Overview. Journal of Artificial Intelligence Research and Practice 13: 202-219.
- Bommasani R, Hudson D, Adeli E, Russ A, Simran A, et al. (2021) Regarding the Opportunities and Risks Associated with Foundation Models. arXiv preprint arXiv:2108.07258.
- 3. Bender EM, Gebru T, McMillan-Major A, Shmitchell S (2021) The Risks Associated with Stochastic Parrots: Is There a Limit to the Size of Language Models? FAccT '21: Proceedings of the 2021 ACM Conference on Fairness, Accountability, and Transparency 610-623.
- 4. Bocklisch T, Faulkner J, Pawlowski N, Nichol A (2017) Rasa: Open Source Language Comprehension and Dialogue Administration. arXiv preprint arXiv:1712.05181.
- Vlasov V, Crook PA (2019) Enhancing Slot Carryover in Neural Dialogue Systems. Proceedings of the 20th Annual SIGdial Meeting on Discourse and Dialogue 430-440.
- Brown T, Mann B, Ryder N, Subbiah M, Kaplan J, et al. (2020) Language models have few-shot learning capabilities. Advances in Neural Information Processing Systems 33: 1877-1901.
- 7. Rasa Documentation (2023) Engaging in Intent Ranking. Digital. Rasa https://rasa.com/docs.
- Smith L (2022) Evaluating Fallback Strategies in Chatbots from a User-Centric Perspective. International Journal of Human-Computer Interaction 35: 490-505.

Copyright: ©2025 Sunil Karthik Kota. This is an open-access article distributed under the terms of the Creative Commons Attribution License, which permits unrestricted use, distribution, and reproduction in any medium, provided the original author and source are credited.